# LINEAR CONVERGENT DECENTRALIZED OPTIMIZATION WITH COMPRESSION

Xiaorui Liu<sup>1</sup>, Yao Li<sup>2,3</sup>, Rongrong Wang<sup>3,2</sup>, Jiliang Tang<sup>1</sup> & Ming Yan<sup>3,2</sup>

- <sup>1</sup> Department of Computer Science and Engineering
- <sup>2</sup> Department of Mathematics
- <sup>3</sup> Department of Computational Mathematics, Science and Engineering Michigan State University, East Lansing, MI 48823, USA {xiaorui,liyao6, wongron6, tangjili, myan}@msu.edu

#### **ABSTRACT**

Communication compression has become a key strategy to speed up distributed optimization. However, existing decentralized algorithms with compression mainly focus on compressing DGD-type algorithms. They are unsatisfactory in terms of convergence rate, stability, and the capability to handle heterogeneous data. Motivated by primal-dual algorithms, this paper proposes the first LinEAr convergent Decentralized algorithm with compression, LEAD. Our theory describes the coupled dynamics of the inexact primal and dual update as well as compression error, and we provide the first consensus error bound in such settings without assuming bounded gradients. Experiments on convex problems validate our theoretical analysis, and empirical study on deep neural nets shows that LEAD is applicable to non-convex problems.

## 1 Introduction

Distributed optimization solves the following optimization problem

$$\boldsymbol{x}^* := \underset{\boldsymbol{x} \in \mathbb{R}^d}{\min} \left[ f(\boldsymbol{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\boldsymbol{x}) \right]$$
 (1)

with n computing agents and a communication network. Each  $f_i(\boldsymbol{x}): \mathbb{R}^d \to \mathbb{R}$  is a local objective function of agent i and typically defined on the data  $\mathcal{D}_i$  settled at that agent. The data distributions  $\{\mathcal{D}_i\}$  can be heterogeneous depending on the applications such as in federated learning. The variable  $\boldsymbol{x} \in \mathbb{R}^d$  often represents model parameters in machine learning. A distributed optimization algorithm seeks an optimal solution that minimizes the overall objective function  $f(\boldsymbol{x})$  collectively. According to the communication topology, existing algorithms can be conceptually categorized into centralized and decentralized ones. Specifically, centralized algorithms require global communication between agents (through central agents or parameter servers). While decentralized algorithms only require local communication between connected agents and are more widely applicable than centralized ones. In both paradigms, the computation can be relatively fast with powerful computing devices; efficient communication is the key to improve algorithm efficiency and system scalability, especially when the network bandwidth is limited.

In recent years, various communication compression techniques, such as quantization and sparsification, have been developed to reduce communication costs. Notably, extensive studies (Seide et al., 2014; Alistarh et al., 2017; Bernstein et al., 2018; Stich et al., 2018; Karimireddy et al., 2019; Mishchenko et al., 2019; Tang et al., 2019b; Liu et al., 2020) have utilized gradient compression to significantly boost communication efficiency for centralized optimization. They enable efficient large-scale optimization while maintaining comparable convergence rates and practical performance with their non-compressed counterparts. This great success has suggested the potential and significance of communication compression in decentralized algorithms.

While extensive attention has been paid to centralized optimization, communication compression is relatively less studied in decentralized algorithms because the algorithm design and analysis are

more challenging in order to cover general communication topologies. There are recent efforts trying to push this research direction. For instance, DCD-SGD and ECD-SGD (Tang et al., 2018a) introduce difference compression and extrapolation compression to reduce model compression error. (Reisizadeh et al., 2019a;b) introduce QDGD and QuanTimed-DSGD to achieve exact convergence with small stepsize. DeepSqueeze (Tang et al., 2019a) directly compresses the local model and compensates the compression error in the next iteration. CHOCO-SGD (Koloskova et al., 2019; 2020) presents a novel quantized gossip algorithm that reduces compression error by difference compression and preserves the model average. Nevertheless, most existing works focus on the compression of primal-only algorithms, i.e., reduce to DGD (Nedic & Ozdaglar, 2009; Yuan et al., 2016) or P-DSGD (Lian et al., 2017). They are unsatisfying in terms of convergence rate, stability, and the capability to handle heterogeneous data. Part of the reason is that they inherit the drawback of DGD-type algorithms, whose convergence rate is slow in heterogeneous data scenarios where the data distributions are significantly different from agent to agent.

In the literature of decentralized optimization, it has been proved that primal-dual algorithms can achieve faster converge rates and better support heterogeneous data (Ling et al., 2015; Shi et al., 2015; Li et al., 2019; Yuan et al., 2020). However, it is unknown whether communication compression is feasible for primal-dual algorithms and how fast the convergence can be with compression. In this paper, we attempt to bridge this gap by investigating the communication compression for primal-dual decentralized algorithms. Our major contributions can be summarized as:

- We delineate two key challenges in the algorithm design for communication compression in decentralized optimization, i.e., data heterogeneity and compression error, and motivated by primal-dual algorithms, we propose a novel decentralized algorithm with compression, LEAD.
- We prove that for LEAD, a constant stepsize in the range  $(0, 2/(\mu + L)]$  is sufficient to ensure linear convergence for strongly convex and smooth objective functions. To the best of our knowledge, LEAD is the first linear convergent decentralized algorithm with compression. Moreover, LEAD provably works with unbiased compression of arbitrary precision.
- We further prove that if the stochastic gradient is used, LEAD converges linearly to the  $O(\sigma^2)$  neighborhood of the optimum with constant stepsize. LEAD is also able to achieve exact convergence to the optimum with diminishing stepsize.
- Extensive experiments on convex problems validate our theoretical analyses, and the empirical study on training deep neural nets shows that LEAD is applicable for nonconvex problems. LEAD achieves state-of-art computation and communication efficiency in all experiments and significantly outperforms the baselines on heterogeneous data. Moreover, LEAD is robust to parameter settings and needs minor effort for parameter tuning.

#### 2 RELATED WORKS

Decentralized optimization can be traced back to the work by Tsitsiklis et al. (1986). DGD (Nedic & Ozdaglar, 2009) is the most classical decentralized algorithm. It is intuitive and simple but converges slowly due to the diminishing stepsize that is needed to obtain the optimal solution (Yuan et al., 2016). Its stochastic version D-PSGD (Lian et al., 2017) has been shown effective for training nonconvex deep learning models. Algorithms based on primal-dual formulations or gradient tracking are proposed to eliminate the convergence bias in DGD-type algorithms and improve the convergence rate, such as D-ADMM (Mota et al., 2013), DLM (Ling et al., 2015), EXTRA (Shi et al., 2015), NIDS (Li et al., 2019),  $D^2$  (Tang et al., 2018b), Exact Diffusion (Yuan et al., 2018), OPTRA(Xu et al., 2020), DIGing (Nedic et al., 2017), GSGT (Pu & Nedić, 2020), etc.

Recently, communication compression is applied to decentralized settings by Tang et al. (2018a). It proposes two algorithms, i.e., DCD-SGD and ECD-SGD, which require compression of high accuracy and are not stable with aggressive compression. Reisizadeh et al. (2019a;b) introduce QDGD and QuanTimed-DSGD to achieve exact convergence with small stepsize and the convergence is slow. DeepSqueeze (Tang et al., 2019a) compensates the compression error to the compression in the next iteration. Motivated by the quantized average consensus algorithms, such as (Carli et al., 2010), the quantized gossip algorithm CHOCO-Gossip (Koloskova et al., 2019) converges linearly to the consensual solution. Combining CHOCO-Gossip and D-PSGD leads to a decentralized algorithm with compression, CHOCO-SGD, which converges sublinearly under the strong convexity and

gradient boundedness assumptions. Its nonconvex variant is further analyzed in (Koloskova et al., 2020). A new compression scheme using the modulo operation is introduced in (Lu & De Sa, 2020) for decentralized optimization. A general algorithmic framework aiming to maintain the linear convergence of distributed optimization under compressed communication is considered in (Magnússon et al., 2020). It requires a contractive property that is not satisfied by many decentralized algorithms including the algorithm in this paper.

#### 3 ALGORITHM

We first introduce notations and definitions used in this work. We use bold upper-case letters such as  $\mathbf{X}$  to define matrices and bold lower-case letters such as  $\mathbf{x}$  to define vectors. Let  $\mathbf{1}$  and  $\mathbf{0}$  be vectors with all ones and zeros, respectively. Their dimensions will be provided when necessary. Given two matrices  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times d}$ , we define their inner product as  $\langle \mathbf{X}, \mathbf{Y} \rangle = \operatorname{tr}(\mathbf{X}^{\top}\mathbf{Y})$  and the norm as  $\|\mathbf{X}\| = \sqrt{\langle \mathbf{X}, \mathbf{X} \rangle}$ . We further define  $\langle \mathbf{X}, \mathbf{Y} \rangle_{\mathbf{P}} = \operatorname{tr}(\mathbf{X}^{\top}\mathbf{P}\mathbf{Y})$  and  $\|\mathbf{X}\|_{\mathbf{P}} = \sqrt{\langle \mathbf{X}, \mathbf{X} \rangle_{\mathbf{P}}}$  for any given symmetric positive semidefinite matrix  $\mathbf{P} \in \mathbb{R}^{n \times n}$ . For simplicity, we will majorly use the matrix notation in this work. For instance, each agent i holds an individual estimate  $\mathbf{x}_i \in \mathbb{R}^d$  of the global variable  $\mathbf{x} \in \mathbb{R}^d$ . Let  $\mathbf{X}^k$  and  $\nabla \mathbf{F}(\mathbf{X}^k)$  be the collections of  $\{\mathbf{x}_i^k\}_{i=1}^n$  and  $\{\nabla f_i(\mathbf{x}_i^k)\}_{i=1}^n$  which are defined below:

$$\mathbf{X}^{k} = \begin{bmatrix} \boldsymbol{x}_{1}^{k}, \dots, \boldsymbol{x}_{n}^{k} \end{bmatrix}^{\top} \in \mathbb{R}^{n \times d}, \quad \nabla \mathbf{F}(\mathbf{X}^{k}) = \begin{bmatrix} \nabla f_{1}(\boldsymbol{x}_{1}^{k}), \dots, \nabla f_{n}(\boldsymbol{x}_{n}^{k}) \end{bmatrix}^{\top} \in \mathbb{R}^{n \times d}.$$
 (2)

We use  $\nabla \mathbf{F}(\mathbf{X}^k; \xi^k)$  to denote the stochastic approximation of  $\nabla \mathbf{F}(\mathbf{X}^k)$ . With these notations, the update  $\mathbf{X}^{k+1} = \mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^k; \xi^k)$  means that  $\mathbf{x}_i^{k+1} = \mathbf{x}_i^k - \eta \nabla f_i(\mathbf{x}_i^k; \xi_i^k)$  for all i. In this paper, we need the average of all rows in  $\mathbf{X}^k$  and  $\nabla \mathbf{F}(\mathbf{X}^k)$ , so we define  $\overline{\mathbf{X}}^k = (\mathbf{1}^{\top} \mathbf{X}^k)/n$  and  $\overline{\nabla} \mathbf{F}(\mathbf{X}^k) = (\mathbf{1}^{\top} \nabla \mathbf{F}(\mathbf{X}^k))/n$ . They are row vectors, and we will take a transpose if we need a column vector. The pseudoinverse of a matrix  $\mathbf{M}$  is denoted as  $\mathbf{M}^{\dagger}$ . The largest, ith-largest, and smallest nonzero eigenvalues of a symmetric matrix  $\mathbf{M}$  are  $\lambda_{\max}(\mathbf{M})$ ,  $\lambda_i(\mathbf{M})$ , and  $\lambda_{\min}(\mathbf{M})$ .

**Assumption 1** (Mixing matrix). The connected network  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  consists of a node set  $\mathcal{V} = \{1, 2, ..., n\}$  and an undirected edge set  $\mathcal{E}$ . The primitive symmetric doubly-stochastic matrix  $\mathbf{W} = [w_{ij}] \in \mathbb{R}^{n \times n}$  encodes the network structure such that  $w_{ij} = 0$  if nodes i and j are not connected and cannot exchange information.

Assumption 1 implies that  $-1 < \lambda_n(\mathbf{W}) \le \lambda_{n-1}(\mathbf{W}) \le \cdots \lambda_2(\mathbf{W}) < \lambda_1(\mathbf{W}) = 1$  and  $\mathbf{W}\mathbf{1} = \mathbf{1}$  (Xiao & Boyd, 2004; Shi et al., 2015). The matrix multiplication  $\mathbf{X}^{k+1} = \mathbf{W}\mathbf{X}^k$  describes that agent i takes a weighted sum from its neighbors and itself, i.e.,  $\mathbf{x}_i^{k+1} = \sum_{j \in \mathbb{N}_i \cup \{i\}} w_{ij} \mathbf{x}_j^k$ , where  $\mathbb{N}_i$  denotes the neighbors of agent i.

#### 3.1 THE PROPOSED ALGORITHM

The proposed algorithm LEAD to solve problem (1) is showed in Alg. 1 with matrix notations for conciseness. We will refer to the line number in the analysis. A complete algorithm description from the agent's perspective can be found in Appendix A. The motivation behind Alg. 1 is to achieve two goals: (a) consensus  $(\boldsymbol{x}_i^k - (\overline{\mathbf{X}}^k)^\top \to \mathbf{0})$  and (b) convergence  $((\overline{\mathbf{X}}^k)^\top \to \boldsymbol{x}^*)$ . We first discuss how goal (a) leads to goal (b) and then explain how LEAD fulfills goal (a).

In essence, LEAD runs the approximate SGD globally and reduces to the exact SGD under consensus. One key property for LEAD is  $\mathbf{1}_{n\times 1}^{\mathsf{T}}\mathbf{D}^k=\mathbf{0}$ , regardless of the compression error in  $\hat{\mathbf{Y}}^k$ . It holds because that for the initialization, we require  $\mathbf{D}^1=(\mathbf{I}-\mathbf{W})\mathbf{Z}$  for some  $\mathbf{Z}\in\mathbb{R}^{n\times d}$ , e.g.,  $\mathbf{D}^1=\mathbf{0}^{n\times d}$ , and that the update of  $\mathbf{D}^k$  ensures  $\mathbf{D}^k\in\mathbf{Range}(\mathbf{I}-\mathbf{W})$  for all k and  $\mathbf{1}_{n\times 1}^{\mathsf{T}}(\mathbf{I}-\mathbf{W})=\mathbf{0}$  as we will explain later. Therefore, multiplying  $(1/n)\mathbf{1}_{n\times 1}^{\mathsf{T}}$  on both sides of Line 7 leads to a global average view of Alg. 1:

$$\overline{\mathbf{X}}^{k+1} = \overline{\mathbf{X}}^k - \eta \overline{\nabla} \mathbf{F}(\mathbf{X}^k; \xi^k), \tag{3}$$

which doesn't contain the compression error. Note that this is an approximate SGD step because, as shown in (2), the gradient  $\nabla \mathbf{F}(\mathbf{X}^k; \xi^k)$  is not evaluated on a global synchronized model  $\overline{\mathbf{X}}^k$ . However, if the solution converges to the consensus solution, i.e.,  $\boldsymbol{x}_i^k - (\overline{\mathbf{X}}^k)^\top \to \mathbf{0}$ , then  $\mathbb{E}_{\boldsymbol{\xi}^k}[\overline{\nabla}\mathbf{F}(\mathbf{X}^k; \xi^k) - \nabla f(\overline{\mathbf{X}}^k; \xi^k)] \to \mathbf{0}$  and (3) gradually reduces to exact SGD.

#### Algorithm 1 LEAD

```
Input: Stepsize \eta, parameter (\alpha, \gamma), \mathbf{X}^0, \mathbf{H}^1, \mathbf{D}^1 = (\mathbf{I} - \mathbf{W})\mathbf{Z} for any \mathbf{Z}
Output: \mathbf{X}^K or 1/n \sum_{i=1}^n \mathbf{X}_i^K
1: \mathbf{H}_w^1 = \mathbf{W}\mathbf{H}^1
                                                                                                                                   9: procedure COMM((\mathbf{Y}, \mathbf{H}, \mathbf{H}_w))
2: \mathbf{X}^1 = \mathbf{X}^0 - n\nabla \mathbf{F}(\mathbf{X}^0; \boldsymbol{\varepsilon}^0)
                                                                                                                                 10:
                                                                                                                                                  \mathbf{Q} = \text{COMPRESS}(\mathbf{Y} - \mathbf{H})
3: for k = 1, 2, \dots, K - 1 do
                                                                                                                                                  \hat{\mathbf{Y}} = \mathbf{H} + \mathbf{Q}
                                                                                                                                 11:
               \mathbf{Y}^k = \mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^k; \xi^k) - \eta \mathbf{D}^k
                                                                                                                                                  \hat{\mathbf{Y}}_w = \mathbf{H}_w + \mathbf{W}\mathbf{Q}
                                                                                                                                 12:
              \hat{\mathbf{Y}}^k, \hat{\mathbf{Y}}_w^k, \mathbf{H}^{k+1}, \mathbf{H}_w^{k+1} = \text{Comm}(\mathbf{Y}^k, \mathbf{H}^k, \mathbf{H}_w^k)
                                                                                                                                                 \mathbf{H} = (1 - \alpha)\mathbf{H} + \alpha \hat{\mathbf{Y}}
                                                                                                                                13:
             \mathbf{D}^{k+1} = \mathbf{D}^k + \frac{\gamma}{2n} (\hat{\mathbf{Y}}^k - \hat{\mathbf{Y}}_w^k)
                                                                                                                                                  \mathbf{H}_w = (1 - \alpha)\mathbf{H}_w + \alpha \hat{\mathbf{Y}}_w
               \mathbf{X}^{k+1} = \mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \eta \mathbf{D}^{k+1}
                                                                                                                                                  Return: \hat{\mathbf{Y}}, \hat{\mathbf{Y}}_w, \mathbf{H}, \mathbf{H}_w
                                                                                                                                 15:
                                                                                                                                 16: end procedure
8: end for
```

With the establishment of how consensus leads to convergence, the obstacle becomes how to achieve consensus under local communication and compression challenges. It requires addressing two issues, i.e., data heterogeneity and compression error. To deal with these issues, existing algorithms, such as DCD-SGD, ECD-SGD, QDGD, DeepSqueeze, Moniqua, and CHOCO-SGD, need a diminishing or constant but small stepsize depending on the total number of iterations. However, these choices unavoidably cause slower convergence and bring in the difficulty of parameter tuning. In contrast, LEAD takes a different way to solve these issues, as explained below.

Data heterogeneity. It is common in distributed settings that there exists data heterogeneity among agents, especially in real-world applications where different agents collect data from different scenarios. In other words, we generally have  $f_i(\boldsymbol{x}) \neq f_j(\boldsymbol{x})$  for  $i \neq j$ . The optimality condition of problem (1) gives  $\mathbf{1}_{n\times 1}^{\top}\nabla\mathbf{F}(\mathbf{X}^*)=\mathbf{0}$ , where  $\mathbf{X}^*=[\boldsymbol{x}^*,\cdots,\boldsymbol{x}^*]$  is a consensual and optimal solution. The data heterogeneity and optimality condition imply that there exist at least two agents i and j such that  $\nabla f_i(\boldsymbol{x}^*) \neq \mathbf{0}$  and  $\nabla f_j(\boldsymbol{x}^*) \neq \mathbf{0}$ . As a result, a simple D-PSGD algorithm cannot converge to the consensual and optimal solution as  $\mathbf{X}^* \neq \mathbf{W}\mathbf{X}^* - \eta \mathbb{E}_{\xi} \nabla \mathbf{F}(\mathbf{X}^*; \xi)$  even when the stochastic gradient variance is zero.

Gradient correction. Primal-dual algorithms or gradient tracking algorithms are able to convergence much faster than DGD-type algorithms by handling the data heterogeneity issue, as introduced in Section 2. Specifically, LEAD is motivated by the design of primal-dual algorithm NIDS (Li et al., 2019) and the relation becomes clear if we consider the two-step reformulation of NIDS adopted in (Li & Yan, 2019):

$$\mathbf{D}^{k+1} = \mathbf{D}^k + \frac{\mathbf{I} - \mathbf{W}}{2\eta} (\mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^k) - \eta \mathbf{D}^k), \tag{4}$$

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^k) - \eta \mathbf{D}^{k+1}, \tag{5}$$

where  $\mathbf{X}^k$  and  $\mathbf{D}^k$  represent the primal and dual variables respectively. The dual variable  $\mathbf{D}^k$  plays the role of gradient correction. As  $k \to \infty$ , we expect  $\mathbf{D}^k \to -\nabla \mathbf{F}(\mathbf{X}^*)$  and  $\mathbf{X}^k$  will converge to  $\mathbf{X}^*$  via the update in (5) since  $\mathbf{D}^{k+1}$  corrects the nonzero gradient  $\nabla \mathbf{F}(\mathbf{X}^k)$  asymptotically. The key design of Alg. 1 is to provide compression for the auxiliary variable defined as  $\mathbf{Y}^k = \mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^k) - \eta \mathbf{D}^k$ . Such design ensures that the dual variable  $\mathbf{D}^k$  lies in  $\mathbf{Range}(\mathbf{I} - \mathbf{W})$ , which is essential for convergence. Moreover, it achieves the implicit error compression as we will explain later. To stabilize the algorithm with inexact dual update, we introduce a parameter  $\gamma$  to control the stepsize in the dual update. Therefore, if we ignore the details of the compression, Alg. 1 can be concisely written as

$$\mathbf{Y}^k = \mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \eta \mathbf{D}^k$$
 (6)

$$\mathbf{D}^{k+1} = \mathbf{D}^k + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) \hat{\mathbf{Y}}^k$$
 (7)

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^k; \boldsymbol{\zeta}^k) - \eta \mathbf{D}^{k+1}$$
(8)

where  $\hat{\mathbf{Y}}^k$  represents the compression of  $\mathbf{Y}^k$  and  $\mathbf{F}(\mathbf{X}^k; \xi^k)$  denote the stochastic gradients.

Nevertheless, how to compress the communication and how fast the convergence we can attain with compression error are unknown. In the following, we propose to carefully control the compression error by difference compression and error compensation such that the inexact dual update (Line 6) and primal update (Line 7) can still guarantee the convergence as proved in Section 4.

**Compression error.** Different from existing works, which typically compress the primal variable  $\mathbf{X}^k$  or its difference, LEAD first construct an intermediate variable  $\mathbf{Y}^k$  and apply compression to obtain its coarse representation  $\hat{\mathbf{Y}}^k$  as shown in the procedure COMM( $\mathbf{Y}, \mathbf{H}, \mathbf{H}_w$ ):

- Compress the difference between Y and the state variable H as Q;
- Q is encoded into the low-bit representation, which enables the efficient local communication step  $\hat{\mathbf{Y}}_w = \mathbf{H}_w + \mathbf{W}\mathbf{Q}$ . It is the only communication step in each iteration.
- Each agent recovers its estimate  $\hat{\mathbf{Y}}$  by  $\hat{\mathbf{Y}} = \mathbf{H} + \mathbf{Q}$  and we have  $\hat{\mathbf{Y}}_w = \mathbf{W}\hat{\mathbf{Y}}$ .
- States **H** and  $\mathbf{H}_w$  are updated based on  $\hat{\mathbf{Y}}$  and  $\hat{\mathbf{Y}}_w$ , respectively. We have  $\mathbf{H}_w = \mathbf{W}\mathbf{H}$ .

By this procedure, we expect when both  $\mathbf{Y}^k$  and  $\mathbf{H}^k$  converge to  $\mathbf{X}^*$ , the compression error vanishes asymptotically due to the assumption we make for the compression operator in Assumption 2.

**Remark 1.** Note that difference compression is also applied in DCD-PSGD (Tang et al., 2018a) and CHOCO-SGD (Koloskova et al., 2019), but their state update is the simple integration of the compressed difference. We find this update is usually too aggressive and cause instability as showed in our experiments. Therefore, we adopt a momentum update  $\mathbf{H} = (1 - \alpha)\mathbf{H} + \alpha\hat{\mathbf{Y}}$  motivated from DIANA (Mishchenko et al., 2019), which reduces the compression error for gradient compression in centralized optimization.

Implicit error compensation. On the other hand, even if the compression error exists, LEAD essentially compensates for the error in the inexact dual update (Line 6), making the algorithm more stable and robust. To illustrate how it works, let  $\mathbf{E}^k = \hat{\mathbf{Y}}^k - \mathbf{Y}^k$  denote the compression error and  $e_i^k$  be its *i*-th row. The update of  $\mathbf{D}^k$  gives

$$\mathbf{D}^{k+1} = \mathbf{D}^k + \frac{\gamma}{2\eta}(\hat{\mathbf{Y}}^k - \hat{\mathbf{Y}}_w^k) = \mathbf{D}^k + \frac{\gamma}{2\eta}(\mathbf{I} - \mathbf{W})\mathbf{Y}^k + \frac{\gamma}{2\eta}(\mathbf{E}^k - \mathbf{W}\mathbf{E}^k)$$

where  $-\mathbf{W}\mathbf{E}^k$  indicates that agent i spreads total compression error  $-\sum_{j\in\mathbb{N}_i\cup\{i\}}w_{ji}e_i^k=-e_i^k$  to all agents and  $\mathbf{E}^k$  indicates that each agent compensates this error locally by adding  $e_i^k$  back. This error compensation also explains why the global view in (3) doesn't involve compression error.

**Remark 2.** Note that in LEAD, the compression error is compensated into the model  $\mathbf{X}^{k+1}$  through Line 6 and Line 7 such that the gradient computation in the next iteration is aware of the compression error. This has some subtle but important difference from the error compensation or error feedback in (Seide et al., 2014; Wu et al., 2018; Stich et al., 2018; Karimireddy et al., 2019; Tang et al., 2019b; Liu et al., 2020; Tang et al., 2019a), where the error is stored in the memory and only compensated after gradient computation and before the compression.

**Remark 3.** The proposed algorithm, LEAD in Alg. 1, recovers NIDS (Li et al., 2019),  $D^2$  (Tang et al., 2018b), Exact Diffusion (Yuan et al., 2018). These connections are established in Appendix B.

#### 4 THEORETICAL ANALYSIS

In this section, we show the convergence rate for the proposed algorithm LEAD. Before showing the main theorem, we make some assumptions, which are commonly used for the analysis of decentralized optimization algorithms. All proofs are provided in Appendix E.

**Assumption 2** (Unbiased and C-contracted operator). The compression operator  $Q: \mathbb{R}^d \to \mathbb{R}^d$  is unbiased, i.e.,  $\mathbb{E}Q(\boldsymbol{x}) = \boldsymbol{x}$ , and there exists  $C \geq 0$  such that  $\mathbb{E}\|\boldsymbol{x} - Q(\boldsymbol{x})\|_2^2 \leq C\|\boldsymbol{x}\|_2^2$  for all  $\boldsymbol{x} \in \mathbb{R}^d$ .

**Assumption 3** (Stochastic gradient). The stochastic gradient  $\nabla f_i(\boldsymbol{x};\xi)$  is unbiased, i.e.,  $\mathbb{E}_{\xi}\nabla f_i(\boldsymbol{x};\xi) = \nabla f_i(\boldsymbol{x})$ , and the stochastic gradient variance is bounded:  $\mathbb{E}_{\xi}\|\nabla f_i(\boldsymbol{x};\xi) - \nabla f_i(\boldsymbol{x})\|_2^2 \leq \sigma_i^2$  for all  $i \in [n]$ . Denote  $\sigma^2 = \frac{1}{n}\sum_{i=1}^n \sigma_i^2$ .

**Assumption 4.** Each  $f_i$  is L-smooth and  $\mu$ -strongly convex with  $L \ge \mu > 0$ , i.e., for i = 1, 2, ..., n and  $\forall x, y \in \mathbb{R}^d$ , we have

$$f_i(oldsymbol{y}) + \langle 
abla f_i(oldsymbol{y}), oldsymbol{x} - oldsymbol{y} 
angle + rac{\mu}{2} \|oldsymbol{x} - oldsymbol{y}\|^2 \leq f_i(oldsymbol{x}) \leq f_i(oldsymbol{y}) + \langle 
abla f_i(oldsymbol{y}), oldsymbol{x} - oldsymbol{y} 
angle + rac{L}{2} \|oldsymbol{x} - oldsymbol{y}\|^2.$$

**Theorem 1** (Constant stepsize). Let  $\{\mathbf{X}^k, \mathbf{H}^k, \mathbf{D}^k\}$  be the sequence generated from Alg. 1 and  $\mathbf{X}^*$  is the optimal solution with  $\mathbf{D}^* = -\nabla \mathbf{F}(\mathbf{X}^*)$ . Under Assumptions 1-4, for any constant stepsize  $\eta \in (0, 2/(\mu + L)]$ , if the compression parameters  $\alpha$  and  $\gamma$  satisfy

$$\gamma \in \left(0, \min\left\{\frac{2}{(3C+1)\beta}, \frac{2\mu\eta(2-\mu\eta)}{[2-\mu\eta(2-\mu\eta)]C\beta}\right\}\right),\tag{9}$$

$$\alpha \in \left[ \frac{C\beta\gamma}{2(1+C)}, \frac{1}{a_1} \min \left\{ \frac{2-\beta\gamma}{4-\beta\gamma}, \mu\eta(2-\mu\eta) \right\} \right], \tag{10}$$

with  $\beta \coloneqq \lambda_{\max}(\mathbf{I} - \mathbf{W})$ . Then, in total expectation we have

$$\frac{1}{n}\mathbb{E}\mathcal{L}^{k+1} \le \rho \frac{1}{n}\mathbb{E}\mathcal{L}^k + \eta^2 \sigma^2,\tag{11}$$

where

$$\mathcal{L}^{k} \coloneqq (1 - a_{1}\alpha) \|\mathbf{X}^{k} - \mathbf{X}^{*}\|^{2} + (2\eta^{2}/\gamma) \mathbb{E} \|\mathbf{D}^{k} - \mathbf{D}^{*}\|_{(\mathbf{I} - \mathbf{W})^{\dagger}}^{2} + a_{1} \|\mathbf{H}^{k} - \mathbf{X}^{*}\|^{2},$$

$$\rho \coloneqq \max \left\{ \frac{1 - \mu \eta (2 - \mu \eta)}{1 - a_{1}\alpha}, 1 - \frac{\gamma}{2\lambda_{\max}((\mathbf{I} - \mathbf{W})^{\dagger})}, 1 - \alpha \right\} < 1, a_{1} \coloneqq \frac{4(1 + C)}{C\beta\gamma + 2}$$

The result holds for  $C \to 0$ .

**Corollary 1** (Complexity bounds). Define the condition numbers of the objective function and communication graph as  $\kappa_f = \frac{L}{\mu}$  and  $\kappa_g = \frac{\lambda_{\max}(\mathbf{I} - \mathbf{W})}{\lambda_{\min}^+(\mathbf{I} - \mathbf{W})}$ , respectively. Under the same setting in Theorem 1, we can choose  $\eta = \frac{1}{L}$ ,  $\gamma = \min\{\frac{1}{C\beta\kappa_f}, \frac{1}{(1+3C)\beta}\}$ , and  $\alpha = \mathcal{O}(\frac{1}{(1+C)\kappa_f})$  such that

$$\rho = \max \left\{ 1 - \mathcal{O}\left(\frac{1}{(1+C)\kappa_f}\right), 1 - \mathcal{O}\left(\frac{1}{(1+C)\kappa_a}\right), 1 - \mathcal{O}\left(\frac{1}{C\kappa_f\kappa_a}\right) \right\}.$$

With full-gradient (i.e.,  $\sigma = 0$ ), we obtain the following complexity bounds:

• LEAD converges to the  $\epsilon$ -accurate solution with the iteration complexity

$$\mathcal{O}\Big(\big((1+C)(\kappa_f+\kappa_g)+C\kappa_f\kappa_g\big)\log\frac{1}{\epsilon}\Big).$$

- When C=0 (i.e., there is no compression), we obtain  $\rho=\max\{1-\mathcal{O}(\frac{1}{\kappa_f}),1-\mathcal{O}(\frac{1}{\kappa_g})\}$ , and the iteration complexity  $\mathcal{O}\left((\kappa_f+\kappa_g)\log\frac{1}{\epsilon}\right)$ . This exactly recovers the convergence rate of NIDS (Li et al., 2019).
- When  $C \leq \frac{\kappa_f + \kappa_g}{\kappa_f \kappa_g + \kappa_f + \kappa_g}$ , the asymptotical complexity is  $\mathcal{O}\left((\kappa_f + \kappa_g) \log \frac{1}{\epsilon}\right)$ , which also recovers that of NIDS (Li et al., 2019) and indicates that the compression doesn't harm the convergence in this case.
- With C=0 (or  $C\leq \frac{\kappa_f+\kappa_g}{\kappa_f\kappa_g+\kappa_f+\kappa_g}$ ) and fully connected communication graph (i.e.,  $\mathbf{W}=\frac{\mathbf{1}\mathbf{1}^\top}{n}$ ), we have  $\beta=1$  and  $\kappa_g=1$ . Therefore, we obtain  $\rho=1-\mathcal{O}(\frac{1}{\kappa_f})$  and the complexity bound  $\mathcal{O}(\kappa_f\log\frac{1}{\epsilon})$ . This recovers the convergence rate of gradient descent (Nesterov, 2013).

**Remark 4.** Under the setting in Theorem 1, LEAD converges linearly to the  $\mathcal{O}(\sigma^2)$  neighborhood of the optimum and converges linearly exactly to the optimum if full gradient is used, e.g.,  $\sigma=0$ . The linear convergence of LEAD holds when  $\eta<2/L$ , but we omit the proof.

**Remark 5** (Arbitrary compression precision). Pick any  $\eta \in (0, 2/(\mu + L)]$ , based on the compression-related constant C and the network-related constant  $\beta$ , we can select  $\gamma$  and  $\alpha$  in certain ranges to achieve the convergence. It suggests that LEAD supports unbiased compression with arbitrary precision, i.e., any C > 0.

**Corollary 2** (Consensus error). Under the same setting in Theorem 1, let  $\bar{x}^k = \frac{1}{n} \sum_{i=1}^n x_i^k$  be the averaged model and  $\mathbf{H}^0 = \mathbf{H}^1$ , then all agents achieve consensus at the rate

$$\frac{1}{n} \sum_{i=1}^{n} \mathbb{E} \left\| \boldsymbol{x}_{i}^{k} - \bar{\boldsymbol{x}}^{k} \right\|^{2} \leq \frac{2\mathcal{L}^{0}}{n} \rho^{k} + \frac{2\sigma^{2}}{1-\rho} \eta^{2}. \tag{12}$$

where  $\rho$  is defined as in Corollary 1 with appropriate parameter settings.

**Theorem 2** (Diminishing stepsize). Let  $\{\mathbf{X}^k, \mathbf{H}^k, \mathbf{D}^k\}$  be the sequence generated from Alg. 1 and  $\mathbf{X}^*$  is the optimal solution with  $\mathbf{D}^* = -\nabla \mathbf{F}(\mathbf{X}^*)$ . Under Assumptions 1-4, if  $\eta_k = \frac{2\theta_5}{\theta_3\theta_4\theta_5k+2}$  and  $\gamma_k = \theta_4\eta_k$ , by taking  $\alpha_k = \frac{C\beta\gamma_k}{2(1+C)}$ , in total expectation we have

$$\frac{1}{n} \sum_{i=1}^{n} \mathbb{E} \left\| \boldsymbol{x}_{i}^{k} - \boldsymbol{x}^{*} \right\|^{2} \lesssim \mathcal{O} \left( \frac{1}{k} \right)$$
(13)

where  $\theta_1, \theta_2, \theta_3, \theta_4$  and  $\theta_5$  are constants defined in the proof. The complexity bound for arriving at the  $\epsilon$ -accurate solution is  $\mathcal{O}(\frac{1}{\epsilon})$ .

**Remark 6.** Compared with CHOCO-SGD, LEAD requires unbiased compression and the convergence under biased compression is not investigated yet. The analysis of CHOCO-SGD relies on the bounded gradient assumptions, i.e.,  $\|\nabla f_i(x)\|^2 \leq G$ , which is restrictive because it conflicts with the strong convexity while LEAD doesn't need this assumption. Moreover, in the theorem of CHOCO-SGD, it requires a specific point set of  $\gamma$  while LEAD only requires  $\gamma$  to be within a rather large range. This may explain the advantages of LEAD over CHOCO-SGD in terms of robustness to parameter setting.

#### 5 Numerical Experiment

We consider three machine learning problems –  $\ell_2$ -regularized linear regression, logistic regression, and deep neural network. The proposed LEAD is compared with QDGD (Reisizadeh et al., 2019a), DeepSqueeze (Tang et al., 2019a), CHOCO-SGD (Koloskova et al., 2019), and two non-compressed algorithms DGD (Yuan et al., 2016) and NIDS (Li et al., 2019).

**Setup.** We consider eight machines connected in a ring topology network. Each agent can only exchange information with its two 1-hop neighbors. The mixing weight is simply set as 1/3. For compression, we use the unbiased b-bits quantization method with  $\infty$ -norm

$$Q_{\infty}(\boldsymbol{x}) := \left( \|\boldsymbol{x}\|_{\infty} 2^{-(b-1)} \operatorname{sign}(\boldsymbol{x}) \right) \cdot \left\lfloor \frac{2^{(b-1)} |\boldsymbol{x}|}{\|\boldsymbol{x}\|_{\infty}} + \boldsymbol{u} \right\rfloor, \tag{14}$$

where  $\cdot$  is the Hadamard product, |x| is the elementwise absolute value of x, and u is a random vector uniformly distributed in  $[0,1]^d$ . Only  $\operatorname{sign}(x)$ , norm  $||x||_{\infty}$ , and integers in the bracket need to be transmitted. Note that this quantization method is similar to the quantization used in QSGD (Alistarh et al., 2017) and CHOCO-SGD (Koloskova et al., 2019), but we use the  $\infty$ -norm scaling instead of the 2-norm. This small change brings significant improvement on compression precision as justified both theoretically and empirically in Appendix C. In this section, we choose 2-bit quantization and quantize the data blockwise (block size = 512).

For all experiments, we tune the stepsize  $\eta$  from  $\{0.01, 0.05, 0.1, 0.5\}$ . For QDGD, CHOCO-SGD and Deepsqueeze,  $\gamma$  is tuned from  $\{0.01, 0.1, 0.2, 0.4, 0.6, 0.8, 1.0\}$ . Note that different notations are used in their original papers. Here we uniformly denote the stepsize as  $\eta$  and the additional parameter in these algorithms as  $\gamma$  for simplicity. For LEAD, we simply fix  $\alpha=0.5$  and  $\gamma=1.0$  for all experiments since we find LEAD is robust to parameter settings as we validate in the parameter sensitivity analysis in Appendix D.1. This indicates the minor effort needed for tuning LEAD. Detailed parameter settings for all experiments are summarized in Appendix D.3.

**Linear regression.** We consider the problem:  $f(x) = \sum_{i=1}^n (\|\mathbf{A}_i x - b_i\|^2 + \lambda \|x\|^2)$ . Data matrices  $\mathbf{A}_i \in \mathbb{R}^{200 \times 200}$  and the true solution x' is randomly synthesized. The values  $b_i$  are generated by adding Gaussian noise to  $\mathbf{A}_i x'$ . We let  $\lambda = 0.1$  and the optimal solution of the linear regression problem be  $x^*$ . We use full-batch gradient to exclude the impact of gradient variance. The performance is showed in Fig. 1. The distance to  $x^*$  in Fig. 1a and the consensus error in Fig. 1c verify

that LEAD converges exponentially to the optimal consensual solution. It significantly outperforms most baselines and matches NIDS well under the same number of iterations. Fig. 1b demonstrates the benefit of compression when considering the communication bits. Fig. 1d shows that the compression error vanishes for both LEAD and CHOCO-SGD while the compression error is pretty large for QDGD and DeepSqueeze because they directly compress the local models.

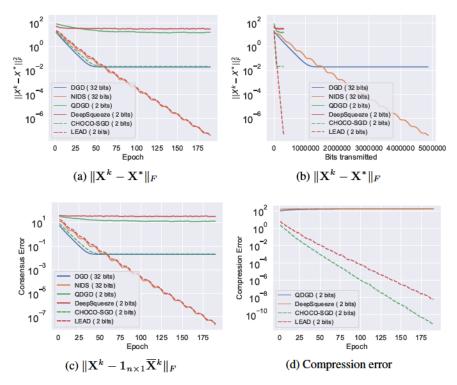


Figure 1: Linear regression problem.

**Logistic regression.** We further consider a logistic regression problem on the MNIST dataset. The regularization parameter is  $10^{-4}$ . We consider both *homogeneous* and *heterogeneous* data settings. In the homogeneous setting, the data samples are randomly shuffled before being uniformly partitioned among all agents such that the data distribution from each agent is very similar. In the heterogeneous setting, the samples are first sorted by their labels and then partitioned among agents. Due to the space limit, we mainly present the results in heterogeneous setting here and defer the homogeneous setting to Appendix D.2. The results using full-batch gradient and mini-batch gradient (the mini-batch size is 512 for each agent) are showed in Fig. 2 and Fig. 3 respectively and both settings shows the faster convergence and higher precision of LEAD.

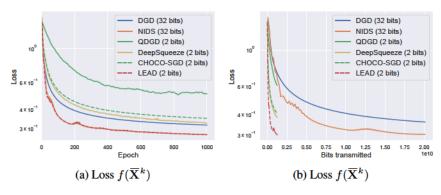


Figure 2: Logistic regression problem in the heterogeneous case (full-batch gradient).

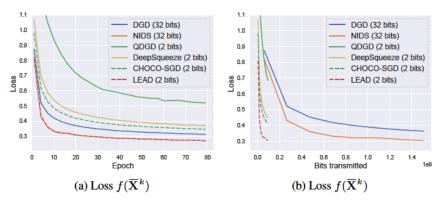


Figure 3: Logistic regression in the heterogeneous case (mini-batch gradient).

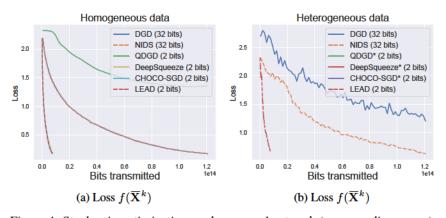


Figure 4: Stochastic optimization on deep neural network (\* means divergence).

Neural network. We empirically study the performance of LEAD in optimizing deep neural network by training AlexNet (240 MB) on CIFAR10 dataset. The mini-batch size is 64 for each agents. Both the homogeneous and heterogeneous case are showed in Fig. 4. In the homogeneous case, CHOCO-SGD, DeepSqueeze and LEAD perform similarly and outperform the non-compressed variants in terms of communication efficiency, but CHOCO-SGD and DeepSqueeze need more efforts for parameter tuning because their convergence is sensitive to the setting of  $\gamma$ . In the heterogeneous cases, LEAD achieves the fastest and most stable convergence. Note that in this setting, sufficient information exchange is more important for convergence because models from different agents are moving to significantly diverse directions. In such case, DGD only converges with smaller stepsize and its communication compressed variants, including QDGD, DeepSqueeze and CHOCO-SGD, diverge in all parameter settings we try.

In summary, our experiments verify our theoretical analysis and show that LEAD is able to handle data heterogeneity very well. Furthermore, the performance of LEAD is robust to parameter settings and needs less effort for parameter tuning, which is critical in real-world applications.

#### 6 CONCLUSION

In this paper, we investigate the communication compression in decentralized optimization. Motivated by primal-dual algorithms, a novel decentralized algorithm with compression, LEAD, is proposed to achieve faster convergence rate and to better handle heterogeneous data while enjoying the benefit of efficient communication. The nontrivial analyses on the coupled dynamics of inexact primal and dual updates as well as compression error establish the linear convergence of LEAD when full gradient is used and the linear convergence to the  $\mathcal{O}(\sigma^2)$  neighborhood of the optimum when stochastic gradient is used. Extensive experiments validate the theoretical analysis and demonstrate the state-of-the-art efficiency and robustness of LEAD. LEAD is also applicable to non-convex problems as empirically verified in the neural network experiments but we leave the non-convex analysis as the future work.

#### **ACKNOWLEDGEMENTS**

Xiaorui Liu and Dr. Jiliang Tang are supported by the National Science Foundation (NSF) under grant numbers CNS-1815636, IIS-1928278, IIS-1714741, IIS-1845081, IIS-1907704, and IIS-1955285. Yao Li and Dr. Ming Yan are supported by NSF grant DMS-2012439 and Facebook Faculty Research Award (Systems for ML). Dr. Rongrong Wang is supported by NSF grant CCF-1909523.

#### REFERENCES

- Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient sgd via gradient quantization and encoding. In Advances in Neural Information Processing Systems, pp. 1709–1720. 2017.
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. SIGNSGD: compressed optimisation for non-convex problems. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 559–568, 2018.
- Ruggero Carli, Fabio Fagnani, Paolo Frasca, and Sandro Zampieri. Gossip consensus algorithms via quantized communication. *Automatica*, 46(1):70–80, 2010.
- Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Urban Stich, and Martin Jaggi. Error feedback fixes SignSGD and other gradient compression schemes. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 3252–3261. PMLR, 2019.
- Anastasia Koloskova, Sebastian U. Stich, and Martin Jaggi. Decentralized stochastic optimization and gossip algorithms with compressed communication. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 3479–3487. PMLR, 2019.
- Anastasia Koloskova, Tao Lin, Sebastian U Stich, and Martin Jaggi. Decentralized deep learning with arbitrary communication compression. In *International Conference on Learning Representations*, 2020.
- Yao Li and Ming Yan. On linear convergence of two decentralized algorithms. *arXiv* preprint *arXiv*:1906.07225, 2019.
- Zhi Li, Wei Shi, and Ming Yan. A decentralized proximal-gradient method with network independent step-sizes and separated convergence rates. *IEEE Transactions on Signal Processing*, 67 (17):4494–4506, 2019.
- Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pp. 5330–5340, 2017.
- Qing Ling, Wei Shi, Gang Wu, and Alejandro Ribeiro. DLM: Decentralized linearized alternating direction method of multipliers. *IEEE Transactions on Signal Processing*, 63(15):4051–4064, 2015.
- Xiaorui Liu, Yao Li, Jiliang Tang, and Ming Yan. A double residual compression algorithm for efficient distributed learning. *The 23rd International Conference on Artificial Intelligence and Statistics*, 2020.
- Yucheng Lu and Christopher De Sa. Moniqua: Modulo quantized communication in decentralized SGD. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- Sindri Magnússon, Hossein Shokri-Ghadikolaei, and Na Li. On maintaining linear convergence of distributed learning and optimization under limited communication. *IEEE Transactions on Signal Processing*, 68:6101–6116, 2020.
- Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.

- Joao FC Mota, Joao MF Xavier, Pedro MQ Aguiar, and Markus Püschel. D-ADMM: A communication-efficient distributed algorithm for separable optimization. *IEEE Transactions on Signal Processing*, 61(10):2718–2723, 2013.
- Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- Angelia Nedic, Alex Olshevsky, and Wei Shi. Achieving geometric convergence for distributed optimization over time-varying graphs. *SIAM Journal on Optimization*, 27(4):2597–2633, 2017.
- Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2013.
- Shi Pu and Angelia Nedić. Distributed stochastic gradient tracking methods. *Mathematical Programming*, pp. 1–49, 2020.
- Amirhossein Reisizadeh, Aryan Mokhtari, Hamed Hassani, and Ramtin Pedarsani. An exact quantized decentralized gradient descent algorithm. *IEEE Transactions on Signal Processing*, 67(19): 4934–4947, 2019a.
- Amirhossein Reisizadeh, Hossein Taheri, Aryan Mokhtari, Hamed Hassani, and Ramtin Pedarsani. Robust and communication-efficient collaborative learning. In *Advances in Neural Information Processing Systems*, pp. 8388–8399, 2019b.
- Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and application to data-parallel distributed training of speech DNNs. In *Interspeech 2014*, September 2014.
- Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. EXTRA: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- Sebastian U. Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD with memory. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 4452–4463, 2018.
- Hanlin Tang, Shaoduo Gan, Ce Zhang, Tong Zhang, and Ji Liu. Communication compression for decentralized training. In Advances in Neural Information Processing Systems, pp. 7652–7662. 2018a.
- Hanlin Tang, Xiangru Lian, Ming Yan, Ce Zhang, and Ji Liu.  $D^2$ : Decentralized training over decentralized data. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 4848–4856, 2018b.
- Hanlin Tang, Xiangru Lian, Shuang Qiu, Lei Yuan, Ce Zhang, Tong Zhang, and Ji Liu. Deepsqueeze: Decentralization meets error-compensated compression. *CoRR*, abs/1907.07346, 2019a. URL http://arxiv.org/abs/1907.07346.
- Hanlin Tang, Chen Yu, Xiangru Lian, Tong Zhang, and Ji Liu. DoubleSqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 6155–6165, 2019b.
- John Tsitsiklis, Dimitri Bertsekas, and Michael Athans. Distributed asynchronous deterministic and stochastic gradient optimization algorithms. *IEEE transactions on automatic control*, 31(9): 803–812, 1986.
- Jiaxiang Wu, Weidong Huang, Junzhou Huang, and Tong Zhang. Error compensated quantized SGD and its applications to large-scale distributed optimization. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 5325–5333, 2018.
- Lin Xiao and Stephen Boyd. Fast linear iterations for distributed averaging. *Systems & Control Letters*, 53(1):65–78, 2004.
- Jinming Xu, Ye Tian, Ying Sun, and Gesualdo Scutari. Accelerated primal-dual algorithms for distributed smooth convex optimization over networks. In *International Conference on Artificial Intelligence and Statistics*, pp. 2381–2391. PMLR, 2020.

- Kun Yuan, Qing Ling, and Wotao Yin. On the convergence of decentralized gradient descent. *SIAM Journal on Optimization*, 26(3):1835–1854, 2016.
- Kun Yuan, Bicheng Ying, Xiaochuan Zhao, and Ali H Sayed. Exact diffusion for distributed optimization and learning—part i: Algorithm development. *IEEE Transactions on Signal Processing*, 67(3):708–723, 2018.
- Kun Yuan, Wei Xu, and Qing Ling. Can primal methods outperform primal-dual methods in decentralized dynamic optimization? *arXiv preprint arXiv:2003.00816*, 2020.

# **Contents of Appendix**

A	LEA	AD in agent's perspective	14
В	Con	nections with exiting works	14
C	Com	npression method	15
	C.1	p-norm b-bits quantization	15
	C.2	Compression error	16
D	Exp	eriments	17
	D.1	Parameter sensitivity	17
	D.2	Experiments in homogeneous setting	17
	D.3	Parameter settings	18
E	Proc	ofs of the theorems	19
	E.1	Illustrative flow	19
	E.2	Two central Lemmas	20
	E.3	Proof of Lemma 1	20
	E.4	Proof of Lemma 2	22
	E.5	Proof of Theorem 1	23
	E.6	Proof of Theorem 2	28

### LEAD IN AGENT'S PERSPECTIVE

In the main paper, we described the algorithm with matrix notations for concision. Here we further provide a complete algorithm description from the agents' perspective.

#### **Algorithm 2** LEAD in Agent's Perspective

**input:** stepsize  $\eta$ , compression parameters  $(\alpha, \gamma)$ , initial values  $x_i^0, h_i^1, z_i, \forall i \in \{1, 2, ..., n\}$ **output:**  $\boldsymbol{x}_i^K, \ \forall i \in \{1, 2, \dots, n\}$  or  $\frac{\sum_{i=1}^n \boldsymbol{x}_i^K}{n}$ 

- 1: for each agent  $i \in \{1, 2, \dots, n\}$  do
- $d_i^1 = z_i \sum_{j \in \mathcal{N}_i \cup \{i\}} w_{ij} z_j$
- $(\boldsymbol{h}_w)_i^1 = \sum_{j \in \mathcal{N}_i \cup \{i\}} w_{ij} (\boldsymbol{h}_w)_j^1$  $\boldsymbol{x}_i^1 = \boldsymbol{x}_i^0 \eta \nabla f_i(\boldsymbol{x}_i^0; \boldsymbol{\xi}_i^0)$

- 6: **for**  $k = 1, 2, \dots, K 1$  **do** in parallel for all agents  $i \in \{1, 2, \dots, n\}$
- compute  $\nabla f_i(\boldsymbol{x}_i^k; \boldsymbol{\xi}_i^k)$

- $\boldsymbol{y}_{i}^{k} = \boldsymbol{x}_{i}^{k} \eta \nabla f_{i}(\boldsymbol{x}_{i}^{k}; \xi_{i}^{k}) \eta \boldsymbol{d}_{i}^{k}$
- $q_i^k = \text{Compress}(y_i^k h_i^k)$
- $\hat{m{y}}_i^k = m{h}_i^k + m{q}_i^k$ 10:
- for neighbors  $j \in \mathbb{N}_i$  do 11:
- Send  $q_i^k$  and receive  $q_i^k$

- 13:

- 16:
- $\begin{aligned} &(\hat{\boldsymbol{y}}_{w})_{i}^{k} = (\boldsymbol{h}_{w})_{i}^{k} + \sum_{j \in \mathbb{N}_{i} \cup \{i\}} w_{ij} \boldsymbol{q}_{j}^{k} \\ &\boldsymbol{h}_{i}^{k+1} = (1 \alpha) \boldsymbol{h}_{i}^{k} + \alpha \hat{\boldsymbol{y}}_{i}^{k} \\ &(\boldsymbol{h}_{w})_{i}^{k+1} = (1 \alpha) (\boldsymbol{h}_{w})_{i}^{k} + \alpha (\hat{\boldsymbol{y}}_{w})_{i}^{k} \\ &\boldsymbol{d}_{i}^{k+1} = \boldsymbol{d}_{i}^{k} + \frac{\gamma}{2\eta} (\hat{\boldsymbol{y}}_{i}^{k} (\hat{\boldsymbol{y}}_{w})_{i}^{k}) \\ &\boldsymbol{x}_{i}^{k+1} = \boldsymbol{x}_{i}^{k} \eta \nabla f_{i}(\boldsymbol{x}_{i}^{k}; \boldsymbol{\xi}_{i}^{k}) \eta \boldsymbol{d}_{i}^{k+1} \end{aligned}$ 18:
- 19: **end for**

#### В CONNECTIONS WITH EXITING WORKS

The non-compressed variant of LEAD in Alg. 1 recovers NIDS (Li et al., 2019),  $D^2$  (Tang et al., 2018b) and Exact Diffusion (Yuan et al., 2018) as shown in Proposition 1. In Corollary 3, we show that the convergence rate of LEAD exactly recovers the rate of NIDS when  $C=0, \gamma=1$  and  $\sigma=0$ .

**Proposition 1** (Connection to NIDS,  $D^2$  and Exact Diffusion). When there is no communication compression (i.e.,  $\hat{\mathbf{Y}}^k = \mathbf{Y}^k$ ) and  $\gamma = 1$ , Alg. 1 recovers  $D^2$ :

$$\mathbf{X}^{k+1} = \frac{\mathbf{I} + \mathbf{W}}{2} \left( 2\mathbf{X}^k - \mathbf{X}^{k-1} - \eta \nabla \mathbf{F}(\mathbf{X}^k; \xi^k) + \eta \nabla \mathbf{F}(\mathbf{X}^{k-1}; \xi^{k-1}) \right). \tag{15}$$

Furthermore, if the stochastic estimator of the gradient  $\nabla \mathbf{F}(\mathbf{X}^k; \xi^k)$  is replaced by the full gradient, it recovers NIDS and Exact Diffusion with specific settings.

**Corollary 3** (Consistency with NIDS). When C=0 (no communication compression),  $\gamma=1$  and  $\sigma = 0$  (full gradient), LEAD has the convergence consistent with NIDS with  $\eta \in (0, 2/(\mu + L)]$ :

$$\mathcal{L}^{k+1} \le \max \left\{ 1 - \mu(2\eta - \mu\eta^2), 1 - \frac{1}{2\lambda_{\max}((\mathbf{I} - \mathbf{W})^{\dagger})} \right\} \mathcal{L}^k.$$
 (16)

See the proof in E.5.

*Proof of Proposition 1.* Let  $\gamma = 1$  and  $\hat{\mathbf{Y}}^k = \mathbf{Y}^k$ . Combing Lines 4 and 6 of Alg. 1 gives

$$\mathbf{D}^{k+1} = \mathbf{D}^k + \frac{\mathbf{I} - \mathbf{W}}{2\eta} (\mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^k; \xi^k) - \eta \mathbf{D}^k). \tag{17}$$

Based on Line 7, we can represent  $\eta \mathbf{D}^k$  from the previous iteration as

$$\eta \mathbf{D}^k = \mathbf{X}^{k-1} - \mathbf{X}^k - \eta \nabla \mathbf{F}(\mathbf{X}^{k-1}; \xi^{k-1}). \tag{18}$$

Eliminating both  $\mathbf{D}^k$  and  $\mathbf{D}^{k+1}$  by substituting (17)-(18) into Line 7, we obtain

$$\mathbf{X}^{k+1} = \mathbf{X}^{k} - \eta \nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k}) - \left(\eta \mathbf{D}^{k} + \frac{\mathbf{I} - \mathbf{W}}{2} (\mathbf{X}^{k} - \eta \nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k}) - \eta \mathbf{D}^{k})\right) \quad (\text{from } (17))$$

$$= \frac{\mathbf{I} + \mathbf{W}}{2} (\mathbf{X}^{k} - \eta \nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k})) - \frac{\mathbf{I} + \mathbf{W}}{2} \eta \mathbf{D}^{k}$$

$$= \frac{\mathbf{I} + \mathbf{W}}{2} (\mathbf{X}^{k} - \eta \nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k})) - \frac{\mathbf{I} + \mathbf{W}}{2} (\mathbf{X}^{k-1} - \mathbf{X}^{k} - \eta \nabla \mathbf{F}(\mathbf{X}^{k-1}; \boldsymbol{\xi}^{k-1})) \quad (\text{from } (18))$$

$$= \frac{\mathbf{I} + \mathbf{W}}{2} (2\mathbf{X}^{k} - \mathbf{X}^{k-1} - \eta \nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k}) + \eta \nabla \mathbf{F}(\mathbf{X}^{k-1}; \boldsymbol{\xi}^{k-1})), \quad (19)$$

which is exactly  $D^2$ . It also recovers Exact Diffusion with  $\mathbf{A} = \frac{\mathbf{I} + \mathbf{W}}{2}$  and  $\mathbf{M} = \eta \mathbf{I}$  in Eq. (97) of (Yuan et al., 2018).

#### C COMPRESSION METHOD

#### C.1 P-NORM B-BITS QUANTIZATION

**Theorem 3** (p-norm b-bit quantization). Let us define the quantization operator as

$$Q_p(\boldsymbol{x}) := \left( \|\boldsymbol{x}\|_p \operatorname{sign}(\boldsymbol{x}) 2^{-(b-1)} \right) \cdot \left| \frac{2^{b-1} |\boldsymbol{x}|}{\|\boldsymbol{x}\|_p} + \boldsymbol{u} \right|$$
(20)

where  $\cdot$  is the Hadamard product,  $|\mathbf{x}|$  is the elementwise absolute value and  $\mathbf{u}$  is a random dither vector uniformly distributed in  $[0,1]^d$ .  $Q_p(\mathbf{x})$  is unbiased, i.e.,  $\mathbb{E}Q_p(\mathbf{x}) = \mathbf{x}$ , and the compression variance is upper bounded by

$$\mathbb{E}\|\boldsymbol{x} - Q_p(\boldsymbol{x})\|^2 \le \frac{1}{4}\|\operatorname{sign}(\boldsymbol{x})2^{-(b-1)}\|^2\|\boldsymbol{x}\|_p^2,\tag{21}$$

which suggests that  $\infty$ -norm provides the smallest upper bound for the compression variance due to  $\|\mathbf{x}\|_p \leq \|\mathbf{x}\|_q$ ,  $\forall \mathbf{x}$  if  $1 \leq q \leq p \leq \infty$ .

Remark 7. For the compressor defined in (20), we have the following the compression constant

$$C = \sup_{\mathbf{x}} \frac{\|\operatorname{sign}(\mathbf{x})2^{-(b-1)}\|^2 \|\mathbf{x}\|_p^2}{4\|\mathbf{x}\|^2}.$$

*Proof.* Let denote  $\mathbf{v} = \|\mathbf{x}\|_p \operatorname{sign}(\mathbf{x}) 2^{-(b-1)}$ ,  $s = \frac{2^{b-1}|\mathbf{x}|}{\|\mathbf{x}\|_p}$ ,  $s_1 = \left\lfloor \frac{2^{b-1}|\mathbf{x}|}{\|\mathbf{x}\|_p} \right\rfloor$  and  $s_2 = \left\lceil \frac{2^{b-1}|\mathbf{x}|}{\|\mathbf{x}\|_p} \right\rceil$ . We can rewrite  $\mathbf{x}$  as  $\mathbf{x} = s \cdot \mathbf{v}$ .

For any coordinate i such that  $s_i = (s_1)_i$ , we have  $Q_p(\boldsymbol{x}_i) = (s_1)_i \boldsymbol{v}_i$  with probability 1. Hence  $\mathbb{E}Q_p(\boldsymbol{x})_i = s_i \boldsymbol{v}_i = \boldsymbol{x}_i$  and

$$\mathbb{E}(\boldsymbol{x}_i - Q_p(\boldsymbol{x})_i)^2 = (\boldsymbol{x}_i - s_i \boldsymbol{v}_i)^2 = 0.$$

For any coordinate i such that  $s_i \neq (s_1)_i$ , we have  $(s_2)_i - (s_1)_i = 1$  and  $Q_p(x)_i$  satisfies

$$Q_p(\boldsymbol{x})_i = \begin{cases} (s_1)_i \boldsymbol{v}_i, & \text{w.p. } (s_2)_i - s_i, \\ (s_2)_i \boldsymbol{v}_i, & \text{w.p. } s_i - (s_1)_i. \end{cases}$$

Thus, we derive

$$\mathbb{E}Q_n(\mathbf{x})_i = \mathbf{v}_i(s_1)_i(s_2 - s)_i + \mathbf{v}_i(s_2)_i(s - s_1)_i = \mathbf{v}_i s_i(s_2 - s_1)_i = \mathbf{v}_i s_i = \mathbf{x}_i,$$

and

$$\begin{split} \mathbb{E}[x_i - Q_p(x)_i]^2 &= (x_i - v_i(s_1)_i)^2 (s_2 - s)_i + (x_i - v_i(s_2)_i)^2 (s - s_1)_i \\ &= (s_2 - s_1)_i x_i^2 + \left( (s_1)_i (s_2)_i (s_1 - s_2)_i + s_i ((s_2)_i^2 - (s_1)_i^2) \right) v_i^2 - 2s_i (s_2 - s_1)_i x_i v_i \\ &= x_i^2 + \left( -(s_1)_i (s_2)_i + s_i (s_2 + s_1)_i \right) v_i^2 - 2s_i x_i v_i \\ &= (x_i - s_i v_i)^2 + \left( -(s_1)_i (s_2)_i + s_i (s_2 + s_1)_i - s_i^2 \right) v_i^2 \\ &= (x_i - s_i v_i)^2 + (s_2 - s)_i (s - s_1)_i v_i^2 \\ &= (s_2 - s)_i (s - s_1)_i v_i^2 \\ &\leq \frac{1}{4} v_i^2. \end{split}$$

Considering both cases, we have  $\mathbb{E}Q(x) = x$  and

$$\begin{split} \mathbb{E}\|x - Q_p(x)\|^2 &= \sum_{\{s_i = (s_1)_i\}} \mathbb{E}[x_i - Q_p(x)_i]^2 + \sum_{\{s_i \neq (s_1)_i\}} \mathbb{E}[x_i - Q_p(x)_i]^2 \\ &\leq 0 + \frac{1}{4} \sum_{\{s_i \neq (s_1)_i\}} v_i^2 \\ &\leq \frac{1}{4} \|v\|^2 \\ &= \frac{1}{4} \|\operatorname{sign}(x) 2^{-(b-1)} \|^2 \|x\|_p^2. \end{split}$$

#### C.2 COMPRESSION ERROR

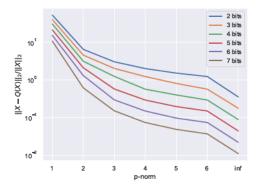


Figure 5: Relative compression error  $\frac{\|\boldsymbol{x}-Q(\boldsymbol{x})\|_2}{\|\boldsymbol{x}\|_2}$  for p-norm b-bit quantization

To verify Theorem 3, we compare the compression error of the quantization method defined in (20) with different norms  $(p=1,2,3,\ldots,6,\infty)$ . Specifically, we uniformly generate 100 random vectors in  $\mathbb{R}^{10000}$  and compute the average compression error. The result shown in Figure 5 verifies our proof in Theorem 3 that the compression error decreases when p increases. This suggests that  $\infty$ -norm provides the best compression precision under the same bit constraint.

Under similar setting, we also compare the compression error with other popular compression methods, such as top-k and random-k sparsification. The x-axes represents the average bits needed to represent each element of the vector. The result is showed in Fig. 6. Note that intuitively top-k methods should perform better than random-k method, but the top-k method needs extra bits to transmitted the index while random-k method can avoid this by using the same random seed. Therefore, top-k method doesn't outperform random-k too much under the same communication budget. The result in Fig. 6 suggests that  $\infty$ -norm b-bits quantization provides significantly better compression precision than others under the same bit constraint.

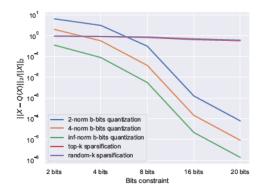


Figure 6: Comparison of compression error  $\frac{\|x-Q(x)\|_2}{\|x\|_2}$  between different compression methods

#### D EXPERIMENTS

#### D.1 PARAMETER SENSITIVITY

In the linear regression problem, the convergence of LEAD under different parameter settings of  $\alpha$  and  $\gamma$  are tested. The result showed in Figure 7 indicates that LEAD performs well in most settings and is robust to the parameter setting. Therefore, in this paper, we simply set  $\alpha=0.5$  and  $\gamma=1.0$  for LEAD in all experiment, which indicates the minor effort needed for parameter tuning.

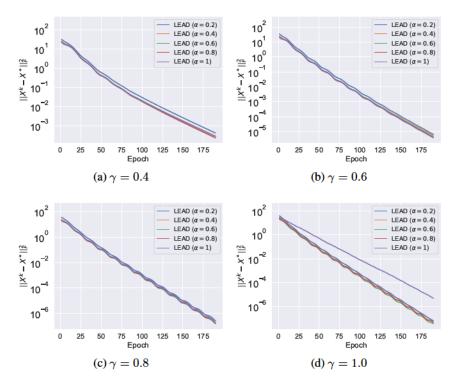


Figure 7: Parameter analysis on linear regression problem.

#### D.2 EXPERIMENTS IN HOMOGENEOUS SETTING

The experiments on logistic regression problem in homogeneous case are showed in Fig. 8 and Fig. 9. It shows that DeepSqueeze, CHOCO-SGD and LEAD converges similarly while Deep-

Squeeze and CHOCO-SGD require to tune a smaller  $\gamma$  for convergence as showed in the parameter setting in Section D.3. Generally, a smaller  $\gamma$  decreases the model propagation between agents since  $\gamma$  changes the effective mixing matrix and this may cause slower convergence. However, in the setting where data from different agents are very similar, the models move to close directions such that the convergence is not affected too much.

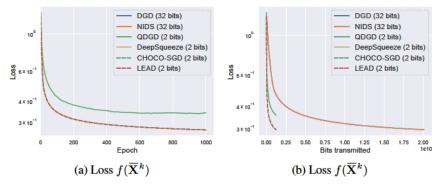


Figure 8: Logistic regression in the homogeneous case (full-batch gradient)

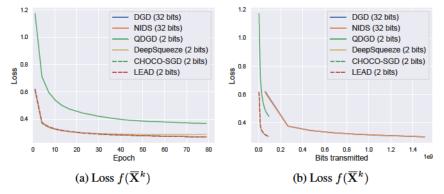


Figure 9: Logistic regression in the homogeneous case (mini-batch gradient)

#### D.3 PARAMETER SETTINGS

The best parameter settings we search for all algorithms and experiments are summarized in Tables 1– 4. QDGD and DeepSqueeze are more sensitive to  $\gamma$  and CHOCO-SGD is slight more robust. LEAD is most robust to parameter settings and it works well for the setting  $\alpha=0.5$  and  $\gamma=1.0$  in all experiments in this paper.

Algorithm	$\eta$	$\gamma$	$\alpha$
DGD	0.1	-	-
NIDS	0.1	-	-
QDGD	0.1	0.2	-
DeepSqueeze	0.1	0.2	-
CHOCO-SGD	0.1	0.8	-
LEAD	0.1	1.0	0.5

Table 1: Parameter settings for the linear regression problem.

Algorithm	$\eta$	$\gamma$	$\alpha$
DGD	0.1	-	-
NIDS	0.1	-	-
QDGD	0.1	0.4	-
DeepSqueeze	0.1	0.4	-
CHOCO-SGD	0.1	0.6	-
LEAD	0.1	1.0	0.5

Algorithm	$\eta$	$\gamma$	$\alpha$
DGD	0.1	-	-
NIDS	0.1	-	-
QDGD	0.1	0.2	-
DeepSqueeze	0.1	0.6	-
CHOCO-SGD	0.1	0.6	-
LEAD	0.1	1.0	0.5

Homogeneous case

Heterogeneous case

Table 2: Parameter settings for the logistic regression problem (full-batch gradient).

Algorithm	$\eta$	$\gamma$	$\alpha$
DGD	0.1	-	-
NIDS	0.1	-	-
QDGD	0.05	0.2	-
DeepSqueeze	0.1	0.6	-
CHOCO-SGD	0.1	0.6	-
LEAD	0.1	1.0	0.5

Algorithm	$\eta$	$\gamma$	$\alpha$
DGD	0.1	-	-
NIDS	0.1	-	-
QDGD	0.05	0.2	-
DeepSqueeze	0.1	0.6	-
CHOCO-SGD	0.1	0.6	-
LEAD	0.1	1.0	0.5

Homogeneous case

Heterogeneous case

Table 3: Parameter settings for the logistic regression problem (mini-batch gradient).

Algorithm	$\eta$	$\gamma$	$\alpha$
DGD	0.1	-	-
NIDS	0.1	-	-
QDGD	0.05	0.1	-
DeepSqueeze	0.1	0.2	-
CHOCO-SGD	0.1	0.6	-
LEAD	0.1	1.0	0.5

Algorithm	$\eta$	$\gamma$	$\alpha$
DGD	0.05	-	-
NIDS	0.1	-	-
QDGD	*	*	-
DeepSqueeze	*	*	-
CHOCO-SGD	*	*	-
LEAD	0.1	1.0	0.5

Homogeneous case

Heterogeneous case

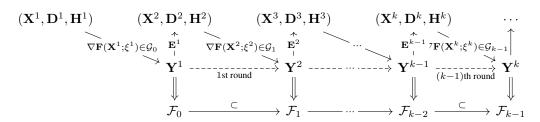
Table 4: Parameter settings for the deep neural network. (\* means divergence for all options we try)

#### E PROOFS OF THE THEOREMS

#### E.1 ILLUSTRATIVE FLOW

The following flow graph depicts the relation between iterative variables and clarifies the range of conditional expectation.  $\{\mathcal{G}_k\}_{k=0}^{\infty}$  and  $\{\mathcal{F}_k\}_{k=0}^{\infty}$  are two  $\sigma$ -algebras generated by the gradient sampling and the stochastic compression respectively. They satisfy

$$\mathcal{G}_0 \subset \mathcal{F}_0 \subset \mathcal{G}_1 \subset \mathcal{F}_1 \subset \cdots \subset \mathcal{G}_k \subset \mathcal{F}_k \subset \cdots$$



The solid and dashed arrows in the top flow illustrate the dynamics of the algorithm, while in the bottom, the arrows stand for the relation between successive  $\mathcal{F}$ - $\sigma$ -algebras. The downward arrows

determine the range of  $\mathcal{F}$ - $\sigma$ -algebras. E.g., up to  $\mathbf{E}^k$ , all random variables are in  $\mathcal{F}_{k-1}$  and up to  $\nabla \mathbf{F}(\mathbf{X}^k; \xi^k)$ , all random variables are in  $\mathcal{G}_{k-1}$  with  $\mathcal{G}_{k-1} \subset \mathcal{F}_{k-1}$ . Throughout the appendix, without specification,  $\mathbb{E}$  is the expectation conditioned on the corresponding stochastic estimators given the context.

#### E.2 TWO CENTRAL LEMMAS

**Lemma 1** (Fundamental equality). Let  $\mathbf{X}^*$  be the optimal solution,  $\mathbf{D}^* := -\nabla \mathbf{F}(\mathbf{X}^*)$  and  $\mathbf{E}^k$  denote the compression error in the kth iteration, that is  $\mathbf{E}^k = \mathbf{Q}^k - (\mathbf{Y}^k - \mathbf{H}^k) = \hat{\mathbf{Y}}^k - \mathbf{Y}^k$ . From Alg. 1, we have

$$\begin{split} &\|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2 + (\eta^2/\gamma)\|\mathbf{D}^{k+1} - \mathbf{D}^*\|_{\mathbf{M}}^2 \\ = &\|\mathbf{X}^k - \mathbf{X}^*\|^2 + (\eta^2/\gamma)\|\mathbf{D}^k - \mathbf{D}^*\|_{\mathbf{M}}^2 - (\eta^2/\gamma)\|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2 - \eta^2\|\mathbf{D}^{k+1} - \mathbf{D}^*\|^2 \\ &- 2\eta\langle\mathbf{X}^k - \mathbf{X}^*, \nabla\mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \nabla\mathbf{F}(\mathbf{X}^*)\rangle + \eta^2\|\nabla\mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \nabla\mathbf{F}(\mathbf{X}^*)\|^2 + 2\eta\langle\mathbf{E}^k, \mathbf{D}^{k+1} - \mathbf{D}^*\rangle, \end{split}$$

where  $\mathbf{M} := 2(\mathbf{I} - \mathbf{W})^{\dagger} - \gamma \mathbf{I}$  and  $\gamma < 2/\lambda_{\max}(\mathbf{I} - \mathbf{W})$  ensures the positive definiteness of  $\mathbf{M}$  over range( $\mathbf{I} - \mathbf{W}$ ).

**Lemma 2** (State inequality). Let the same assumptions in Lemma 1 hold. From Alg. 1, if we take the expectation over the compression operator conditioned on the k-th iteration, we have

$$\begin{split} \mathbb{E}\|\mathbf{H}^{k+1} - \mathbf{X}^*\|^2 &\leq (1-\alpha)\|\mathbf{H}^k - \mathbf{X}^*\|^2 + \alpha \mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2 + \alpha \eta^2 \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|^2 \\ &+ \frac{2\alpha\eta^2}{\gamma} \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2 + \alpha^2 \mathbb{E}\|\mathbf{E}^k\|^2 - \alpha\gamma \mathbb{E}\|\mathbf{E}^k\|_{\mathbf{I}-\mathbf{W}}^2 - \alpha(1-\alpha)\|\mathbf{Y}^k - \mathbf{H}^k\|^2. \end{split}$$

#### E.3 PROOF OF LEMMA 1

Before proving Lemma 1, we let  $\mathbf{E}^k = \hat{\mathbf{Y}}^k - \mathbf{Y}^k$  and introduce the following three Lemmas.

**Lemma 3.** Let  $X^*$  be the consensus solution. Then, from Line 4-7 of Alg. 1, we obtain

$$\frac{\mathbf{I} - \mathbf{W}}{2\eta} (\mathbf{X}^{k+1} - \mathbf{X}^*) = \left(\frac{I}{\gamma} - \frac{\mathbf{I} - \mathbf{W}}{2}\right) (\mathbf{D}^{k+1} - \mathbf{D}^k) - \frac{\mathbf{I} - \mathbf{W}}{2\eta} \mathbf{E}^k. \tag{22}$$

*Proof.* From the iterations in Alg. 1, we have

$$\begin{split} \mathbf{D}^{k+1} &= \mathbf{D}^k + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) \hat{\mathbf{Y}}^k & \text{(from Line 6)} \\ &= \mathbf{D}^k + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) (\mathbf{Y}^k + \mathbf{E}^k) \\ &= \mathbf{D}^k + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) (\mathbf{X}^k - \eta \nabla \mathbf{F} (\mathbf{X}^k; \boldsymbol{\xi}^k) - \eta \mathbf{D}^k + \mathbf{E}^k) & \text{(from Line 4)} \\ &= \mathbf{D}^k + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) (\mathbf{X}^k - \eta \nabla \mathbf{F} (\mathbf{X}^k; \boldsymbol{\xi}^k) - \eta \mathbf{D}^{k+1} - \mathbf{X}^* + \eta (\mathbf{D}^{k+1} - \mathbf{D}^k) + \mathbf{E}^k) \\ &= \mathbf{D}^k + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) (\mathbf{X}^{k+1} - \mathbf{X}^*) + \frac{\gamma}{2} (\mathbf{I} - \mathbf{W}) (\mathbf{D}^{k+1} - \mathbf{D}^k) + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) \mathbf{E}^k, \end{split}$$

where the fourth equality holds due to  $(\mathbf{I} - \mathbf{W})\mathbf{X}^* = \mathbf{0}$  and the last equality comes from Line 7 of Alg. 1. Rewriting this equality, and we obtain (22).

**Lemma 4.** Let  $\mathbf{D}^* = -\nabla \mathbf{F}(\mathbf{X}^*) \in \mathbf{span}\{\mathbf{I} - \mathbf{W}\}$ , we have

$$\langle \mathbf{X}^{k+1} - \mathbf{X}^*, \mathbf{D}^{k+1} - \mathbf{D}^k \rangle = \frac{\eta}{\gamma} \| \mathbf{D}^{k+1} - \mathbf{D}^k \|_{\mathbf{M}}^2 - \langle \mathbf{E}^k, \mathbf{D}^{k+1} - \mathbf{D}^k \rangle, \tag{23}$$

$$\langle \mathbf{X}^{k+1} - \mathbf{X}^*, \mathbf{D}^{k+1} - \mathbf{D}^* \rangle = \frac{\eta}{\gamma} \langle \mathbf{D}^{k+1} - \mathbf{D}^k, \mathbf{D}^{k+1} - \mathbf{D}^* \rangle_{\mathbf{M}} - \langle \mathbf{E}^k, \mathbf{D}^{k+1} - \mathbf{D}^* \rangle, \tag{24}$$

where  $\mathbf{M} = 2(\mathbf{I} - \mathbf{W})^{\dagger} - \gamma \mathbf{I}$  and  $\gamma < 2/\lambda_{max}(\mathbf{I} - \mathbf{W})$  ensures the positive definiteness of  $\mathbf{M}$  over span $\{\mathbf{I} - \mathbf{W}\}$ .

*Proof.* Since  $\mathbf{D}^{k+1} \in \mathbf{span}\{\mathbf{I} - \mathbf{W}\}\$ for any k, we have

$$\begin{split} &\langle \mathbf{X}^{k+1} - \mathbf{X}^*, \mathbf{D}^{k+1} - \mathbf{D}^k \rangle \\ = &\langle (\mathbf{I} - \mathbf{W})(\mathbf{X}^{k+1} - \mathbf{X}^*), (\mathbf{I} - \mathbf{W})^{\dagger}(\mathbf{D}^{k+1} - \mathbf{D}^k) \rangle \\ = &\left\langle \frac{\eta}{\gamma} (2\mathbf{I} - \gamma(\mathbf{I} - \mathbf{W}))(\mathbf{D}^{k+1} - \mathbf{D}^k) - (\mathbf{I} - \mathbf{W})\mathbf{E}^k, (\mathbf{I} - \mathbf{W})^{\dagger}(\mathbf{D}^{k+1} - \mathbf{D}^k) \right\rangle \qquad \text{(from (22))} \\ = &\left\langle \frac{\eta}{\gamma} (2(\mathbf{I} - \mathbf{W})^{\dagger} - \gamma \mathbf{I})(\mathbf{D}^{k+1} - \mathbf{D}^k) - \mathbf{E}^k, \mathbf{D}^{k+1} - \mathbf{D}^k \right\rangle \\ = &\frac{\eta}{\gamma} \|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2 - \langle \mathbf{E}^k, \mathbf{D}^{k+1} - \mathbf{D}^k \rangle. \end{split}$$

Similarly, we have

$$\begin{split} &\langle \mathbf{X}^{k+1} - \mathbf{X}^*, \mathbf{D}^{k+1} - \mathbf{D}^* \rangle \\ = &\langle (\mathbf{I} - \mathbf{W})(\mathbf{X}^{k+1} - \mathbf{X}^*), (\mathbf{I} - \mathbf{W})^\dagger (\mathbf{D}^{k+1} - \mathbf{D}^*) \rangle \\ = &\left\langle \frac{\eta}{\gamma} (2\mathbf{I} - \gamma (\mathbf{I} - \mathbf{W}))(\mathbf{D}^{k+1} - \mathbf{D}^k) - (\mathbf{I} - \mathbf{W})\mathbf{E}^k, (\mathbf{I} - \mathbf{W})^\dagger (\mathbf{D}^{k+1} - \mathbf{D}^*) \right\rangle \\ = &\left\langle \frac{\eta}{\gamma} (2(\mathbf{I} - \mathbf{W})^\dagger - \mathbf{I})(\mathbf{D}^{k+1} - \mathbf{D}^k) - \mathbf{E}^k, \mathbf{D}^{k+1} - \mathbf{D}^* \right\rangle \\ = &\frac{\eta}{\gamma} \langle \mathbf{D}^{k+1} - \mathbf{D}^k, \mathbf{D}^{k+1} - \mathbf{D}^* \rangle_{\mathbf{M}} - \langle \mathbf{E}^k, \mathbf{D}^{k+1} - \mathbf{D}^* \rangle. \end{split}$$

To make sure that M is positive definite over span $\{I - W\}$ , we need  $\gamma < 2/\lambda_{max}(I - W)$ .

Lemma 5. Taking the expectation conditioned on the compression in the kth iteration, we have

$$2\eta \mathbb{E} \langle \mathbf{E}^{k}, \mathbf{D}^{k+1} - \mathbf{D}^{*} \rangle = 2\eta \mathbb{E} \left\langle \mathbf{E}^{k}, \mathbf{D}^{k} + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) \mathbf{Y}^{k} + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) \mathbf{E}^{k} - \mathbf{D}^{*} \right\rangle$$

$$= \gamma \mathbb{E} \langle \mathbf{E}^{k}, (\mathbf{I} - \mathbf{W}) \mathbf{E}^{k} \rangle = \gamma \mathbb{E} ||\mathbf{E}^{k}||_{\mathbf{I} - \mathbf{W}}^{2},$$

$$2\eta \mathbb{E} \langle \mathbf{E}^{k}, \mathbf{D}^{k+1} - \mathbf{D}^{k} \rangle = 2\eta \mathbb{E} \left\langle \mathbf{E}^{k}, \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) \mathbf{Y}^{k} + \frac{\gamma}{2\eta} (\mathbf{I} - \mathbf{W}) \mathbf{E}^{k} \right\rangle$$

$$= \gamma \mathbb{E} \langle \mathbf{E}^{k}, (\mathbf{I} - \mathbf{W}) \mathbf{E}^{k} \rangle = \gamma \mathbb{E} ||\mathbf{E}^{k}||_{\mathbf{I} - \mathbf{W}}^{2}.$$

Proof. The proof is straightforward and omitted here.

Proof of Lemma 1. From Alg. 1, we have

$$2\eta\langle\mathbf{X}^{k}-\mathbf{X}^{*},\nabla\mathbf{F}(\mathbf{X}^{k};\boldsymbol{\xi}^{k})-\nabla\mathbf{F}(\mathbf{X}^{*})\rangle$$

$$=2\langle\mathbf{X}^{k}-\mathbf{X}^{*},\eta\nabla\mathbf{F}(\mathbf{X}^{k};\boldsymbol{\xi}^{k})-\eta\nabla\mathbf{F}(\mathbf{X}^{*})\rangle$$

$$=2\langle\mathbf{X}^{k}-\mathbf{X}^{*},\mathbf{X}^{k}-\mathbf{X}^{k+1}-\eta(\mathbf{D}^{k+1}-\mathbf{D}^{*})\rangle \quad (\text{from Line 7})$$

$$=2\langle\mathbf{X}^{k}-\mathbf{X}^{*},\mathbf{X}^{k}-\mathbf{X}^{k+1}\rangle-2\eta\langle\mathbf{X}^{k}-\mathbf{X}^{*},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle$$

$$=2\langle\mathbf{X}^{k}-\mathbf{X}^{*},\mathbf{X}^{k}-\mathbf{X}^{k+1}\rangle-2\eta\langle\mathbf{X}^{k}-\mathbf{X}^{k+1},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle-2\eta\langle\mathbf{X}^{k+1}-\mathbf{X}^{*},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle$$

$$=2\langle\mathbf{X}^{k}-\mathbf{X}^{*}-\eta(\mathbf{D}^{k+1}-\mathbf{D}^{*}),\mathbf{X}^{k}-\mathbf{X}^{k+1}\rangle-2\eta\langle\mathbf{X}^{k+1}-\mathbf{X}^{*},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle$$

$$=2\langle\mathbf{X}^{k+1}-\mathbf{X}^{*}+\eta(\nabla\mathbf{F}(\mathbf{X}^{k};\boldsymbol{\xi}^{k})-\nabla\mathbf{F}(\mathbf{X}^{*})),\mathbf{X}^{k}-\mathbf{X}^{k+1}\rangle-2\eta\langle\mathbf{X}^{k+1}-\mathbf{X}^{*},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle \quad (\text{from Line 7})$$

$$=2\langle\mathbf{X}^{k+1}-\mathbf{X}^{*},\mathbf{X}^{k}-\mathbf{X}^{k+1}\rangle+2\eta\langle\nabla\mathbf{F}(\mathbf{X}^{k};\boldsymbol{\xi}^{k})-\nabla\mathbf{F}(\mathbf{X}^{*}),\mathbf{X}^{k}-\mathbf{X}^{k+1}\rangle$$

$$-2\eta\langle\mathbf{X}^{k+1}-\mathbf{X}^{*},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle. \quad (25)$$

Then we consider the terms on the right hand side of (25) separately. Using  $2\langle \mathbf{A} - \mathbf{B}, \mathbf{B} - \mathbf{C} \rangle = \|\mathbf{A} - \mathbf{C}\|^2 - \|\mathbf{B} - \mathbf{C}\|^2 - \|\mathbf{A} - \mathbf{B}\|^2$ , we have

$$2\langle \mathbf{X}^{k+1} - \mathbf{X}^*, \mathbf{X}^k - \mathbf{X}^{k+1} \rangle = 2\langle \mathbf{X}^* - \mathbf{X}^{k+1}, \mathbf{X}^{k+1} - \mathbf{X}^k \rangle$$
$$= \|\mathbf{X}^k - \mathbf{X}^*\|^2 - \|\mathbf{X}^{k+1} - \mathbf{X}^k\|^2 - \|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2. \tag{26}$$

Using 
$$2\langle \mathbf{A}, \mathbf{B} \rangle = \|\mathbf{A}\|^2 + \|\mathbf{B}\|^2 - \|\mathbf{A} - \mathbf{B}\|^2$$
, we have

$$2\eta \langle \nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k}) - \nabla \mathbf{F}(\mathbf{X}^{*}), \mathbf{X}^{k} - \mathbf{X}^{k+1} \rangle$$

$$= \eta^{2} \|\nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k}) - \nabla \mathbf{F}(\mathbf{X}^{*})\|^{2} + \|\mathbf{X}^{k} - \mathbf{X}^{k+1}\|^{2} - \|\mathbf{X}^{k} - \mathbf{X}^{k+1} - \eta(\nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k}) - \nabla \mathbf{F}(\mathbf{X}^{*}))\|^{2}$$

$$= \eta^{2} \|\nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k}) - \nabla \mathbf{F}(\mathbf{X}^{*})\|^{2} + \|\mathbf{X}^{k} - \mathbf{X}^{k+1}\|^{2} - \eta^{2} \|\mathbf{D}^{k+1} - \mathbf{D}^{*}\|^{2}. \quad \text{(from Line 7)}$$
(27)

Combining (25), (26), (27), and (23), we obtain

$$\begin{split} &2\eta\langle\mathbf{X}^{k}-\mathbf{X}^{*},\nabla\mathbf{F}(\mathbf{X}^{k};\boldsymbol{\xi}^{k})-\nabla\mathbf{F}(\mathbf{X}^{*})\rangle\\ &=\underbrace{\|\mathbf{X}^{k}-\mathbf{X}^{*}\|^{2}-\|\mathbf{X}^{k+1}-\mathbf{X}^{k}\|^{2}-\|\mathbf{X}^{k+1}-\mathbf{X}^{*}\|^{2}}_{2\langle\mathbf{X}^{k+1}-\mathbf{X}^{*},\mathbf{X}^{k}-\mathbf{X}^{k+1}\rangle}\\ &+\underbrace{\eta^{2}\|\nabla\mathbf{F}(\mathbf{X}^{k};\boldsymbol{\xi}^{k})-\nabla\mathbf{F}(\mathbf{X}^{*})\|^{2}+\|\mathbf{X}^{k}-\mathbf{X}^{k+1}\|^{2}-\eta^{2}\|\mathbf{D}^{k+1}-\mathbf{D}^{*}\|^{2}}_{2\eta\langle\nabla\mathbf{F}(\mathbf{X}^{k};\boldsymbol{\xi}^{k})-\nabla\mathbf{F}(\mathbf{X}^{*}),\mathbf{X}^{k}-\mathbf{X}^{k+1}\rangle}\\ &-\underbrace{\left(\frac{2\eta^{2}}{\gamma}\langle\mathbf{D}^{k+1}-\mathbf{D}^{k},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle_{\mathbf{M}}-2\eta\langle\mathbf{E}^{k},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle\right)}_{2\eta\langle\mathbf{X}^{k+1}-\mathbf{X}^{*},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle}\\ &=\|\mathbf{X}^{k}-\mathbf{X}^{*}\|^{2}-\|\mathbf{X}^{k+1}-\mathbf{X}^{k}\|^{2}-\|\mathbf{X}^{k+1}-\mathbf{X}^{*}\|^{2}\\ &+\eta^{2}\|\nabla\mathbf{F}(\mathbf{X}^{k};\boldsymbol{\xi}^{k})-\nabla\mathbf{F}(\mathbf{X}^{*})\|^{2}+\|\mathbf{X}^{k}-\mathbf{X}^{k+1}\|^{2}-\eta^{2}\|\mathbf{D}^{k+1}-\mathbf{D}^{*}\|^{2}\\ &+\frac{\eta^{2}}{\gamma}\underbrace{\left(\|\mathbf{D}^{k}-\mathbf{D}^{*}\|_{\mathbf{M}}^{2}-\|\mathbf{D}^{k+1}-\mathbf{D}^{*}\|_{\mathbf{M}}^{2}-\|\mathbf{D}^{k+1}-\mathbf{D}^{k}\|_{\mathbf{M}}^{2}\right)}_{-2\langle\mathbf{D}^{k+1}-\mathbf{D}^{k},\mathbf{D}^{k+1}-\mathbf{D}^{*}\rangle_{\mathbf{M}}}\end{aligned}$$

where the last equality holds because

$$2\langle \mathbf{D}^k - \mathbf{D}^{k+1}, \mathbf{D}^{k+1} - \mathbf{D}^* \rangle_{\mathbf{M}} = \|\mathbf{D}^k - \mathbf{D}^*\|_{\mathbf{M}}^2 - \|\mathbf{D}^{k+1} - \mathbf{D}^*\|_{\mathbf{M}}^2 - \|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2$$

Thus, we reformulate it as

$$\begin{split} &\|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2 + \frac{\eta^2}{\gamma} \|\mathbf{D}^{k+1} - \mathbf{D}^*\|_{\mathbf{M}}^2 \\ = &\|\mathbf{X}^k - \mathbf{X}^*\|^2 + \frac{\eta^2}{\gamma} \|\mathbf{D}^k - \mathbf{D}^*\|_{\mathbf{M}}^2 - \frac{\eta^2}{\gamma} \|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2 - \eta^2 \|\mathbf{D}^{k+1} - \mathbf{D}^*\|^2 \\ &- 2\eta \langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle + \eta^2 \|\nabla \mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \nabla \mathbf{F}(\mathbf{X}^*)\|^2 + 2\eta \langle \mathbf{E}^k, \mathbf{D}^{k+1} - \mathbf{D}^* \rangle \end{split}$$

#### E.4 PROOF OF LEMMA 2

which completes the proof.

Proof of Lemma 2. From Alg. 1, we take the expectation conditioned on kth compression and obtain

$$\mathbb{E}\|\mathbf{H}^{k+1} - \mathbf{X}^*\|^2$$

$$= \mathbb{E}\|(1-\alpha)(\mathbf{H}^k - \mathbf{X}^*) + \alpha(\mathbf{Y}^k - \mathbf{X}^*) + \alpha\mathbf{E}^k\|^2 \qquad \text{(from Line 13)}$$

$$= \|(1-\alpha)(\mathbf{H}^k - \mathbf{X}^*) + \alpha(\mathbf{Y}^k - \mathbf{X}^*)\|^2 + \alpha^2 \mathbb{E}\|\mathbf{E}^k\|^2$$

$$= (1-\alpha)\|\mathbf{H}^k - \mathbf{X}^*\|^2 + \alpha\|\mathbf{Y}^k - \mathbf{X}^*\|^2 - \alpha(1-\alpha)\|\mathbf{H}^k - \mathbf{Y}^k\|^2 + \alpha^2 \mathbb{E}\|\mathbf{E}^k\|^2. \qquad (28)$$

In the second equality, we used the unbiasedness of the compression, i.e.,  $\mathbb{E}\mathbf{E}^k = \mathbf{0}$ . The last equality holds because of

$$\|(1-\alpha)\mathbf{A} + \alpha\mathbf{B}\|^2 = (1-\alpha)\|\mathbf{A}\|^2 + \alpha\|\mathbf{B}\|^2 - \alpha(1-\alpha)\|\mathbf{A} - \mathbf{B}\|^2$$

In addition, by taking the conditional expectation on the compression, we have

$$\|\mathbf{Y}^{k} - \mathbf{X}^{*}\|^{2} = \|\mathbf{X}^{k} - \eta \nabla \mathbf{F}(\mathbf{X}^{k}; \boldsymbol{\xi}^{k}) - \eta \mathbf{D}^{k} - \mathbf{X}^{*}\|^{2} \quad \text{(from Line 4)}$$

$$= \mathbb{E}\|\mathbf{X}^{k+1} + \eta \mathbf{D}^{k+1} - \eta \mathbf{D}^{k} - \mathbf{X}^{*}\|^{2} \quad \text{(from Line 7)}$$

$$= \mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^{*}\|^{2} + \eta^{2}\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^{k}\|^{2} + 2\eta\mathbb{E}\langle\mathbf{X}^{k+1} - \mathbf{X}^{*}, \mathbf{D}^{k+1} - \mathbf{D}^{k}\rangle$$

$$= \mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^{*}\|^{2} + \eta^{2}\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^{k}\|^{2}$$

$$+ \frac{2\eta^{2}}{\gamma}\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^{k}\|^{2}_{\mathbf{M}} - 2\eta\mathbb{E}\langle\mathbf{E}^{k}, \mathbf{D}^{k+1} - \mathbf{D}^{k}\rangle. \quad \text{(from (23))}$$

$$= \mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^{*}\|^{2} + \eta^{2}\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^{k}\|^{2}$$

$$+ \frac{2\eta^{2}}{\gamma}\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^{k}\|^{2}_{\mathbf{M}} - \gamma\mathbb{E}\|\mathbf{E}^{k}\|^{2}_{\mathbf{I}-\mathbf{W}}. \quad \text{(from Line 6)}$$

Combing the above two equations (28) and (29) together, we have

$$\mathbb{E}\|\mathbf{H}^{k+1} - \mathbf{X}^*\|^2$$

$$\leq (1 - \alpha)\|\mathbf{H}^k - \mathbf{X}^*\|^2 + \alpha \mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2 + \alpha \eta^2 \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|^2 + \frac{2\alpha\eta^2}{\gamma} \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2$$

$$- \alpha\gamma \mathbb{E}\|\mathbf{E}^k\|_{\mathbf{I}-\mathbf{W}}^2 + \alpha^2 \mathbb{E}\|\mathbf{E}^k\|^2 - \alpha(1 - \alpha)\|\mathbf{Y}^k - \mathbf{H}^k\|^2, \tag{30}$$

which completes the proof.

#### E.5 PROOF OF THEOREM 1

*Proof of Theorem 1.* Combining Lemmas 1, 2, and 5, we have the expectation conditioned on the compression satisfying

(30)

$$\mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2 + \frac{\eta^2}{\gamma} \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^*\|_{\mathbf{M}}^2 + a_1 \mathbb{E}\|\mathbf{H}^{k+1} - \mathbf{X}^*\|^2$$

$$\leq \|\mathbf{X}^k - \mathbf{X}^*\|^2 + \frac{\eta^2}{\gamma} \|\mathbf{D}^k - \mathbf{D}^*\|_{\mathbf{M}}^2 - \frac{\eta^2}{\gamma} \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2 - \eta^2 \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^*\|^2$$

$$- 2\eta \langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle + \eta^2 \|\nabla \mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \nabla \mathbf{F}(\mathbf{X}^*)\|^2 + \gamma \mathbb{E}\|\mathbf{E}^k\|_{\mathbf{I}-\mathbf{W}}^2$$

$$+ a_1 (1 - \alpha) \|\mathbf{H}^k - \mathbf{X}^*\|^2 + a_1 \alpha \mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2 + a_1 \alpha \eta^2 \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|^2$$

$$+ \frac{2a_1 \alpha \eta^2}{\gamma} \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2 + a_1 \alpha^2 \mathbb{E}\|\mathbf{E}^k\|^2 - a_1 \alpha \gamma \mathbb{E}\|\mathbf{E}^k\|_{\mathbf{I}-\mathbf{W}}^2 - a_1 \alpha (1 - \alpha)\|\mathbf{Y}^k - \mathbf{H}^k\|^2$$

$$= \underbrace{\|\mathbf{X}^k - \mathbf{X}^*\|^2 - 2\eta \langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle + \eta^2 \|\nabla \mathbf{F}(\mathbf{X}^k; \boldsymbol{\xi}^k) - \nabla \mathbf{F}(\mathbf{X}^*)\|^2}_{\mathcal{A}}$$

$$+ a_1 \alpha \mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2 + \frac{\eta^2}{\gamma}\|\mathbf{D}^k - \mathbf{D}^*\|_{\mathbf{M}}^2 - \eta^2 \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^*\|^2$$

$$+ a_1 (1 - \alpha)\|\mathbf{H}^k - \mathbf{X}^*\|^2 - (1 - 2a_1 \alpha)\frac{\eta^2}{\gamma} \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2 + a_1 \alpha \eta^2 \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|^2$$

$$+ \underbrace{a_1 \alpha^2 \mathbb{E}\|\mathbf{E}^k\|^2 + (1 - a_1 \alpha)\gamma \mathbb{E}\|\mathbf{E}^k\|_{\mathbf{I}-\mathbf{W}}^2 - a_1 \alpha (1 - \alpha)\|\mathbf{Y}^k - \mathbf{H}^k\|^2}_{\mathcal{B}}, \tag{31}$$

where  $a_1$  is a non-negative number to be determined. Then we deal with the three terms on the right hand side separately. We want the terms  $\mathcal{B}$  and  $\mathcal{C}$  to be nonpositive. First, we consider  $\mathcal{B}$ . Note that  $\mathbf{D}^k \in \mathbf{Range}(\mathbf{I} - \mathbf{W})$ . If we want  $\mathcal{B} \leq 0$ , then, we need  $1 - 2a_1\alpha > 0$ , i.e.,  $a_1\alpha < 1/2$ . Therefore we have

$$\mathcal{B} = -(1 - 2a_1\alpha)\frac{\eta^2}{\gamma}\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|_{\mathbf{M}}^2 + a_1\alpha\eta^2\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|^2$$
$$\leq \left(a_1\alpha - \frac{(1 - 2a_1\alpha)\lambda_{n-1}(\mathbf{M})}{\gamma}\right)\eta^2\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^k\|^2,$$

where  $\lambda_{n-1}(\mathbf{M}) > 0$  is the second smallest eigenvalue of  $\mathbf{M}$ . It means that we also need

$$a_1\alpha + \frac{(2a_1\alpha - 1)\lambda_{n-1}(\mathbf{M})}{\gamma} \le 0,$$

which is equivalent to

$$a_1 \alpha \le \frac{\lambda_{n-1}(\mathbf{M})}{\gamma + 2\lambda_{n-1}(\mathbf{M})} < 1/2. \tag{32}$$

Then we look at C. We have

$$C = a_1 \alpha^2 \mathbb{E} \|\mathbf{E}^k\|^2 + (1 - a_1 \alpha) \gamma \mathbb{E} \|\mathbf{E}^k\|_{\mathbf{I} - \mathbf{W}}^2 - a_1 \alpha (1 - \alpha) \|\mathbf{Y}^k - \mathbf{H}^k\|^2$$

$$\leq ((1 - a_1 \alpha) \beta \gamma + a_1 \alpha^2) \mathbb{E} \|\mathbf{E}^k\|^2 - a_1 \alpha (1 - \alpha) \|\mathbf{Y}^k - \mathbf{H}^k\|^2$$

$$\leq C((1 - a_1 \alpha) \beta \gamma + a_1 \alpha^2) \|\mathbf{Y}^k - \mathbf{H}^k\|^2 - a_1 \alpha (1 - \alpha) \|\mathbf{Y}^k - \mathbf{H}^k\|^2$$

Because we have  $1 - a_1 \alpha > 1/2$ , so we need

$$C((1 - a_1 \alpha)\beta \gamma + a_1 \alpha^2) - a_1 \alpha (1 - \alpha) = (1 + C)a_1 \alpha^2 - a_1 (C\beta \gamma + 1)\alpha + C\beta \gamma \le 0.$$
 (33)

That is

$$\alpha \ge \frac{a_1(C\beta\gamma + 1) - \sqrt{a_1^2(C\beta\gamma + 1)^2 - 4(1+C)Ca_1\beta\gamma}}{2(1+C)a_1} =: \alpha_0, \tag{34}$$

$$\alpha \le \frac{a_1(C\beta\gamma + 1) + \sqrt{a_1^2(C\beta\gamma + 1)^2 - 4(1+C)Ca_1\beta\gamma}}{2(1+C)a_1} =: \alpha_1.$$
 (35)

Next, we look at A. Firstly, by the bounded variance assumption, we have the expectation conditioned on the gradient sampling in kth iteration satisfying

$$\mathbb{E}\|\mathbf{X}^k - \mathbf{X}^*\|^2 - 2\eta \mathbb{E}\langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k; \xi^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle + \eta^2 \mathbb{E}\|\nabla \mathbf{F}(\mathbf{X}^k; \xi^k) - \nabla \mathbf{F}(\mathbf{X}^*)\|^2$$

$$\leq \|\mathbf{X}^k - \mathbf{X}^*\|^2 - 2\eta \langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle + \eta^2 \|\nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*)\|^2 + n\eta^2 \sigma^2$$

Then with the smoothness and strong convexity from Assumptions 4, we have the co-coercivity of  $\nabla g_i(\boldsymbol{x})$  with  $g_i(\boldsymbol{x}) := f_i(\boldsymbol{x}) - \frac{u}{2} \|\boldsymbol{x}\|_2^2$ , which gives

$$\langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle \ge \frac{\mu L}{\mu + L} \|\mathbf{X}^k - \mathbf{X}^*\|^2 + \frac{1}{\mu + L} \|\nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*)\|^2.$$

When  $\eta \leq 2/(\mu + L)$ , we have

$$\begin{split} &\langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle \\ &= \left(1 - \frac{\eta(\mu + L)}{2}\right) \langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle + \frac{\eta(\mu + L)}{2} \langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle \\ &\geq \left(\mu - \frac{\eta\mu(\mu + L)}{2} + \frac{\eta\mu L}{2}\right) \|\mathbf{X}^k - \mathbf{X}^*\|^2 + \frac{\eta}{2} \|\nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*)\|^2 \\ &= \mu \left(1 - \frac{\eta\mu}{2}\right) \|\mathbf{X}^k - \mathbf{X}^*\|^2 + \frac{\eta}{2} \|\nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*)\|^2. \end{split}$$

Therefore, we obtain

$$-2\eta \langle \mathbf{X}^k - \mathbf{X}^*, \nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*) \rangle$$
  

$$\leq -\eta^2 \|\nabla \mathbf{F}(\mathbf{X}^k) - \nabla \mathbf{F}(\mathbf{X}^*)\|^2 - \mu(2\eta - \mu\eta^2) \|\mathbf{X}^k - \mathbf{X}^*\|^2.$$
(36)

Conditioned on the kthe iteration, (i.e., conditioned on the gradient sampling in kth iteration), the inequality (31) becomes

$$\mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2 + \frac{\eta^2}{\gamma} \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^*\|_{\mathbf{M}}^2 + a_1 \mathbb{E}\|\mathbf{H}^{k+1} - \mathbf{X}^*\|^2$$

$$\leq \left(1 - \mu(2\eta - \mu\eta^2)\right) \|\mathbf{X}^k - \mathbf{X}^*\|^2 + a_1\alpha \mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^*\|^2$$

$$+ \frac{\eta^2}{\gamma} \|\mathbf{D}^k - \mathbf{D}^*\|_{\mathbf{M}}^2 - \eta^2 \mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^*\|^2 + a_1(1 - \alpha)\|\mathbf{H}^k - \mathbf{X}^*\|^2 + n\eta^2\sigma^2, \tag{37}$$

if the step size satisfies  $\eta \leq \frac{2}{\mu + L}$ . Rewriting (37), we have

$$(1 - a_{1}\alpha)\mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^{*}\|^{2} + \frac{\eta^{2}}{\gamma}\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^{*}\|_{\mathbf{M}}^{2} + \eta^{2}\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^{*}\|^{2} + a_{1}\mathbb{E}\|\mathbf{H}^{k+1} - \mathbf{X}^{*}\|^{2}$$

$$\leq \left(1 - \mu(2\eta - \mu\eta^{2})\right)\|\mathbf{X}^{k} - \mathbf{X}^{*}\|^{2} + \frac{\eta^{2}}{\gamma}\|\mathbf{D}^{k} - \mathbf{D}^{*}\|_{\mathbf{M}}^{2} + a_{1}(1 - \alpha)\|\mathbf{H}^{k} - \mathbf{X}^{*}\|^{2} + n\eta^{2}\sigma^{2},$$
(38)

and thus

$$(1 - a_{1}\alpha)\mathbb{E}\|\mathbf{X}^{k+1} - \mathbf{X}^{*}\|^{2} + \frac{\eta^{2}}{\gamma}\mathbb{E}\|\mathbf{D}^{k+1} - \mathbf{D}^{*}\|_{\mathbf{M}+\gamma\mathbf{I}}^{2} + a_{1}\mathbb{E}\|\mathbf{H}^{k+1} - \mathbf{X}^{*}\|^{2}$$

$$\leq \left(1 - \mu(2\eta - \mu\eta^{2})\right)\|\mathbf{X}^{k} - \mathbf{X}^{*}\|^{2} + \frac{\eta^{2}}{\gamma}\|\mathbf{D}^{k} - \mathbf{D}^{*}\|_{\mathbf{M}}^{2} + a_{1}(1 - \alpha)\|\mathbf{H}^{k} - \mathbf{X}^{*}\|^{2} + n\eta^{2}\sigma^{2}.$$
(39)

With the definition of  $\mathcal{L}^k$  in (12), we have

$$\mathbb{E}\mathcal{L}^{k+1} \le \rho \mathcal{L}^k + n\eta^2 \sigma^2,\tag{40}$$

with

$$\rho = \max \left\{ \frac{1 - \mu(2\eta - \mu\eta^2)}{1 - a_1\alpha}, \frac{\lambda_{\max}(\mathbf{M})}{\gamma + \lambda_{\max}(\mathbf{M})}, 1 - \alpha \right\}.$$

where

$$\lambda_{\max}(\mathbf{M}) = 2\lambda_{\max}((\mathbf{I} - \mathbf{W})^{\dagger}) - \gamma.$$

Recall all the conditions on the parameters  $a_1, \alpha$ , and  $\gamma$  to make sure that  $\rho < 1$ :

$$a_1 \alpha \le \frac{\lambda_{n-1}(\mathbf{M})}{\gamma + 2\lambda_{n-1}(\mathbf{M})},$$
 (41)

$$a_1 \alpha \le \mu (2\eta - \mu \eta^2), \tag{42}$$

$$\alpha \ge \frac{a_1(C\beta\gamma + 1) - \sqrt{a_1^2(C\beta\gamma + 1)^2 - 4(1+C)Ca_1\beta\gamma}}{2(1+C)a_1} =: \alpha_0, \tag{43}$$

$$\alpha \le \frac{a_1(C\beta\gamma + 1) + \sqrt{a_1^2(C\beta\gamma + 1)^2 - 4(1 + C)Ca_1\beta\gamma}}{2(1 + C)a_1} =: \alpha_1.$$
 (44)

In the following, we show that there exist parameters that satisfy these conditions.

Since we can choose any  $a_1$ , we let

$$a_1 = \frac{4(1+C)}{C\beta\gamma + 2},$$

such that

$$a_1^2(C\beta\gamma + 1)^2 - 4(1+C)Ca_1\beta\gamma = a_1^2$$

Then we have

$$\alpha_0 = \frac{C\beta\gamma}{2(1+C)} \to 0, \quad \text{as } \gamma \to 0,$$

$$\alpha_1 = \frac{C\beta\gamma + 2}{2(1+C)} \to \frac{1}{1+C}, \quad \text{as } \gamma \to 0.$$

Conditions (43) and (44) show

$$a_1 \alpha \in \left[ \frac{2C\beta\gamma}{C\beta\gamma + 2}, 2 \right] \to [0, 2], \text{ if } C = 0 \text{ or } \gamma \to 0.$$

Hence in order to make (41) and (42) satisfied, it's sufficient to make

$$\frac{2C\beta\gamma}{C\beta\gamma+2} \le \min\left\{\frac{\lambda_{n-1}(\mathbf{M})}{\gamma+2\lambda_{n-1}(\mathbf{M})}, \mu(2\eta-\mu\eta^2)\right\} = \min\left\{\frac{\frac{2}{\beta}-\gamma}{\frac{4}{\beta}-\gamma}, \mu(2\eta-\mu\eta^2)\right\}. \tag{45}$$

where we use  $\lambda_{n-1}(\mathbf{M}) = \frac{2}{\lambda_{\max}(\mathbf{I} - \mathbf{W})} - \gamma = \frac{2}{\beta} - \gamma$ .

When C > 0, the condition (45) is equivalent to

$$\gamma \le \min \left\{ \frac{(3C+1) - \sqrt{(3C+1)^2 - 4C}}{C\beta}, \frac{2\mu\eta(2-\mu\eta)}{[2-\mu\eta(2-\mu\eta)]C\beta} \right\}. \tag{46}$$

The first term can be simplified using

$$\frac{(3C+1) - \sqrt{(3C+1)^2 - 4C}}{C\beta} \ge \frac{2}{(3C+1)\beta}$$

due to  $\sqrt{1-x} \le 1 - \frac{x}{2}$  when  $x \in (0,1)$ .

Therefore, for a given stepsize  $\eta$ , if we choose

$$\gamma \in \left(0, \min\left\{\frac{2}{(3C+1)\beta}, \frac{2\mu\eta(2-\mu\eta)}{[2-\mu\eta(2-\mu\eta)]C\beta}\right\}\right)$$

and

$$\alpha \in \left[\frac{C\beta\gamma}{2(1+C)}, \min\left\{\frac{C\beta\gamma+2}{2(1+C)}, \frac{2-\beta\gamma}{4-\beta\gamma}\frac{C\beta\gamma+2}{4(1+C)}, \mu\eta(2-\mu\eta)\frac{C\beta\gamma+2}{4(1+C)}\right\}\right],$$

then, all conditions (41)-(44) hold.

Note that  $\gamma < \frac{2}{(3C+1)\beta}$  implies  $\gamma < \frac{2}{\beta}$ , which ensures the positive definiteness of  $\mathbf{M}$  over  $\mathbf{span}\{\mathbf{I} - \mathbf{W}\}$  in Lemma 4.

Note that  $\eta \leq \frac{2}{\mu + L}$  ensures

$$\mu \eta (2 - \mu \eta) \frac{C\beta \gamma + 2}{4(1+C)} \le \frac{C\beta \gamma + 2}{2(1+C)}.$$
 (47)

So, we can simplify the bound for  $\alpha$  as

$$\alpha \in \left[ \frac{C\beta\gamma}{2(1+C)}, \min\left\{ \frac{2-\beta\gamma}{4-\beta\gamma} \frac{C\beta\gamma+2}{4(1+C)}, \mu\eta(2-\mu\eta) \frac{C\beta\gamma+2}{4(1+C)} \right\} \right].$$

Lastly, taking the total expectation on both sides of (40) and using tower property, we complete the proof for  ${\cal C}>0$ .

Proof of Corollary 1. Let's first define  $\kappa_f = \frac{L}{\mu}$  and  $\kappa_g = \frac{\lambda_{\max}(\mathbf{I} - \mathbf{W})}{\lambda_{\min}^+(\mathbf{I} - \mathbf{W})} = \lambda_{\max}(\mathbf{I} - \mathbf{W})\lambda_{\max}((\mathbf{I} - \mathbf{W})^{\dagger}).$ 

We can choose the stepsize  $\eta=\frac{1}{L}$  such that the upper bound of  $\gamma$  is

$$\gamma_{\text{upper}} = \min\left\{\frac{2}{(3C+1)\beta}, \frac{\frac{2}{\kappa_f}\left(2 - \frac{1}{\kappa_f}\right)}{\left[2 - \frac{1}{\kappa_f}\left(2 - \frac{1}{\kappa_f}\right)\right]C\beta}, \frac{2}{\beta}\right\} \ge \min\left\{\frac{2}{(3C+1)\beta}, \frac{1}{\kappa_f C\beta}\right\},$$

due to  $\frac{x(2-x)}{2-x(2-x)} \ge \frac{x}{2-x} \ge x$  when  $x \in (0,1)$ .

Hence we can take  $\gamma = \min\{\frac{1}{(3C+1)\beta}, \frac{1}{\kappa_f C \beta}\}$ .

The bound of  $\alpha$  is

$$\alpha \in \left[ \frac{C\beta\gamma}{2(1+C)}, \min\left\{ \frac{2-\beta\gamma}{4-\beta\gamma} \frac{C\beta\gamma+2}{4(1+C)}, \frac{1}{\kappa_f} (2-\frac{1}{\kappa_f}) \frac{C\beta\gamma+2}{4(1+C)} \right\} \right]$$

When  $\gamma$  is chosen as  $\frac{1}{\kappa_f C \beta}$ , pick

$$\alpha = \frac{C\beta\gamma}{2(1+C)} = \frac{1}{2(1+C)\kappa_f}.$$
(48)

When  $\frac{1}{(3C+1)\beta} \leq \frac{1}{\kappa_f C\beta}$ , the upper bound of  $\alpha$  is

$$\begin{split} \alpha_{\text{upper}} &= \min \left\{ \frac{2 - \beta \gamma}{4 - \beta \gamma} \frac{C \beta \gamma + 2}{4(1 + C)}, \frac{1}{\kappa_f} (2 - \frac{1}{\kappa_f}) \frac{C \beta \gamma + 2}{4(1 + C)} \right\} \\ &= \min \left\{ \frac{6C + 1}{12C + 3}, \frac{1}{\kappa_f} (2 - \frac{1}{\kappa_f}) \right\} \frac{7C + 2}{4(C + 1)(3C + 1)} \\ &\geq \min \left\{ \frac{6C + 1}{12C + 3}, \frac{1}{\kappa_f} \right\} \frac{7C + 2}{4(C + 1)(3C + 1)}. \end{split}$$

In this case, we pick

$$\alpha = \min\left\{\frac{6C+1}{12C+3}, \frac{1}{\kappa_f}\right\} \frac{7C+2}{4(C+1)(3C+1)}.$$
 (49)

Note  $\alpha = \mathcal{O}\left(\frac{1}{(1+C)\kappa_f}\right)$  since  $\frac{6C+1}{12C+3}$  is lower bounded by  $\frac{1}{3}$ . Hence in both cases (Eq. (48) and Eq. (49)),  $\alpha = \mathcal{O}\left(\frac{1}{(1+C)\kappa_f}\right)$ , and the third term of  $\rho$  is upper bounded by

$$1 - \alpha \le \max\left\{1 - \frac{1}{2(1+C)\kappa_f}, 1 - \min\left\{\frac{6C+1}{12C+3}, \frac{1}{\kappa_f}\right\} \frac{7C+2}{4(1+C)(3C+1)}\right\}$$

In two cases of  $\gamma$ , the second term of  $\rho$  becomes

$$1 - \frac{\gamma}{2\lambda_{\max}((\mathbf{I} - \mathbf{W})^{\dagger})} = \max\left\{1 - \frac{1}{2C\kappa_f\kappa_g}, 1 - \frac{1}{(1 + 3C)\kappa_g}\right\}$$

Before analysing the first term of  $\rho$ , we look at  $a_1\alpha$  in two cases of  $\gamma$ . When  $\gamma = \frac{1}{\kappa_f C\beta}$ ,

$$a_1 \alpha = \frac{2C\beta\gamma}{C\beta\gamma + 2} = \frac{2}{2\kappa_f + 1} \le \frac{1}{\kappa_f}.$$

When  $\gamma = \frac{1}{(3C+1)\beta}$ ,

$$a_1 \alpha = \min \left\{ \frac{6C+1}{(12C+3)}, \frac{1}{\kappa_f} \right\} \le \frac{1}{\kappa_f}.$$

In both cases,  $a_1 \alpha \leq \frac{1}{\kappa_f}$ . Therefore, the first term of  $\rho$  becomes

$$\frac{1 - \mu \eta(2 - \mu \eta)}{1 - a_1 \alpha} \le \frac{1 - \frac{1}{\kappa_f} (2 - \frac{1}{\kappa_f})}{1 - \frac{1}{\kappa_f}} = 1 - \frac{1 - \frac{1}{\kappa_f}}{\kappa_f - 1} = 1 - \frac{1}{\kappa_f}.$$

To summarize, we have

$$\rho \leq 1 - \min\left\{\frac{1}{\kappa_f}, \frac{1}{2C\kappa_f\kappa_g}, \frac{1}{(1+3C)\kappa_g}, \frac{1}{2(1+C)\kappa_f}, \min\left\{\frac{6C+1}{12C+3}, \frac{1}{\kappa_f}\right\} \frac{7C+2}{4(1+C)(3C+1)}\right\}$$

and therefore

$$\rho = \max\left\{1 - \mathcal{O}\left(\frac{1}{(1+C)\kappa_f}\right), 1 - \mathcal{O}\left(\frac{1}{(1+C)\kappa_g}\right), 1 - \mathcal{O}\left(\frac{1}{C\kappa_f\kappa_g}\right)\right\}.$$

With full-gradient (i.e.,  $\sigma = 0$ ), we get  $\epsilon$ -accuracy solution with the total number of iterations

$$k \ge \widetilde{\mathcal{O}}((1+C)(\kappa_f + \kappa_q) + C\kappa_f \kappa_q).$$

When C=0, i.e., there is no compression, the iteration complexity recovers that of NIDS,  $\widetilde{\mathcal{O}}\left(\kappa_f+\kappa_g\right)$ .

When  $C \leq \frac{\kappa_f + \kappa_g}{\kappa_f \kappa_g + \kappa_f + \kappa_g}$ , the complexity is improved to that of NIDS, i.e., the compression doesn't harm the convergence in terms of the order of the coefficients.

Proof of Corollary 2. Note that  $(\bar{x}^k)^{\top} = \bar{X}^k$  and  $\mathbf{1}_{n \times 1} \bar{X}^* = X^*$ , then

$$\sum_{i=1}^{n} \mathbb{E} \| \mathbf{x}_{i}^{k} - \overline{\mathbf{x}}^{k} \|^{2} = \mathbb{E} \| \mathbf{X}^{k} - \mathbf{1}_{n \times 1} \overline{\mathbf{X}}^{k} \|^{2}$$

$$= \mathbb{E} \| \mathbf{X}^{k} - \mathbf{X}^{*} + \mathbf{X}^{*} - \mathbf{1}_{n \times 1} \overline{\mathbf{X}}^{k} \|^{2}$$

$$= \mathbb{E} \| \mathbf{X}^{k} - \mathbf{X}^{*} - \frac{\mathbf{1}_{n \times 1} \mathbf{1}_{n \times 1}^{\top}}{n} (\mathbf{X}^{k} - \mathbf{X}^{*}) \|$$

$$\leq \mathbb{E} \| \mathbf{X}^{k} - \mathbf{X}^{*} \|^{2}$$

$$\leq \frac{\rho \mathbb{E} \mathcal{L}^{k-1} + n\eta^{2} \sigma^{2} (1 - \rho)^{-1}}{1 - a_{1} \alpha}$$

$$\leq 2\rho^{k} \mathcal{L}^{0} + 2 \frac{n\eta^{2} \sigma^{2}}{1 - \rho}.$$
(50)

The last inequality holds because we have  $a_1 \alpha \leq 1/2$ .

*Proof of Corollary 3.* From the proof of Theorem 1, when C=0, we can set  $\gamma=1$ ,  $\alpha=1$ , and  $a_1=0$ . Plug those values into  $\rho$ , and we obtain the convergence rate for NIDS.

#### E.6 Proof of Theorem 2

Proof of Theorem 2. In order to get exact convergence, we pick diminishing step-size, set  $\alpha = \frac{C\beta\gamma}{2(1+C)}$ ,  $a_1\alpha = \frac{2C\beta\gamma_k}{C\beta\gamma_k+2}$ ,  $\theta_1 = \frac{1}{2\lambda_{\max}((\mathbf{I}-\mathbf{W})^\dagger)}$  and  $\theta_2 = \frac{C\beta}{2(1+C)}$ , then

$$\rho_k = \max \left\{ 1 - \frac{\mu \eta_k (2 - \mu \eta_k) - a_1 \alpha}{1 - a_1 \alpha}, 1 - \theta_1 \gamma_k, 1 - \theta_2 \gamma_k \right\}$$

If we further pick diminishing  $\eta_k$  and  $\gamma_k$  such that  $\mu \eta_k (2 - \mu \eta_k) - a_1 \alpha \ge a_1 \alpha$ , then

$$\frac{\mu\eta_k(2-\mu\eta_k)-a_1\alpha}{1-a_1\alpha} \ge \frac{a_1\alpha}{1-a_1\alpha} = \frac{2C\beta\gamma_k}{2-C\beta\gamma_k} \ge C\beta\gamma_k.$$

Notice that  $C\beta\gamma\leq \frac{2}{3}$  since  $(3C+1)-\sqrt{(3C+1)^2-4C}$  is increasing in C>0 with limit  $\frac{2}{3}$  at  $\infty$ .

In this case we only need,

$$\gamma_k \in \left(0, \min\left\{\frac{(3C+1) - \sqrt{(3C+1)^2 - 4C}}{C\beta}, \frac{2\mu\eta_k(2 - \mu\eta_k)}{[4 - \mu\eta_k(2 - \mu\eta_k)]C\beta}, \frac{2}{\beta}\right\}\right). \tag{51}$$

And

$$\rho_k \le \max\left\{1 - C\beta\gamma_k, 1 - \theta_1\gamma_k, 1 - \theta_2\gamma_k\right\} \le 1 - \theta_3\gamma_k$$

if  $\theta_3 = \min\{\theta_1, \theta_2\}$  and note that  $\theta_2 \leq C\beta$ .

We define

$$\mathcal{L}^k := (1 - a_1 \alpha_k) \|\mathbf{X}^k - \mathbf{X}^*\|^2 + (2\eta_k^2/\gamma_k) \mathbb{E} \|\mathbf{D}^{k+1} - \mathbf{D}^*\|_{(\mathbf{I} - \mathbf{W})^{\dagger}}^2 + a_1 \|\mathbf{H}^k - \mathbf{X}^*\|^2.$$

Hence

$$\mathbb{E}\mathcal{L}^{k+1} \le (1 - \theta_3 \gamma_k) \mathbb{E}\mathcal{L}^k + n\sigma^2 \eta_k^2$$

From  $a_1 \alpha \leq \frac{\mu \eta_k (2 - \mu \eta_k)}{2}$ , we get

$$\frac{4C\beta\gamma_k}{C\beta\gamma_k+2} \le \mu\eta_k(2-\mu\eta_k).$$

If we pick  $\gamma_k = \theta_4 \eta_k$ , then it's sufficient to let

$$2C\beta\theta_4\eta_k \leq \mu\eta_k(2-\mu\eta_k).$$

Hence if  $\theta_4 < \frac{\mu}{C\beta}$  and let  $\eta_* = \frac{2(\mu - C\beta\theta_4)}{\mu^2}$ , then  $\eta_k = \frac{\gamma_k}{\theta_4} \in (0, \eta_*)$  guarantees the above discussion and

$$\mathbb{E}\mathcal{L}^{k+1} \le (1 - \theta_3 \theta_4 \eta_k) \mathbb{E}\mathcal{L}^k + n\sigma^2 \eta_k^2.$$

So far all restrictions for  $\eta_k$  are

$$\eta_k \le \min\left\{\frac{2}{\mu + L}, \eta_*\right\}$$

and

$$\eta_k \le \frac{1}{\theta_4} \min \left\{ \frac{(3C+1) - \sqrt{(3C+1)^2 - 4C}}{C\beta}, \frac{2}{\beta} \right\}$$

Let  $\theta_5 = \min\left\{\frac{2}{\mu+L}, \eta_*, \frac{(3C+1)-\sqrt{(3C+1)^2-4C}}{C\beta\theta_4}, \frac{2}{\beta\theta_4}\right\}, \eta_k = \frac{1}{Bk+A} \text{ and } D = \max\left\{A\mathcal{L}^0, \frac{2n\sigma^2}{\theta_3\theta_4}\right\},$ 

we claim that if we pick  $B = \frac{\theta_3 \theta_4}{2}$  and some A, by setting  $\eta_k = \frac{2}{\theta_3 \theta_4 k + 2A}$ , we get

$$\mathbb{E}\mathcal{L}^k \le \frac{D}{Bk+A}.$$

Induction:

When k = 0, it's obvious. Suppose previous k inequalities hold. Then

$$\mathbb{E}\mathcal{L}^{k+1} \le \left(1 - \frac{2\theta_3\theta_4}{\theta_3\theta_4k + 2A}\right) \frac{2D}{\theta_3\theta_4k + 2A} + \frac{4n\sigma^2}{(\theta_3\theta_4k + 2A)^2}.$$

Multiply  $M := (\theta_3\theta_4k + \theta_3\theta_4 + 2A)(\theta_3\theta_4k + 2A)(2D)^{-1}$  on both sides, we get

$$\begin{split} M\mathbb{E}\mathcal{L}^{k+1} &\leq \left(1 - \frac{2\theta_3\theta_4}{\theta_3\theta_4k + 2A}\right) (\theta_3\theta_4k + \theta_3\theta_4 + 2A) + \frac{4n\sigma^2(\theta_3\theta_4k + \theta_3\theta_4 + 2A)}{2D(\theta_3\theta_4k + 2A)} \\ &= \frac{2D(\theta_3\theta_4k + 2A - 2\theta_3\theta_4)(\theta_3\theta_4k + \theta_3\theta_4 + 2A) + 4n\sigma^2(\theta_3\theta_4k + \theta_3\theta_4 + 2A)}{2D(\theta_3\theta_4k + 2A)} \\ &= \frac{2D(\theta_3\theta_4k + 2A)^2 + 4n\sigma^2(\theta_3\theta_4k + 2A) - 4D\theta_3\theta_4(\theta_3\theta_4k + 2A) + 2D\theta_3\theta_4(\theta_3\theta_4k + 2A)}{2D(\theta_3\theta_4k + 2A)} \\ &+ \frac{-4D(\theta_3\theta_4)^2 + 4n\sigma^2\theta_3\theta_4}{2D(\theta_3\theta_4k + 2A)} \\ &\leq \theta_3\theta_4k + 2A. \end{split}$$

Hence

$$\mathbb{E}\mathcal{L}^{k+1} \le \frac{2D}{\theta_3\theta_4(k+1) + 2A}$$

This induction holds for any A such that  $\eta_k$  is feasible, i.e.

$$\eta_0 = \frac{1}{A} \le \theta_5.$$

Here we summarize the definition of constant numbers:

$$\theta_1 = \frac{1}{2\lambda_{\max}((\mathbf{I} - \mathbf{W})^{\dagger})}, \ \theta_2 = \frac{C\beta}{2(1+C)},$$
 (52)

$$\theta_3 = \min\{\theta_1, \theta_2\}, \ \theta_4 \in \left(0, \frac{\mu}{C\beta}\right), \ \eta_* = \frac{2(\mu - C\beta\theta_4)}{\mu^2}, \tag{53}$$

$$\theta_5 = \min\left\{\frac{2}{\mu + L}, \eta_*, \frac{(3C+1) - \sqrt{(3C+1)^2 - 4C}}{C\beta\theta_4}, \frac{2}{\beta\theta_4}\right\}.$$
 (54)

Therefore, let  $A=\frac{1}{\theta_5}$  and  $\eta_k=\frac{2\theta_5}{\theta_3\theta_4\theta_5k+2},$  we get

$$\frac{1}{n} \mathbb{E} \mathcal{L}^k \le \frac{2 \max\left\{\frac{1}{n} \mathcal{L}^0, \frac{2\sigma^2 \theta_5}{\theta_3 \theta_4}\right\}}{\theta_3 \theta_4 \theta_5 k + 2}.$$

Since  $1 - a_1 \alpha_k \ge 1/2$ , we complete the proof.