
Decentralized Strongly Convex Optimization with Less Local First-Order Oracle Complexity

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we study the decentralized optimization problem of minimizing the
2 strongly convex objective that is the sum of smooth convex functions stored across
3 a network with m local agents. We propose an efficient algorithm for finding an ϵ -
4 suboptimal solution within at most $\mathcal{O}((m + \sqrt{m\kappa}) \log(1/\epsilon))$ local first-order oracle
5 calls and $\tilde{\mathcal{O}}(\sqrt{\kappa/\alpha} \log(1/\epsilon))$ communication rounds, where κ is the condition
6 number of the objective and α is the spectral gap of the gossip matrix. Both of our
7 local first-order oracle complexity and communication complexity nearly match
8 the corresponding lower bounds. The proposed algorithm allows only few of the
9 agents compute their local gradients during one iteration, which significantly re-
10 duces the total computational cost. In contrast, the existing decentralized convex
11 optimization algorithms require all of the agents compute their local gradients
12 during every iteration, which leads to at least $\Omega(m\sqrt{\kappa} \log(1/\epsilon))$ local first-order
13 oracle complexity totally.

14 1 Introduction

15 In this paper, we focus on solving the decentralized optimization problem

$$\min_{x \in \mathbb{R}^d} f(x) \triangleq \frac{1}{m} \sum_{i=1}^m f_i(x) \quad (1.1)$$

16 on an undirected connected network with m agents, where the objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$
17 is μ -strongly-convex, and the local function $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ on the i -th agent is L -smooth and convex.
18 Decentralized algorithms desire all of the m agents to solve the optimization problem cooperatively
19 and each of the agents is only allowed to communicate with its neighbors.

20 First-order algorithms for decentralized convex optimization have been extensively studied in re-
21 cent years [11, 13, 15, 16, 19–21, 25, 27, 29, 31, 32, 37–39]. Scaman et al. [29] showed that
22 achieving an ϵ -suboptimal solution of problem (1.1) requires at least $\Omega(\sqrt{\kappa} \log(1/\epsilon))$ gradient steps
23 and $\Omega(\sqrt{\kappa/\alpha} \log(1/\epsilon))$ communication steps, where κ is the condition number of the objective
24 and α is the spectral gap of the gossip matrix. They attempted to match these lower bounds by
25 proposing multi-step dual accelerated (MSDA) method. However, the iteration of MSDA relies on
26 accessing the dual gradients of local functions, which may be intractable. We are more interested in
27 dual-free methods [7, 11, 12, 14, 32, 37] that only require the local gradient calls during the itera-
28 tions. In particular, Kovalev et al. [11] applied the idea of primal dual framework [3, 4, 6, 18] and
29 Chebyshev acceleration [17, 28] to design optimal proximal alternating predictor-corrector (OPAPC)

Table 1: We summarize local first-order oracle complexity and communication complexity of proposed KNOT and previous work. We use $\tilde{\mathcal{O}}(\cdot)$ and $\tilde{\Omega}(\cdot)$ to hide logarithmic factors of κ and m .

Methods	# Local First-Order Oracle	# Communication
APM-C [15]	$\mathcal{O}(m\sqrt{\kappa}\log(1/\varepsilon))$	$\mathcal{O}(\sqrt{\kappa/\alpha}\log^2(1/\varepsilon))$
OPAPC [11]	$\mathcal{O}(m\sqrt{\kappa}\log(1/\varepsilon))$	$\mathcal{O}(\sqrt{\kappa/\alpha}\log(1/\varepsilon))$
Acc-GT+CA [13]	$\mathcal{O}(m\sqrt{\kappa}\log(1/\varepsilon))$	$\mathcal{O}(\sqrt{\kappa/\alpha}\log(1/\varepsilon))$
KNOT (Theorem 1)	$\mathcal{O}((m + \sqrt{m\kappa})\log(1/\varepsilon))$	$\tilde{\mathcal{O}}(\sqrt{\kappa/\alpha}\log(1/\varepsilon))$
Lower Bounds [29, 34]	$\tilde{\Omega}(m + \sqrt{m\kappa}\log(1/\varepsilon))$	$\Omega(\sqrt{\kappa/\alpha}\log(1/\varepsilon))$

30 methods, which avoid dual gradient computation and the later one matches both of the lower bounds
31 for gradient steps and communication steps provided by Scaman et al. [29]. Li and Lin [13] es-
32 tablished the algorithm by incorporating gradient tracking [16, 21, 25, 31, 35, 36] and Chebyshev
33 acceleration [17, 28] into Nesterov’s acceleration [22] (Acc-GT+CA), which achieves the same com-
34 putation and communication complexities.

35 We notice that the existing statements on optimality of gradient steps for above first-order decen-
36 tralized algorithms can be refined [7, 11, 14, 37]. Concretely, the existing “optimal” first-order
37 algorithms for decentralized convex optimization require all of the agents computing their local
38 gradient during every iteration. Hence, each “gradient step” of these algorithm contains m local
39 gradient calls, resulting at least $\Omega(m\sqrt{\kappa}\log(1/\varepsilon))$ local gradient complexity in total.¹ In fact, the
40 atomic operation of the first-order decentralized algorithm is computing the gradient of one local
41 function, which implies allowing only few of the agents computing their local gradients during one
42 iteration potentially makes the algorithm be more computation efficient. In practice, decentralized
43 optimization are usually applied to networks with limited computational resources (e.g. mobile de-
44 vices [33], wireless sensors [26] and smart home appliances [9]), which also encourages us to design
45 decentralized algorithms with less local computational cost to reduce the energy consumption.

46 We consider the problem of minimizing objective function in problem (1.1) on a single machine. It
47 is well known that accelerated stochastic gradient methods [1, 10, 24] can achieve an ε -suboptimal
48 solution of such finite-sum problem within $\mathcal{O}((m + \sqrt{m\kappa})\log(1/\varepsilon))$ individual component gradient
49 calls. Since the objective of decentralized optimization problem also has the finite-sum structure, it
50 implies the known local gradient oracle complexity $\mathcal{O}(m\sqrt{\kappa}\log(1/\varepsilon))$ of existing first-order algo-
51 rithms [7, 11, 14, 32, 37] maybe not optimal. This naturally leads to the following question

52 *Can we design a decentralized first-order algorithm with less local gradient calls?*

53 In this paper, we give a positive answer to above question by proposing Katyusha-type Near-Optimal
54 decenTralized algorithm (KNOT). Our method allows the agents to skip the step of computing lo-
55 cal gradient during most of the iterations, which significantly improves the total computational ef-
56 ficiency. We prove that KNOT can achieve an ε -suboptimal solution of problem (1.1) within at
57 most $\mathcal{O}((m + \sqrt{m\kappa})\log(1/\varepsilon))$ local gradient oracle complexity and $\tilde{\mathcal{O}}(\sqrt{\kappa/\alpha}\log(1/\varepsilon))$ commu-
58 nication complexity, which nearly matches the corresponding lower bounds [29, 34]. We compare
59 the theoretical results of proposed KNOT and previous methods in Table 1.

60 **Paper Organization** In section 2, we introduce the notations and settings throughout this paper.
61 In section 3, we propose a new decentralized optimization algorithm and provide its convergence
62 analysis. In section 4, we give a discussion for the optimality of proposed algorithm. In section 5,
63 we provide numerical experiments to validate our theory. We conclude our work in section 6. All
64 proofs are deferred to appendix.

¹Recall that finding an ε -suboptimal solution of smooth and strongly convex function on single machine requires at least $\Omega(\sqrt{\kappa}\log(1/\varepsilon))$ gradient calls in terms of the objective function [22].

Algorithm 1 AccGossip (\mathbf{v}_0, K)

- 1: $\mathbf{v}^{-1} = \mathbf{v}^0$
 - 2: $\beta = \frac{1 - \sqrt{1 - \lambda_2^2(W)}}{1 + \sqrt{1 - \lambda_2^2(W)}}$
 - 3: **for** $k = 0, \dots, K$
 - 4: $\mathbf{v}^{k+1} = (1 + \beta)W\mathbf{v}^k - \beta\mathbf{v}^{k-1}$
 - 5: **end for**
 - 6: **Output:** \mathbf{v}^K
-

65 2 Preliminaries

66 We use $\|\cdot\|$ to present the Frobenius norm of the matrix and the Euclidean norm of the vector. We
67 introduce the aggregated notations

$$\mathbf{x} = [x_1, \dots, x_m]^\top, \quad \nabla F(\mathbf{x}) = [\nabla f_1(x_1), \dots, \nabla f_m(x_m)]^\top \quad \text{and} \quad \bar{x} = \frac{1}{m} \mathbf{1}^\top \mathbf{x},$$

68 where x_i is the local variable on the i -th agent and $\nabla f_i(x_i)$ is the corresponding local gradient.

69 We impose the following assumptions on the decentralized optimization problem (1.1).

70 **Assumption 1.** We assume each local function $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth, i.e., there exists constant
71 $L > 0$ such that

$$f_i(y) - f_i(x) \leq \langle \nabla f_i(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$$

72 for any $x, y \in \mathbb{R}^d$.

73 **Assumption 2.** We assume each local function $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex, i.e., we have

$$f_i(y) - f_i(x) \geq \langle \nabla f_i(x), y - x \rangle$$

74 for any $x, y \in \mathbb{R}^d$.

75 **Assumption 3.** We assume the global function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex, i.e., there exists
76 constant $\mu > 0$ such that

$$f(y) - f(x) \geq \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2$$

77 for any $x, y \in \mathbb{R}^d$.

78 Assumption 1 implies $f(\cdot)$ is L -smooth and we define $\kappa \triangleq L/\mu$ as its condition number.

79 The strong convexity means the objective $f(\cdot)$ has unique minimizer x^* . We say $\hat{x} = [\hat{x}_1, \dots, \hat{x}_m]^\top$
80 is an ϵ -suboptimal solution of the decentralized optimization problem (1.1) if $f(\hat{x}_i) - f(x^*) \leq \epsilon$
81 holds for any $i = 1, \dots, m$.

82 We let $W \in \mathbb{R}^{m \times m}$ be the gossip matrix associated with the network of m agents and it satisfies the
83 following assumption.

84 **Assumption 4.** We assume the gossip matrix W is symmetric and $W_{i,j} \neq 0$ if and only if the i -th
85 and the j -th agents are connected in the network. We also assume W satisfies $\mathbf{0} \preceq W \preceq I, W\mathbf{1} = \mathbf{1}$
86 and $\text{null}(I - W) = \text{span}(\mathbf{1})$.

87 We define the spectral gap of the gossip matrix as $\alpha \triangleq 1 - \lambda_2(W)$ which describes the connectivity
88 of the network, where $\lambda_2(W)$ is the second largest eigenvalue of W . Each communication step
89 can be performed by a multiplication of W by an aggregated variable. It is popular to reduce the
90 consensus error by Chebyshev acceleration [17, 28]. We present the details in Algorithm 1 and its
91 has the following property.

92 **Proposition 1.** Let \mathbf{v}^0 and \mathbf{v}^t be the input and output of Algorithm 1 respectively, and $\bar{\mathbf{v}} = \frac{1}{m} \mathbf{1}^\top \mathbf{v}^0$.
 93 Then we have $\bar{\mathbf{v}} = \frac{1}{m} \mathbf{1}^\top \mathbf{x}^K$ and $\|\mathbf{v}^t - \mathbf{1}\bar{\mathbf{v}}\| \leq (1 - \sqrt{1 - \lambda_2(W)})^K \|\mathbf{v}^0 - \mathbf{1}\bar{\mathbf{v}}\|$.

94 3 The Algorithm and Main Results

95 In this section, we introduce the insight and the design of Katyusha-type Near-Optimal decenTralized
 96 algorithm (KNOT). We also provide complexity analysis to show the advantage of KNOT
 97 formally.

98 3.1 Motivation

99 Before studying the decentralized optimization, we first give a brief review of the algorithms for
 100 solving the finite-sum optimization problem

$$\min_{\hat{\mathbf{x}} \in \mathbb{R}^d} \hat{f}(\hat{\mathbf{x}}) \triangleq \frac{1}{m} \sum_{i=1}^m \hat{f}_i(\hat{\mathbf{x}}) \quad (3.1)$$

101 on a single machine, where the objective function $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly-convex, and each
 102 individual component function $\hat{f}_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth and convex. It is well known that accel-
 103 erated gradient descent (AGD) [22] achieves the optimal full gradient complexity $\mathcal{O}(\sqrt{\kappa} \log(1/\epsilon))$
 104 for solving problem (3.1), while each of its iteration requires m incremental first-order oracle (IFO)
 105 calls. A popular way to reduce the iteration is stochastic gradient descent (SGD), while it only con-
 106 verges sublinearly. Variance reduction [1, 5, 8, 10, 23, 24, 30] is a widely used technique to improve
 107 the convergence rate of SGD, e.g. stochastic variance reduced gradient (SVRG) method iterate with
 108 gradient estimator

$$\hat{\mathbf{v}}^t = \nabla \hat{f}(\hat{\mathbf{w}}) + \frac{1}{b} \sum_{j \in \mathcal{S}_t} \left(\nabla \hat{f}_j(\hat{\mathbf{x}}^t) - \nabla \hat{f}_j(\hat{\mathbf{w}}) \right), \quad (3.2)$$

109 where $\hat{\mathbf{w}}$ is a snapshot point that is updated infrequently and \mathcal{S}_t is a random subset of $\{1, \dots, m\}$
 110 with candidate b . Katyusha method [1] iterates with variance reduced estimator $\hat{\mathbf{v}}^t$ by involving the
 111 negative momentum and achieves the near optimal IFO complexity of $\mathcal{O}((m + \sqrt{m\kappa}) \log(1/\epsilon))$,
 112 which is better than $\mathcal{O}(m\sqrt{\kappa} \log(1/\epsilon))$ of AGD.

113 Note that the objective function in decentralized optimization problem (1.1) also has the finite-sum
 114 structure, which motivates us to improve the computational efficiency by introducing some gradient
 115 estimator like variance reduced estimator shown in (3.2).

116 3.2 The Algorithm

117 We propose Katyusha-type Near-Optimal decenTralized algorithm (KNOT) in Algorithm 2. The
 118 gradient tracking steps (line 10 and 21) indicates the algorithm performs the following update

$$\begin{cases} \bar{\mathbf{x}}^t = \theta_1 \bar{\mathbf{z}}^t + \theta_2 \bar{\mathbf{w}}^t + (1 - \theta_1 - \theta_2) \bar{\mathbf{y}}^t \\ \bar{\mathbf{s}}^t = \bar{\mathbf{v}}^t = \bar{\mathbf{u}}^t + \frac{1}{m} \sum_{i=1}^m \frac{\xi_i^t}{q} (\nabla f_i(\mathbf{x}_i^t) - \nabla f_i(\mathbf{w}_i^t)) \\ \bar{\mathbf{z}}^{t+1} = \frac{1}{1 + \eta\sigma} \left(\eta\sigma \bar{\mathbf{x}}^t + \bar{\mathbf{z}}^t - \frac{\eta}{L} \bar{\mathbf{s}}^t \right) \\ \bar{\mathbf{y}}^{t+1} = \bar{\mathbf{x}}^t + \theta_1 (\bar{\mathbf{z}}^{t+1} - \bar{\mathbf{z}}^t) \\ \bar{\mathbf{w}}^{t+1} = \begin{cases} \bar{\mathbf{y}}^t, & \text{if } \zeta^{t+1} = 1 \\ \bar{\mathbf{w}}^t, & \text{if } \zeta^{t+1} = 0 \end{cases} \\ \bar{\mathbf{u}}^{t+1} = \bar{\mathbf{g}}^{t+1} = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{w}_i^{t+1}) \end{cases}$$

Algorithm 2 Katyusha-type Near-Optimal decentralized algorithm (KNOT)

1: **Input:** initial point \bar{w}^0 , probabilities p and q , number of consensus steps K and K_{out} , total iteration numbers T , parameters L, μ, θ_1 and θ_2

2: $\mathbf{y}^0 = \mathbf{z}^0 = \mathbf{w}^0 = \mathbf{1}\bar{w}_0$, $\mathbf{v}^{-1} = \mathbf{s}^{-1} = \mathbf{0}$, $\mathbf{g} = \mathbf{u}^0 = \nabla F(\mathbf{w}^0)$

3: $\eta = 1/(13\theta_1)$, $\sigma = \mu/L$

4: **for** $t = 0, \dots, T$

5: $\mathbf{x}^t = \text{AccGossip}(\theta_1 \mathbf{z}^t + \theta_2 \mathbf{w}^t + (1 - \theta_1 - \theta_2) \mathbf{y}^t, K)$

6: **parallel for** $i = 1, \dots, m$ **do**

7: draw $\xi_i^t \sim \text{Bernoulli}(q)$

8: $v_i^t = u_i^t + \frac{\xi_i^t}{q} (\nabla f_i(x_i^t) - \nabla f_i(w_i^t))$

9: **end parallel for**

10: $\mathbf{s}^t = \text{AccGossip}(\mathbf{s}^{t-1} + \mathbf{v}^t - \mathbf{v}^{t-1}, K)$

11: $\mathbf{z}^{t+1} = \text{AccGossip}\left(\frac{1}{1 + \eta\sigma} \left(\eta\sigma \mathbf{x}^t + \mathbf{z}^t - \frac{\eta}{L} \mathbf{s}^t\right), K\right)$

12: $\mathbf{y}^{t+1} = \text{AccGossip}(\mathbf{x}^t + \theta_1(\mathbf{z}^{t+1} - \mathbf{z}^t), K)$

13: draw $\zeta^{t+1} \sim \text{Bernoulli}(p)$

14: **parallel for** $i = 1, \dots, m$ **do**

15: $\tilde{w}_i^{t+1} = \begin{cases} y_i^t, & \text{if } \zeta^{t+1} = 1 \\ w_i^t, & \text{if } \zeta^{t+1} = 0 \end{cases}$

16: **end parallel for**

17: $\mathbf{w}^{t+1} = \text{AccGossip}(\tilde{\mathbf{w}}^{t+1}, K)$

18: **parallel for** $i = 1, \dots, m$ **do**

19: $g_i^{t+1} = \begin{cases} \nabla f_i(w_i^{t+1}), & \text{if } \zeta^{t+1} = 1 \\ g_i^t, & \text{if } \zeta^{t+1} = 0 \end{cases}$

20: **end parallel for**

21: $\mathbf{u}^{t+1} = \text{AccGossip}(\mathbf{u}^t + \mathbf{g}^{t+1} - \mathbf{g}^t, K)$

22: **end for**

23: **Output:** $\mathbf{x}_{\text{out}} = \text{AccGossip}(\mathbf{x}_T, K_{\text{out}})$.

119 on the mean vectors, which is similar to mini-batch version of Loopless Katyusha (L-Katyusha) [1,
 120 10, 24]. The consensus error of the variables can be bounded by the property (Proposition 1) of the
 121 subroutine `AccGossip` (Algorithm 1), which encourages KNOT achieves the similar convergence
 122 result to Katyusha.

123 The computational efficiency of KNOT mainly comes from the local gradient estimator

$$v_i^t = u_i^t + \frac{\xi_i^t}{q} (\nabla f_i(x_i^t) - \nabla f_i(w_i^t)),$$

124 where w_i^t is an estimator of the local gradient at the snapshot point w_i^t and ξ_i^t is a random variable
 125 drawn from Bernoulli distribution with parameter q . If we set $q = b/m$ for some $b \in \{1, \dots, m\}$, the
 126 mean of local gradient estimators v_1^t, \dots, v_m^t can be rewritten as

$$\bar{v}^t = \bar{u}^t + \frac{1}{b} \sum_{i \in \mathcal{I}^t} (\nabla f_i(x_i^t) - \nabla f_i(w_i^t)),$$

127 where $\mathcal{I}^t = \{i : \xi_i^t = 1\}$ and \bar{u}^t can be regarded as an estimator of the global gradient at point \bar{w}^t .
128 This implies the computations of local gradient $\nabla f_i(x_i^t)$ and $\nabla f_i(w_i^t)$ are performed on $|\mathcal{I}^t|$ agents
129 with $\mathbb{E}[|\mathcal{I}^t|] = b$. Hence, the mean vector \bar{v}^t plays a similar role to the variance reduced gradient
130 estimator with batch-size $b = mq$, which leads to the algorithm requires less local gradient compu-
131 tation. Additionally, KNOT only computes the local gradient for all of the m agents when $\zeta^{t+1} = 1$
132 (line 19), which follows the idea of loopless framework for variance reduction [10, 24]. Since, we
133 draw ζ^{t+1} from Bernoulli distribution with parameter p , the small p leads to $\zeta^{t+1} = 1$ occurs infre-
134 quently and the computational cost for this case is not expensive in expectation. KNOT also enjoys
135 the parallel speed up property like Katyusha [1], which means to the appropriate settings for p and q
136 can reduce the number of total iterations that corresponds to less communication rounds in total.

137 3.3 Convergence Analysis

138 The convergence analysis of KNOT (Algorithm 2) is based on the Lyapunov function [24] as follows

$$V^t \triangleq \mathcal{Z}^t + \mathcal{Y}^t + \mathcal{W}^t, \quad (3.3)$$

139 where

$$\mathcal{Z}^t \triangleq \frac{L(1+\eta\sigma)}{2\eta} \|\bar{z}^t - x^*\|^2, \quad \mathcal{Y}^t \triangleq \frac{1}{\theta_1} (f(\bar{y}^t) - f(x^*)) \quad \text{and} \quad \mathcal{W}^t \triangleq \frac{\theta_2}{p\gamma\theta_1} (f(\bar{w}^t) - f(x^*)).$$

140 The parameters $\eta, \sigma, \theta_1, \theta_2$ in (3.3) follow the notations of Algorithm 2 and we let $\gamma \in (1/2, 1)$.

141 The decentralized setting leads to we cannot directly follow the analysis of existing Katyusha-type
142 algorithms [1, 10, 24]. Different from previous work, the recursion on Lyapunov function V^t for
143 KNOT contains the additional terms of consensus error, which is shown in the following lemma.

144 **Lemma 1.** *Under Assumption 1, 2, 3 and 4, if we choose $\eta = 1/(13\theta_1)$, Algorithm 2 holds that*

$$\begin{aligned} \mathbb{E}[V^{t+1}] \leq & \max\left(\frac{1}{1+\eta\sigma}, 1 - \left(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma}\right), 1 - p(1-\gamma)\right) V^t \\ & + \sqrt{\frac{2\eta LV^t}{(1+\eta\sigma)m}} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\| + \frac{L}{3m^2q\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{L}{4m^2q\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2. \end{aligned}$$

145 We consider the consensus error by introducing the vector

$$r^t = \frac{L}{m} \left[\|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2, \frac{\eta^2}{L^2} \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2, \frac{\eta^2}{L^2} \|\mathbf{s}^t - \mathbf{1}\bar{s}^t\|^2, \|\mathbf{z}^t - \mathbf{1}\bar{z}^t\|^2, \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2, \|\mathbf{y}^t - \mathbf{1}\bar{y}^t\|^2 \right]^\top.$$

146 We describe the convergence of r_t by a linear system as follows.

147 **Lemma 2.** *Under the settings of Lemma 1, we run Algorithm 2 by taking*

$$K = \left\lceil \frac{\log(1/\rho)}{\sqrt{1 - \lambda_2(W)}} \right\rceil$$

148 with

$$\rho \leq \frac{q\theta_1^3}{9} \min\left\{ \frac{1}{2}\eta^2\sigma^2, \frac{1}{2}\left(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma}\right), \frac{1}{2}p(1-\gamma) \right\}.$$

149 Then it holds that

$$\mathbb{E}[r^{t+1}] \leq \rho^2 (\mathbf{B} \cdot \mathbf{A} \cdot r^t + e^t)$$

150 for some matrix $\mathbf{A} \in \mathbb{R}^{6 \times 6}$, elementary matrix $\mathbf{B} \in \mathbb{R}^{6 \times 6}$ and vector $e^t \in \mathbb{R}^6$ satisfy²

$$\|\mathbf{A}\| \leq \frac{4}{q\theta_1^2}, \quad \|\mathbf{B}\| \leq 2 \quad \text{and} \quad \|e^t\| < \frac{2}{3q\theta_1^2} (V^{t+1} + V^t).$$

²The expressions of \mathbf{A} , \mathbf{B} and e_t are very complicated and we present them in appendix.

151 By connecting above two lemmas, we obtain the main results.

152 **Theorem 1** (main result). *Under Assumption 1, 2, 3 and 4, we run Algorithm 2 with*

$$p = \max \left\{ \frac{1}{\sqrt{m}}, \frac{1}{\sqrt{\kappa}} \right\}, \quad q = \min \left\{ \frac{1}{\sqrt{m}}, \frac{\sqrt{\kappa}}{m} \right\}, \quad \gamma \in \left(\frac{2}{3}, 1 \right)$$

$$\theta_2 = \frac{1}{2mq}, \quad \theta_1 = \min \left\{ \sqrt{\frac{mq}{\kappa p}} \theta_2, \theta_2 \right\}, \quad \eta = \frac{1}{13\theta_1}$$

153 and take K by following the setting of Lemma 2. Then it holds that

$$\mathbb{E} [V^t] \leq \left(1 - \min \left\{ \eta\sigma, \frac{\theta_1 + \theta_2 - \theta_2/\gamma}{2}, \frac{p(1-\gamma)}{2} \right\} \right)^t (V^0 + \|\mathbf{r}^0\|)$$

154 and

$$\mathbb{E} \left[\frac{L}{m} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 \right] \leq \left(\frac{8}{243} + 2^{-t} \right) \left(1 - \min \left\{ \eta\sigma, \frac{\theta_1 + \theta_2 - \theta_2/\gamma}{2}, \frac{p(1-\gamma)}{2} \right\} \right)^t \cdot (V^0 + \|\mathbf{r}^0\|).$$

155 Theorem 1 establish the linear convergence for the function value at the point of mean vectors. The
 156 Bernoulli variables in the algorithm indicate each iteration has $\mathcal{O}(m/\sqrt{\kappa} + \sqrt{m})$ local gradient
 157 calls in expectation. Hence, we obtain the upper bounds of local gradient oracle complexity and
 158 communication complexity for finding an ϵ -suboptimal solution.

159 **Corollary 2.** *Under the settings of Theorem 1, Algorithm 2 can achieve an ϵ -suboptimal solution
 160 by taking $T = \mathcal{O}(\sqrt{\kappa} \log(1/\epsilon))$ and $K_{\text{out}} = \tilde{\mathcal{O}}(\sqrt{1/\alpha})$, which requires the local first-order oracle
 161 complexity of $\mathcal{O}((m + \sqrt{m\kappa}) \log(1/\epsilon))$ and the communication complexity of $\tilde{\mathcal{O}}(\sqrt{\kappa/\alpha} \log(1/\epsilon))$
 162 in expectation.*

163 4 Discussion for the Optimality

164 In this section, we verify the optimality of the proposed algorithms. We first provide follow the
 165 statement for the lower bounds of decentralized strongly convex optimization provided by Kovalev
 166 et al. [11], which is a direct application of Corollary 2 from Scaman et al. [29] but does not include
 167 the dual gradient oracle.

168 **Proposition 2.** *For any $m \geq 2$ and $\alpha > 1$, there exist a gossip matrix $W \in \mathbb{R}^{m \times m}$ satisfying
 169 $1 - \lambda_2(W) = \alpha$ and a family of smooth strongly convex functions $\{f_i : \mathbb{R}^d \rightarrow \mathbb{R}\}_{i=1}^m$ with condition
 170 number κ such that the following holds: for any $\epsilon > 0$, any first-order decentralized algorithm
 171 requires at least $\Omega(\sqrt{\kappa/\alpha} \log(1/\epsilon))$ communication rounds and at least $\Omega(\sqrt{\kappa} \log(1/\epsilon))$ gradient
 172 steps to output $\mathbf{x} = [x_1, \dots, x_m]^\top$ such that $f(x_i) - f(x^*) \leq \epsilon$ for all $i = 1, \dots, m$, where x^* is
 173 the minimizer of $f(x) = \frac{1}{m} \sum_{i=1}^m f_i(x)$.*

174 It is worth pointing out that the concept “gradient step” described in Proposition 2 only requires the
 175 gradient computation should depend on the history local points of the corresponding agent, while it
 176 does not contain any requirement on the number of agents that participate into their local gradients
 177 computation. This implies the “gradient steps” lower bound of $\Omega(\sqrt{\kappa} \log(1/\epsilon))$ described in this
 178 proposition corresponds to the iterations number of proposed KNOT (Algorithm 2), rather than the
 179 number of local gradient calls. Hence, the result of Corollary 2 means KNOT matches the “gradient
 180 steps” lower bound and nearly matches the communication lower bound provided by Proposition 2.

181 Compared with the number of “gradient steps”, we are more interested in the number of local gra-
 182 dient calls, which essentially reflects the totally computational cost of a decentralized optimization
 183 algorithm. The lower bound of local gradient calls can be established by considering the IFO calls
 184 for the finite-sum optimization problem on single machine. Woodworth and Srebro [34] provide the
 185 following lower bound for solving the finite-sum optimization problem by randomized first-order
 186 (non-distributed) methods.

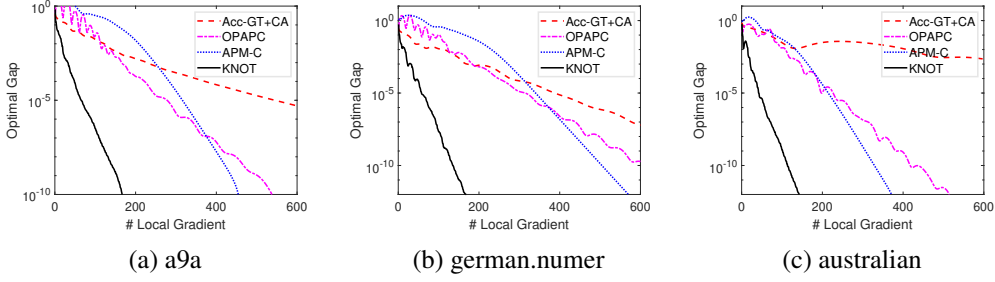


Figure 1: Comparison for the number of local gradient calls vs. optimal gap.

187 **Proposition 3.** For any $m \geq 2$ and $\kappa > 161m$, there exist a family of smooth strongly convex
 188 functions $\{\hat{f}_i : \mathbb{R}^d \rightarrow \mathbb{R}\}_{i=1}^m$ with condition number κ such that the following holds: for any
 189 $\epsilon > 0$, any randomized algorithm require at least $\tilde{\Omega}(m + \sqrt{m\kappa} \log(1/\epsilon))$ IFO calls to output x
 190 such that $\mathbb{E}[f(\hat{x}) - f(x^*)] < \epsilon$, where x^* is the minimizer of $\hat{f}(x) = \frac{1}{m} \sum_{i=1}^m \hat{f}_i(x)$.

191 We can view the individual functions $\{\hat{f}_i\}_{i=1}^m$ in Proposition 3 as the local functions in decentralized
 192 optimization on a fully connected network, then the IFO lower bound $\tilde{\Omega}(m + \sqrt{m\kappa} \log(1/\epsilon))$ just
 193 corresponds to the local gradient lower bound in our decentralized optimization problem. Hence,
 194 Proposition 3 implies the local gradient oracle complexity of proposed KNOT is near optimal.

195 5 Experiments

196 In this section, we provide the numerical experiments to evaluate the performance of proposed
 197 KNOT. We consider the ℓ_2 -regularized logistic regression for binary classification. We formulate
 198 this model by the optimization problem

$$\min_{x \in \mathbb{R}^d} f(x) \triangleq \frac{1}{m} \sum_{i=1}^m f_i(x)$$

199 with

$$f_i(x) = \frac{1}{n} \sum_{j=1}^n \log(1 + \exp(-b_{ij} a_{ij}^\top x)) + \frac{\mu}{2} \|x\|^2,$$

200 where $a_{ij} \in \mathbb{R}^d$ is the feature vector of the j -th sample on agent i , $b_{ij} \in \{-1, 1\}$ is the corresponding
 201 label and $\mu > 0$ is the hyperparameter.

202 We conduct our experiments on three real-world datasets “a9a”, “german.number” and “australian”
 203 which can be found in LIBSVM repository [2]. We let $m = 300$ and $\mu = 0.01$. We set the
 204 gossip matrix W by corresponding a random graph that each pair in the network is connected with
 205 probability $1/30$, which leads to $1 - \lambda_2(W) \approx 0.0382$.

206 We compare the proposed method KNOT with baseline algorithms ACC-GT+CA [13], OPAPC [11]
 207 and APM-C [15]. For KNOT, we set the parameters p, q, θ_1, θ_2 and η by following the settings of
 208 Theorem 1 and tune K from $\{1, 5, 10\}$. For the baseline algorithms, we also select their parameters
 209 by following the corresponding theoretical analysis.

210 We present the experimental results for the computational cost and communication cost in Figures 1
 211 and 2 respectively, where the y -axis represents the optimal gap which is defined as

$$\frac{1}{m} \sum_{i=1}^m f(x_i) - f(x^*).$$

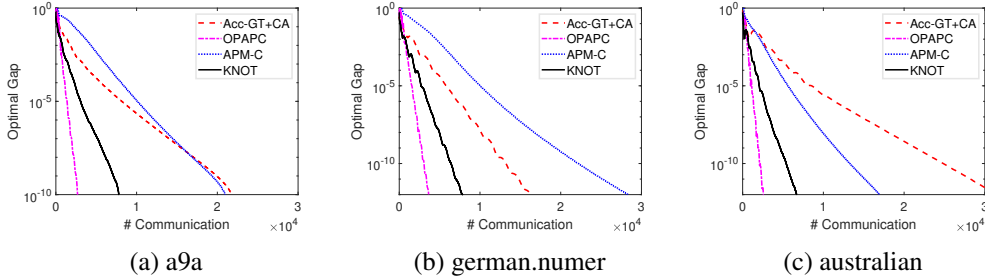


Figure 2: Comparison for the number of communication rounds vs. optimal gap.

212 We observe that the proposed KNOT always has significantly better computational efficiency than
 213 than all of baseline methods. For the communication complexity, the result of KNOT is comparable
 214 to OPAPC and better than other baseline methods.

215 6 Conclusion

216 In this paper, we study decentralized strongly convex optimization and propose a novel method
 217 called Katyusha-type Near-Optimal decENtralized algorithm (KNOT), which avoids computing all
 218 of the local gradients in one iteration. The theoretical analysis shows that our method is near optimal
 219 to both the local first-order oracle complexity and the communication complexity. The empirical
 220 studies on regularized logistic regression problem also supports our theoretical results. We believe
 221 the idea of KNOT is not limited to first-order optimization for convex problems. It is possible to
 222 extend the framework of KNOT to solve variational inequalities. We can also try to design the
 223 second-order decentralized algorithms with less local Hessian calls.

224 References

- 225 [1] Zeyuan Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. In
 226 *STOC*, 2017.
- 227 [2] Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM*
 228 *transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011. URL <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>.
- 230 [3] Peijun Chen, Jianguo Huang, and Xiaoqun Zhang. A primal–dual fixed point algorithm for
 231 convex separable minimization with applications to image restoration. *Inverse Problems*, 29
 232 (2):025011, 2013.
- 233 [4] Laurent Condat, Daichi Kitahara, Andrés Contreras, and Akira Hirabayashi. Proximal splitting
 234 algorithms: Relax them all. *arXiv preprint arXiv:1912.00137*, 2019.
- 235 [5] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient
 236 method with support for non-strongly convex composite objectives. In *NIPS*, 2014.
- 237 [6] Yoel Drori, Shoham Sabach, and Marc Teboulle. A simple algorithm for a class of nonsmooth
 238 convex–concave saddle-point problems. *Operations Research Letters*, 43(2):209–214, 2015.
- 239 [7] Darina Dvinskikh and Alexander Gasnikov. Decentralized and parallel primal and dual accel-
 240 erated methods for stochastic convex programming problems. *Journal of Inverse and Ill-posed*
 241 *Problems*, 29(3):385–405, 2021.
- 242 [8] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive vari-
 243 ance reduction. In *NIPS*, 2013.

- 244 [9] Il-Young Joo and Dae-Hyun Choi. Distributed optimization framework for energy management
245 of multiple smart homes with distributed energy resources. *IEEE Access*, 5:15551–15560,
246 2017.
- 247 [10] Dmitry Kovalev, Samuel Horváth, and Peter Richtárik. Don’t jump through hoops and remove
248 those loops: SVRG and Katyusha are better without the outer loop. In *ALT*, 2020.
- 249 [11] Dmitry Kovalev, Adil Salim, and Peter Richtárik. Optimal and practical algorithms for smooth
250 and strongly convex decentralized optimization. In *NeurIPS*, 2020.
- 251 [12] Huan Li and Zhouchen Lin. Revisiting EXTRA for smooth distributed optimization. *SIAM*
252 *Journal on Optimization*, 30(3):1795–1821, 2020.
- 253 [13] Huan Li and Zhouchen Lin. Accelerated gradient tracking over time-varying graphs for de-
254 centralized optimization. *arXiv preprint arXiv:2104.02596*, 2021.
- 255 [14] Huan Li, Cong Fang, Wotao Yin, and Zhouchen Lin. A sharp convergence rate analysis for
256 distributed accelerated gradient methods. *arXiv preprint arXiv:1810.01053*, 2018.
- 257 [15] Huan Li, Cong Fang, Wotao Yin, and Zhouchen Lin. Decentralized accelerated gradient meth-
258 ods with increasing penalty parameters. *IEEE transactions on Signal Processing*, 68:4855–
259 4870, 2020.
- 260 [16] Zhi Li, Wei Shi, and Ming Yan. A decentralized proximal-gradient method with network inde-
261 pendent step-sizes and separated convergence rates. *IEEE Transactions on Signal Processing*,
262 67(17):4494–4506, 2019.
- 263 [17] Ji Liu and A. Stephen Morse. Accelerated linear iterations for distributed averaging. *Annual*
264 *Reviews in Control*, 35(2):160–165, 2011.
- 265 [18] Ignace Loris and Caroline Verhoeven. On a generalization of the iterative soft-thresholding
266 algorithm for the case of non-separable penalty. *Inverse Problems*, 27(12):125007, 2011.
- 267 [19] Angelia Nedic. Asynchronous broadcast-based convex optimization over a network. *IEEE*
268 *Transactions on Automatic Control*, 56(6):1337–1351, 2010.
- 269 [20] Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent opti-
270 mization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- 271 [21] Angelia Nedic, Alex Olshevsky, and Wei Shi. Achieving geometric convergence for distributed
272 optimization over time-varying graphs. *SIAM Journal on Optimization*, 27(4):2597–2633,
273 2017.
- 274 [22] Yurii Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018.
- 275 [23] Lam M. Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. SARAH: A novel method for
276 machine learning problems using stochastic recursive gradient. In *ICML*, 2017.
- 277 [24] Xun Qian, Zheng Qu, and Peter Richtárik. L-SVRG and L-Katyusha with arbitrary sampling.
278 *Journal of Machine Learning Research*, 22(1):4991–5039, 2021.
- 279 [25] Guannan Qu and Na Li. Harnessing smoothness to accelerate distributed optimization. *IEEE*
280 *Transactions on Control of Network Systems*, 5(3):1245–1260, 2017.
- 281 [26] Michael G. Rabbat and Robert D. Nowak. Decentralized source localization and tracking
282 [wireless sensor networks]. In *ICASSP*, 2004.
- 283 [27] S. Sundhar Ram, A. Nedich, and Venugopal V. Veeravalli. Distributed stochastic subgradient
284 projection algorithms for convex optimization. *Journal of Optimization Theory and Applica-*
285 *tions*, 147:516–545, 2010.

- 286 [28] Youcef Saad. Chebyshev acceleration techniques for solving nonsymmetric eigenvalue prob-
287 lems. *Mathematics of Computation*, 42(166):567–588, 1984.
- 288 [29] Kevin Scaman, Francis Bach, Sébastien Bubeck, Yin Tat Lee, and Laurent Massoulié. Optimal
289 algorithms for smooth and strongly convex distributed optimization in networks. In *ICML*,
290 2017.
- 291 [30] Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic
292 average gradient. *Mathematical Programming*, 162:83–112, 2017.
- 293 [31] Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. EXTRA: An exact first-order algorithm for
294 decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- 295 [32] Zhuoqing Song, Lei Shi, Shi Pu, and Ming Yan. Optimal gradient tracking for decentralized
296 optimization. *arXiv preprint arXiv:2110.05282*, 2021.
- 297 [33] Mu Wang, Changqiao Xu, Xingyan Chen, Lujie Zhong, Zhonghui Wu, and Dapeng Oliver
298 Wu. Bc-mobile device cloud: A blockchain-based decentralized truthful framework for mobile
299 device cloud. *IEEE Transactions on Industrial Informatics*, 17(2):1208–1219, 2020.
- 300 [34] Blake E. Woodworth and Nati Srebro. Tight complexity bounds for optimizing composite
301 objectives. In *NIPS*, 2016.
- 302 [35] Ran Xin and Usman A. Khan. A linear algorithm for optimization over directed graphs with
303 geometric convergence. *IEEE Control Systems Letters*, 2(3):315–320, 2018.
- 304 [36] Jinming Xu, Shanying Zhu, Yeng Chai Soh, and Lihua Xie. Augmented distributed gradient
305 methods for multi-agent optimization under uncoordinated constant stepsizes. In *CDC*, 2015.
- 306 [37] Haishan Ye, Luo Luo, Ziang Zhou, and Tong Zhang. Multi-consensus decentralized acceler-
307 ated gradient descent. *arXiv preprint arXiv:2005.00797*, 2020.
- 308 [38] Haishan Ye, Ziang Zhou, Luo Luo, and Tong Zhang. Decentralized accelerated proximal
309 gradient descent. In *NeurIPS*, 2020.
- 310 [39] Kun Yuan, Qing Ling, and Wotao Yin. On the convergence of decentralized gradient descent.
311 *SIAM Journal on Optimization*, 26(3):1835–1854, 2016.

312 The theoretical analysis of KNOT is organized as follows.

313 • We first consider the mean vector and provide the convergence of $f(\bar{x}^t) - f(x^*)$. The result is
 314 shown in Lemma 1, and its detailed proof is shown in Appendix A.

315 • We then consider the consensus error, which is characterized by vector

$$r^t = \frac{L}{m} \left[\|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2, \frac{\eta^2}{L^2} \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2, \frac{\eta^2}{L^2} \|\mathbf{s}^t - \mathbf{1}\bar{s}^t\|^2, \|\mathbf{z}^t - \mathbf{1}\bar{z}^t\|^2, \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2, \|\mathbf{y}^t - \mathbf{1}\bar{y}^t\|^2 \right]^\top.$$

316 The recursion of r_t is established in Lemma 2, and its detailed proof is shown in Appendix B.
 317 Especially, we present the expression of \mathbf{A} , \mathbf{B} and e_t in the statement of Lemma 15.

318 • We finally apply Lemma 1 and 2 to obtain our final results Theorem 1 and Corollary 2, whose
 319 detailed proofs are shown in Appendix C and D respectively.

320 A Proof of Lemma 1

321 In this section, we focus on analyzing Lyapunov function

$$V^t \triangleq \mathcal{Z}^t + \mathcal{Y}^t + \mathcal{W}^t,$$

322 where

$$\mathcal{Z}^t \triangleq \frac{L(1 + \eta\sigma)}{2\eta} \|\bar{z}^t - x^*\|^2, \quad \mathcal{Y}^t \triangleq \frac{1}{\theta_1} (f(\bar{y}^t) - f(x^*)) \quad \text{and} \quad \mathcal{W}^t \triangleq \frac{\theta_2}{p\gamma\theta_1} (f(\bar{w}^t) - f(x^*)).$$

323 The analysis is more complicated than the counterpart of L-Katyusha [10, 24] because of the con-
 324 sensus error aroused from the decentralized setting.

325 We substitute $q = b/m$ in the following proof, where b can be regarded as expected mini-batch size.

326 Let us begin with a useful lemma of L -smooth and convex functions for our further analysis.

327 **Lemma 3** ([22]). *Under Assumption 1 and 3, it holds that*

$$\frac{1}{2L} \|\nabla f_i(x) - \nabla f_i(y)\|^2 \leq f_i(x) - f_i(y) - \langle \nabla f_i(y), x - y \rangle, \quad (\text{A.1})$$

328 for all i and x, y .

329 We show that the average of local gradient trackers can approximate $\nabla f(\bar{x}^t)$ well first.

330 **Lemma 4.** *Under the settings of Lemma 1, Algorithm 2 holds that*

$$\bar{s}^t = \frac{1}{m} \sum_i^m v_i^t \quad \text{and} \quad \mathbb{E}[\bar{s}^t] = \frac{1}{m} \sum_i^m \nabla f_i(\mathbf{x}_i^t).$$

331 Furthermore, we have

$$\|\nabla f(\bar{x}^t) - \mathbb{E}[\bar{s}^t]\| \leq \frac{L}{\sqrt{m}} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|.$$

332 *Proof.* We have

$$\begin{aligned} \|\nabla f(\bar{x}^t) - \mathbb{E}[\bar{s}^t]\|^2 &= \left\| \frac{1}{m} \sum_{i=1}^m (\nabla f_i(\mathbf{x}_i^t) - \nabla f_i(\bar{x}^t)) \right\|^2 \\ &\leq \frac{1}{m} \sum_{i=1}^m \|\nabla f_i(\mathbf{x}_i^t) - \nabla f_i(\bar{x}^t)\|^2 \\ &\leq \frac{L^2}{m} \sum_{i=1}^m \|\mathbf{x}_i^t - \bar{x}^t\|^2 \\ &= \frac{L^2}{m} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2, \end{aligned}$$

333 where the first equality is due to $\bar{s}^t = \bar{v}^t$ and

$$\mathbb{E}[\bar{v}^t] = \frac{1}{m} \sum_{i=1}^m \nabla f_i(x_i^t)$$

334 hold for Algorithm 2. □

335 Then we provide some lemmas for the mean vectors.

336 **Lemma 5.** *Under the settings of Lemma 1, it holds that*

$$\begin{aligned} \mathbb{E} [\|\bar{s}^t - \nabla f(\bar{x}^t)\|^2] &\leq \frac{12L}{b} (f(\bar{w}^t) - f(\bar{x}^t) - \langle \nabla f(\bar{x}^t), \bar{w}^t - \bar{x}^t \rangle) + \frac{8L^2}{mb} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 \\ &\quad + \frac{6L^2}{mb} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2. \end{aligned} \quad (\text{A.2})$$

337 *Proof.* We have

$$\begin{aligned} &\mathbb{E} \left[\|\bar{s}^t - \nabla f(\bar{x}^t)\|^2 \right] \\ &\stackrel{\text{Alg. 2}}{=} \mathbb{E} \left[\left\| \bar{u}^t + \frac{1}{qm} \sum_{j=1}^m \xi_j (\nabla f_j(x_j^t) - \nabla f_j(w_j^t)) - \nabla f(\bar{x}^t) \right\|^2 \right] \\ &= \mathbb{E} \left[\left\| \frac{1}{m} \sum_{j=1}^m \frac{\xi_j}{q} (\nabla f_j(x_j^t) - \nabla f_j(w_j^t)) - \mathbb{E} [\nabla f_j(x_j^t) - \nabla f_j(w_j^t)] + (\mathbb{E}[\bar{s}^t] - \nabla f(\bar{x}^t)) \right\|^2 \right] \\ &\leq \frac{2}{m} \mathbb{E} \left[\left\| \frac{\xi_j}{q} (\nabla f_j(x_j^t) - \nabla f_j(w_j^t)) \right\|^2 \right] + 2\mathbb{E} [\|\mathbb{E}[\bar{s}^t] - \nabla f(\bar{x}^t)\|^2] \\ &= \frac{2}{mq} \mathbb{E} [\|\nabla f_j(x_j^t) - \nabla f_j(w_j^t)\|^2] + 2\mathbb{E} [\|\mathbb{E}[\bar{s}^t] - \nabla f(\bar{x}^t)\|^2] \\ &\leq \frac{6}{mq} \mathbb{E} [\|\nabla f_j(x_j^t) - \nabla f_j(\bar{x}^t)\|^2 + \|\nabla f_j(\bar{x}^t) - \nabla f_j(\bar{w}^t)\|^2 + \|\nabla f_j(\bar{w}^t) - \nabla f_j(w_j^t)\|^2] \\ &\quad + 2\mathbb{E} [\|\mathbb{E}[\bar{s}^t] - \nabla f(\bar{x}^t)\|^2] \\ &\leq \frac{12L}{mq} (f(\bar{w}^t) - f(\bar{x}^t) - \langle \nabla f(\bar{x}^t), \bar{w}^t - \bar{x}^t \rangle) + \frac{8L^2}{m^2q} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{6L^2}{m^2q} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \\ &= \frac{12L}{b} (f(\bar{w}^t) - f(\bar{x}^t) - \langle \nabla f(\bar{x}^t), \bar{w}^t - \bar{x}^t \rangle) + \frac{8L^2}{mb} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{6L^2}{mb} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2, \end{aligned}$$

338 where the first equality is because of the fact that $\bar{u}^t = \mathbb{E} [\nabla f_j(x_j^t)]$ and the first inequality is
 339 because of the fact that $\mathbb{E}[\|z - \mathbb{E}[z]\|^2] \leq \mathbb{E}[\|z\|^2]$ and the property of variance; the last inequality
 340 is because of Lemma 4.

341 □

342 **Lemma 6.** *Under the settings of Lemma 1, we have*

$$\langle \bar{s}^t, x^* - \bar{z}^{t+1} \rangle + \frac{\mu}{2} \|\bar{x}^t - x^*\|^2 \geq \frac{L}{2\eta} \|\bar{z}^t - \bar{z}^{t+1}\|^2 + \mathcal{Z}^{t+1} - \frac{1}{1 + \eta\sigma} \mathcal{Z}^t. \quad (\text{A.3})$$

343 *Proof.* We start with the definition of \mathbf{z}^{t+1}

$$\mathbf{z}^{t+1} \stackrel{\text{Alg. 2}}{=} \frac{1}{1 + \eta\sigma} \left(\eta\sigma x^t + \mathbf{z}^t - \frac{\eta}{L} \mathbf{s}^t \right),$$

344 which means

$$\frac{\eta}{L} \bar{s}^t = \eta\sigma(\bar{x}^t - \bar{z}^{t+1}) + (\bar{z}^t - \bar{z}^{t+1}).$$

345 It further implies that

$$\begin{aligned}
& \langle \bar{s}^t, \bar{z}^{t+1} - x^* \rangle \\
&= \mu \langle \bar{x}^t - \bar{z}^{t+1}, \bar{z}^{t+1} - x^* \rangle + \frac{L}{\eta} \langle \bar{z}^t - \bar{z}^{t+1}, \bar{z}^{t+1} - x^* \rangle \\
&= \frac{\mu}{2} \left(\|\bar{x}^t - x^*\|^2 - \|\bar{x}^t - \bar{z}^{t+1}\|^2 - \|\bar{z}^{t+1} - x^*\|^2 \right) \\
&\quad + \frac{L}{2\eta} \left(\|\bar{z}^t - x^*\|^2 - \|\bar{z}^t - \bar{z}^{t+1}\|^2 - \|\bar{z}^{t+1} - x^*\|^2 \right) \\
&\leq \frac{\mu}{2} \|\bar{x}^t - x^*\|^2 + \frac{L}{2\eta} \left(\|\bar{z}^t - x^*\|^2 - (1 + \eta\sigma) \|\bar{z}^{t+1} - x^*\|^2 \right) - \frac{L}{2\eta} \|\bar{z}^t - \bar{z}^{t+1}\|^2.
\end{aligned}$$

346

□

347 **Lemma 7.** Under the settings of Lemma 1, we have

$$\frac{1}{\theta_1} (f(\bar{y}^{t+1}) - f(\bar{x}^t)) - \frac{1}{24L\theta_1} \|\bar{s}^t - \nabla f(\bar{x}^t)\|^2 \leq \frac{L}{2\eta} \|\bar{z}^{t+1} - \bar{z}^t\|^2 + \langle \bar{s}^t, \bar{z}^{t+1} - \bar{z}^t \rangle. \quad (\text{A.4})$$

348 *Proof.* We have

$$\begin{aligned}
& \frac{L}{2\eta} \|\bar{z}^{t+1} - \bar{z}^t\|^2 + \langle \bar{s}^t, \bar{z}^{t+1} - \bar{z}^t \rangle \\
&= \frac{1}{\theta_1} \left(\frac{L}{2\eta\theta_1} \|\theta_1(\bar{z}^{t+1} - \bar{z}^t)\|^2 + \langle \bar{s}^t, \theta_1(\bar{z}^{t+1} - \bar{z}^t) \rangle \right) \\
&\stackrel{\text{Alg. 2}}{=} \frac{1}{\theta_1} \left(\frac{L}{2\eta\theta_1} \|\bar{y}^{t+1} - \bar{x}^t\|^2 + \langle \bar{s}^t, \bar{y}^{t+1} - \bar{x}^t \rangle \right) \\
&= \frac{1}{\theta_1} \left(\frac{L}{2\eta\theta_1} \|\bar{y}^{t+1} - \bar{x}^t\|^2 + \langle \nabla f(\bar{x}^t), \bar{y}^{t+1} - \bar{x}^t \rangle + \langle \bar{s}^t - \nabla f(\bar{x}^t), \bar{y}^{t+1} - \bar{x}^t \rangle \right) \\
&= \frac{1}{\theta_1} \left(\frac{L}{2} \|\bar{y}^{t+1} - \bar{x}^t\|^2 + \langle \nabla f(\bar{x}^t), \bar{y}^{t+1} - \bar{x}^t \rangle + \frac{L}{2} \left(\frac{1}{\eta\theta_1} - 1 \right) \|\bar{y}^{t+1} - \bar{x}^t\|^2 + \langle \bar{s}^t - \nabla f(\bar{x}^t), \bar{y}^{t+1} - \bar{x}^t \rangle \right) \\
&\geq \frac{1}{\theta_1} \left(f(\bar{y}^{t+1}) - f(\bar{x}^t) + \frac{L}{2} \left(\frac{1}{\eta\theta_1} - 1 \right) \|\bar{y}^{t+1} - \bar{x}^t\|^2 + \langle \bar{s}^t - \nabla f(\bar{x}^t), \bar{y}^{t+1} - \bar{x}^t \rangle \right) \\
&\geq \frac{1}{\theta_1} \left(f(\bar{y}^{t+1}) - f(\bar{x}^t) - \frac{\eta\theta_1}{2L(1 - \eta\theta_1)} \|\bar{s}^t - \nabla f(\bar{x}^t)\|^2 \right) \\
&= \frac{1}{\theta_1} \left(f(\bar{y}^{t+1}) - f(\bar{x}^t) - \frac{1}{24L} \|\bar{s}^t - \nabla f(\bar{x}^t)\|^2 \right),
\end{aligned}$$

349 where the last inequality uses the Young's inequality in the form of

$$\langle a, b \rangle \geq -\frac{\|a\|^2}{2\beta} - \frac{\beta\|b\|^2}{2} \quad \text{with } \beta = \frac{\eta\theta_1}{L(1 - \eta\theta_1)}$$

350 and the last equality is because of the setting $\eta = 1/(13\theta_1)$.

□

351 **Lemma 8.** Under the settings of Lemma 1, we have

$$\mathbb{E} [\mathcal{W}^{t+1}] = (1 - p)\mathcal{W}^t + \frac{\theta_2}{\gamma}\mathcal{Y}^t.$$

352 *Proof.* From Algorithm 2, we know that

$$\mathbb{E} [f(\bar{w}^{t+1})] = (1 - p)f(\bar{w}^t) + pf(\bar{y}^t).$$

353 Then from the definition of \mathcal{W}^t and \mathcal{Y}^t , the lemma naturally holds.

□

354 Using the above lemmas, we prove Lemma 1 as follows.

355 *Proof for Lemma 1.* Combining Lemma 4, 5, 6, 7 and 8, we obtain

$$\begin{aligned}
f(x^*) &\stackrel{\text{Asm. 3}}{\geq} f(\bar{x}^t) + \langle \nabla f(\bar{x}^t), x^* - \bar{x}^t \rangle + \frac{\mu}{2} \|\bar{x}^t - x^*\|^2 \\
&= f(\bar{x}^t) + \frac{\mu}{2} \|\bar{x}^t - x^*\|^2 + \langle \nabla f(\bar{x}^t), x^* - \bar{z}^t + \bar{z}^t - \bar{x}^t \rangle \\
&\stackrel{\text{Alg. 2}}{=} f(\bar{x}^t) + \frac{\mu}{2} \|\bar{x}^t - x^*\|^2 + \langle \nabla f(\bar{x}^t), x^* - \bar{z}^t \rangle + \frac{\theta_2}{\theta_1} \langle \nabla f(\bar{x}^t), \bar{x}^t - \bar{w}^t \rangle \\
&\quad + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} \langle \nabla f(\bar{x}^t), \bar{x}^t - \bar{y}^t \rangle \\
&\geq f(\bar{x}^t) + \frac{\theta_2}{\theta_1} \langle \nabla f(\bar{x}^t), \bar{x}^t - \bar{w}^t \rangle + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(\bar{x}^t) - f(\bar{y}^t)) \\
&\quad + \mathbb{E} \left[\frac{\mu}{2} \|\bar{x}^t - x^*\|^2 + \langle \bar{s}^t, x^* - z^{t+1} \rangle + \langle \bar{s}^t, z^{t+1} - z^t \rangle \right] + \langle \nabla f(\bar{x}^t) - \mathbb{E}[\bar{s}^t], x^* - \bar{z}^t \rangle \\
&\stackrel{(A.3)}{\geq} f(\bar{x}^t) + \frac{\theta_2}{\theta_1} \langle \nabla f(\bar{x}^t), \bar{x}^t - \bar{w}^t \rangle + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(\bar{x}^t) - f(\bar{y}^t)) \\
&\quad + \mathbb{E} \left[\mathcal{Z}^{t+1} - \frac{1}{1 + \eta\sigma} \mathcal{Z}^t \right] + \mathbb{E} \left[\langle \bar{s}^t, z^{t+1} - z^t \rangle + \frac{L}{2\eta} \|z^t - z^{t+1}\|^2 \right] \\
&\quad + \langle \nabla f(\bar{x}^t) - \mathbb{E}[\bar{s}^t], x^* - \bar{z}^t \rangle \\
&\stackrel{(A.4)}{\geq} f(\bar{x}^t) + \frac{\theta_2}{\theta_1} \langle \nabla f(\bar{x}^t), \bar{x}^t - \bar{w}^t \rangle + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(\bar{x}^t) - f(\bar{y}^t)) \\
&\quad + \mathbb{E} \left[\mathcal{Z}^{t+1} - \frac{1}{1 + \eta\sigma} \mathcal{Z}^t \right] + \mathbb{E} \left[\frac{1}{\theta_1} (f(\bar{y}^{t+1}) - f(\bar{x}^t)) - \frac{1}{24L\theta_1} \|\bar{s}^t - \nabla f(\bar{x}^t)\|^2 \right] \\
&\quad + \langle \nabla f(\bar{x}^t) - \mathbb{E}[\bar{s}^t], x^* - \bar{z}^t \rangle \\
&\stackrel{(A.2)}{\geq} f(\bar{x}^t) + \frac{\theta_2}{\theta_1} \langle \nabla f(\bar{x}^t), \bar{x}^t - \bar{w}^t \rangle + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(\bar{x}^t) - f(\bar{y}^t)) + \mathbb{E} \left[\mathcal{Z}^{t+1} - \frac{1}{1 + \eta\sigma} \mathcal{Z}^t \right] \\
&\quad + \mathbb{E} \left[\frac{1}{\theta_1} (f(\bar{y}^{t+1}) - f(\bar{x}^t)) - \frac{\theta_2}{\theta_1} (f(\bar{w}^t) - f(\bar{x}^t) - \langle \nabla f(\bar{x}^t), \bar{w}^t - \bar{x}^t \rangle) \right] \\
&\quad + \langle \nabla f(\bar{x}^t) - \mathbb{E}[\bar{s}^t], x^* - \bar{z}^t \rangle - \frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 - \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \\
&= f(\bar{x}^t) + \frac{(1 - \theta_1 - \theta_2)}{\theta_1} (f(\bar{x}^t) - f(\bar{y}^t)) - \frac{1}{1 + \eta\sigma} \mathcal{Z}^t - \frac{\theta_2}{\theta_1} (f(\bar{w}^t) - f(\bar{x}^t)) \\
&\quad + \mathbb{E} \left[\mathcal{Z}^{t+1} + \frac{1}{\theta_1} (f(\bar{y}^{t+1}) - f(\bar{x}^t)) \right] \\
&\quad + \langle \nabla f(\bar{x}^t) - \mathbb{E}[\bar{s}^t], x^* - \bar{z}^t \rangle - \frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 - \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2,
\end{aligned}$$

356 where in the second inequality we use the convexity of $f(\cdot)$. The procedure of the algorithm means

$$\begin{aligned}
\mathbf{x}^t &\stackrel{\text{Alg. 2}}{=} \theta_1 \mathbf{z}^t + \theta_2 \mathbf{w}^t + (1 - \theta_1 - \theta_2) \mathbf{y}^t, \\
\mathbf{z}^t - \mathbf{x}^t &\stackrel{\text{Alg. 2}}{=} \frac{\theta_2}{\theta_1} (\mathbf{x}^t - \mathbf{w}^t) + \frac{1 - \theta_1 - \theta_2}{\theta_1} (\mathbf{x}^t - \mathbf{y}^t).
\end{aligned}$$

357 Combining above results, we obtain

$$\begin{aligned}
\mathbb{E} [\mathcal{Z}^{t+1} + \mathcal{Y}^{t+1}] &\leq \frac{1}{1 + \eta\sigma} \mathcal{Z}^t + (1 - \theta_1 - \theta_2) \mathcal{Y}^t + \frac{\theta_2}{\theta_1} (f(\bar{w}^t) - f^*) \\
&\quad - \langle \nabla f(\bar{x}^t) - \mathbb{E}[\bar{s}^t], x^* - \bar{z}^t \rangle + \frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2.
\end{aligned}$$

358 Using definition of \mathcal{W}^t , we get

$$\begin{aligned}
\mathbb{E} [\mathcal{Z}^{t+1} + \mathcal{Y}^{t+1}] &\leq \frac{1}{1 + \eta\sigma} \mathcal{Z}^t + (1 - \theta_1 - \theta_2) \mathcal{Y}^t + p\gamma \mathcal{W}^t \\
&\quad + \|\nabla f(\bar{x}^t) - \mathbb{E}[\bar{s}^t]\| \|x^* - \bar{z}^t\| + \frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \\
&\stackrel{(4)}{\leq} \frac{1}{1 + \eta\sigma} \mathcal{Z}^t + (1 - \theta_1 - \theta_2) \mathcal{Y}^t + p\gamma \mathcal{W}^t \\
&\quad + \frac{L}{\sqrt{m}} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\| \|x^* - \bar{z}^t\| + \frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2.
\end{aligned}$$

359 Finally, we use Lemma 8 to achieve

$$\begin{aligned}
&\mathbb{E} [\mathcal{Z}^{t+1} + \mathcal{Y}^{t+1} + \mathcal{W}^{t+1}] \\
&\leq \frac{1}{1 + \eta\sigma} \mathcal{Z}^t + (1 - \theta_1 - \theta_2) \mathcal{Y}^t + p\gamma \mathcal{W}^t + (1 - p) \mathcal{W}^t + \frac{\theta_2}{\gamma} \mathcal{Y}^t \\
&\quad + \frac{L}{\sqrt{m}} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\| \|x^* - \bar{z}^t\| + \frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \\
&= \frac{1}{1 + \eta\sigma} \mathcal{Z}^t + \left(1 - \left(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma}\right)\right) \mathcal{Y}^t + (1 - p(1 - \gamma)) \mathcal{W}^t \\
&\quad + \frac{L}{\sqrt{m}} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\| \|x^* - \bar{z}^t\| + \frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \\
&\leq \frac{1}{1 + \eta\sigma} \mathcal{Z}^t + \left(1 - \left(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma}\right)\right) \mathcal{Y}^t + (1 - p(1 - \gamma)) \mathcal{W}^t \\
&\quad + \sqrt{\frac{2\eta LV^t}{(1 + \eta\sigma)m}} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\| + \frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2,
\end{aligned}$$

360 where the last inequality is obtained by the definition of V_t . □

361 B Proof of Lemma 2

362 We first provide some lemmas to bound the consensus error.

363 **Lemma 9.** *Letting*

$$r^t = \frac{L}{m} [\|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2, \frac{\eta^2}{L^2} \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2, \frac{\eta^2}{L^2} \|\mathbf{s}^t - \mathbf{1}\bar{s}^t\|^2, \|\mathbf{z}^t - \mathbf{1}\bar{z}^t\|^2, \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2, \|\mathbf{y}^t - \mathbf{1}\bar{y}^t\|^2]^\top,$$

364 *then under the settings of Lemma 2, we have*

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{x}^{t+1} - \mathbf{1}\bar{x}^{t+1}\|^2 \right] &\leq 3\rho^2 \theta_1^2 \|\mathbf{z}^{t+1} - \mathbf{1}\bar{z}^{t+1}\|^2 + 3\rho^2 \theta_2^2 \|\mathbf{w}^{t+1} - \mathbf{1}\bar{w}^{t+1}\|^2 \\
&\quad + 3\rho^2 (1 - \theta_1 - \theta_2)^2 \|\mathbf{y}^{t+1} - \mathbf{1}\bar{y}^{t+1}\|^2, \\
\mathbb{E} \left[\|\mathbf{u}^{t+1} - \mathbf{1}\bar{u}^{t+1}\|^2 \right] &\leq 2\rho^2 \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2 + 2\rho^2 \|\mathbf{g}^{t+1} - \mathbf{g}^t\|^2, \\
\mathbb{E} \left[\|\mathbf{s}^{t+1} - \mathbf{1}\bar{s}^{t+1}\|^2 \right] &\leq 2\rho^2 \|\mathbf{s}^t - \mathbf{1}\bar{s}^t\|^2 + 2\rho^2 \|\mathbf{v}^{t+1} - \mathbf{v}^t\|^2, \\
\mathbb{E} \left[\|\mathbf{z}^{t+1} - \mathbf{1}\bar{z}^{t+1}\|^2 \right] &\leq \frac{3\rho^2 \eta^2 \sigma^2}{(1 + \eta\sigma)^2} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + 3\rho^2 \frac{1}{(1 + \eta\sigma)^2} \|\mathbf{z}^t - \mathbf{1}\bar{z}^t\|^2 \\
&\quad + \frac{3\rho^2 \eta^2}{(1 + \eta\sigma)^2 L^2} \|\mathbf{s}^t - \mathbf{1}\bar{s}^t\|^2, \\
\mathbb{E} \left[\|\mathbf{y}^{t+1} - \mathbf{1}\bar{y}^{t+1}\|^2 \right] &\leq 3\rho^2 \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + 3\rho^2 \theta_1^2 \|\mathbf{z}^{t+1} - \mathbf{1}\bar{z}^{t+1}\|^2 + 3\rho^2 \theta_1^2 \|\mathbf{z}^t - \mathbf{1}\bar{z}^t\|^2,
\end{aligned}$$

365 where $\rho > 0$ is the parameter such that $K = \frac{\log(1/\rho)}{\sqrt{1 - \lambda_2(W)}}$.

366 **Lemma 10.** *Under the settings of Lemma 2, it holds that*

$$\mathbb{E} \left[\left\| \mathbf{w}^{t+1} - \mathbf{1}\bar{w}^{t+1} \right\|^2 \right] \leq \rho^2 p \left\| \mathbf{y}^t - \mathbf{1}\bar{y}^t \right\|^2 + \rho^2 (1-p) \left\| \mathbf{w}^t - \mathbf{1}\bar{w}^t \right\|^2.$$

367 *Proof.* From Lemma 1, we have that

$$\begin{aligned} \mathbb{E} \left[\left\| \mathbf{w}^{t+1} - \mathbf{1}\bar{w}^{t+1} \right\|^2 \right] &\leq \mathbb{E} \left[\rho^2 \left\| \tilde{\mathbf{w}}^{t+1} - \mathbf{1}\tilde{\bar{w}}^{t+1} \right\|^2 \right] \\ &\stackrel{\text{Alg. 2}}{=} \rho^2 p \left\| \mathbf{y}^t - \mathbf{1}\bar{y}^t \right\|^2 + \rho^2 (1-p) \left\| \mathbf{w}^t - \mathbf{1}\bar{w}^t \right\|^2. \end{aligned}$$

368 □

369 Then we provide upper bound for $\left\| \mathbf{g}^{t+1} - \mathbf{g}^t \right\|^2$ and $\left\| \mathbf{v}^{t+1} - \mathbf{v}^t \right\|^2$.

370 **Lemma 11.** *Under the settings of Lemma 2, it holds that*

$$\begin{aligned} \left\| \mathbf{g}^{t+1} - \mathbf{g}^t \right\|^2 &\leq 4L^2 \left\| \mathbf{w}^{t+1} - \mathbf{1}\bar{w}^{t+1} \right\|^2 + \frac{8Lmp\gamma\theta_1}{\theta_2} \mathcal{W}^{t+1} \\ &\quad + 4L^2 \left\| \mathbf{w}^t - \mathbf{1}\bar{w}^t \right\|^2 + \frac{8Lmp\gamma\theta_1}{\theta_2} \mathcal{W}^t. \end{aligned}$$

371 *Proof.* We have

$$\begin{aligned} \left\| \mathbf{g}^{t+1} - \mathbf{g}^t \right\|^2 &= \sum_{i=1}^m \left\| g_i^{t+1} - g_i^t \right\|^2 \\ &\leq 2 \sum_{i=1}^m \left\| \nabla f_i(w_i^{t+1}) - \nabla f_i(x^*) \right\|^2 + 2 \sum_{i=1}^m \left\| \nabla f_i(w_i^t) - \nabla f_i(x^*) \right\|^2. \end{aligned}$$

372 We can also obtain that

$$\begin{aligned} \left\| \nabla f_i(w_i^t) - \nabla f_i(x^*) \right\|^2 &= \left\| \nabla f_i(w_i^t) - \nabla f_i(\bar{w}^t) + \nabla f_i(\bar{w}^t) - \nabla f_i(x^*) \right\|^2 \\ &\leq 2 \left\| \nabla f_i(w_i^t) - \nabla f_i(\bar{w}^t) \right\|^2 + 2 \left\| \nabla f_i(\bar{w}^t) - \nabla f_i(x^*) \right\|^2 \\ &\stackrel{(1)}{\leq} 2L^2 \left\| w_i^t - \bar{w}^t \right\|^2 + 2 \left\| \nabla f_i(\bar{w}^t) - \nabla f_i(x^*) \right\|^2 \\ &\stackrel{(2)}{\leq} 2L^2 \left\| w_i^t - \bar{w}^t \right\|^2 + 4L(f_i(\bar{w}^t) - f_i(x^*)). \end{aligned}$$

373 Combining above results, we achieve

$$\begin{aligned} \left\| \mathbf{g}^{t+1} - \mathbf{g}^t \right\|^2 &\leq 2 \sum_{i=1}^m \left\| \nabla f_i(w_i^{t+1}) - \nabla f_i(x^*) \right\|^2 + 2 \sum_{i=1}^m \left\| \nabla f_i(w_i^t) - \nabla f_i(x^*) \right\|^2 \\ &\leq \sum_{i=1}^m 4L^2 \left\| w_i^{t+1} - \bar{w}^{t+1} \right\|^2 + 8L(f_i(\bar{w}^{t+1}) - f_i(x^*)) \\ &\quad + \sum_{i=1}^m 4L^2 \left\| w_i^t - \bar{w}^t \right\|^2 + 8L(f_i(\bar{w}^t) - f_i(x^*)) \\ &= 4L^2 \left\| \mathbf{w}^{t+1} - \mathbf{1}\bar{w}^{t+1} \right\|^2 + 8Lm(f(\bar{w}^{t+1}) - f(x^*)) \\ &\quad + 4L^2 \left\| \mathbf{w}^t - \mathbf{1}\bar{w}^t \right\|^2 + 8Lm(f(\bar{w}^t) - f(x^*)) \\ &= 4L^2 \left\| \mathbf{w}^{t+1} - \mathbf{1}\bar{w}^{t+1} \right\|^2 + \frac{8Lmp\gamma\theta_1}{\theta_2} \mathcal{W}^{t+1} \\ &\quad + 4L^2 \left\| \mathbf{w}^t - \mathbf{1}\bar{w}^t \right\|^2 + \frac{8Lmp\gamma\theta_1}{\theta_2} \mathcal{W}^t. \end{aligned}$$

374 □

375 Next, we target to bound $\|\mathbf{v}^t - \mathbf{1}\bar{v}^t\|^2$. We first give two auxiliary lemmas.

376 **Lemma 12.** *Under the settings of Lemma 2, it holds that*

$$\sum_{i=1}^m \|u_i^t\|^2 \leq 3 \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2 + 3L^2 \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 + \frac{6Lmp\gamma\theta_1}{\theta_2} \mathcal{W}^t.$$

377 *Proof.* We have

$$\begin{aligned} \sum_{i=1}^m \|u_i^t\|^2 &= \sum_{i=1}^m \|(u_i^t - \bar{u}^t) + (\bar{u}^t - \nabla f(\bar{w}^t)) + (\nabla f(\bar{w}^t) - \nabla f(x^*))\|^2 \\ &\leq \sum_{i=1}^m \left[3 \|u_i^t - \bar{u}^t\|^2 + 3 \|\bar{u}^t - \nabla f(\bar{w}^t)\|^2 + 3 \|\nabla f(\bar{w}^t) - \nabla f(x^*)\|^2 \right] \\ &\stackrel{(A.1)}{\leq} 3 \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2 + 3 \sum_{i=1}^m \left\| \frac{1}{m} \sum_{j=1}^m (\nabla f_j(w_j^t) - \nabla f_j(\bar{w}^t)) \right\|^2 + 3 \sum_{i=1}^m (2L(f_i(\bar{w}^t) - f_i(x^*))) \\ &\leq 3 \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2 + \frac{3}{m} \sum_{i=1}^m \sum_{j=1}^m \|\nabla f_j(w_j^t) - \nabla f_j(\bar{w}^t)\|^2 + 3 \sum_{i=1}^m (2L(f_i(\bar{w}^t) - f_i(x^*))) \\ &\stackrel{\text{Asm. 1}}{\leq} 3 \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2 + \frac{3L^2}{m} \sum_{i=1}^m \sum_{j=1}^m \|w_j^t - \bar{w}^t\|^2 + 3 \sum_{i=1}^m (2L(f_i(\bar{w}^t) - f_i(x^*))) \\ &= 3 \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2 + 3L^2 \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 + \frac{6Lmp\gamma\theta_1}{\theta_2} \mathcal{W}^t. \end{aligned}$$

378

□

379 **Lemma 13.** *Under the settings of Lemma 2, it holds that*

$$\begin{aligned} &\sum_{i=1}^m \mathbb{E} \left[\left\| \frac{\xi_i}{q} (\nabla f_i(x_i^t) - \nabla f_i(w_i^t)) \right\|^2 \right] \\ &\leq \frac{8\eta mL\theta_1}{(1+\eta\sigma)q} \mathcal{Z}^t + \frac{8mLp\gamma\theta_1}{q\theta_2} (1+\theta_2) \mathcal{W}^t + \frac{8mL\theta_1}{q} (1-\theta_1-\theta_2) \mathcal{Y}^t \\ &\quad + \frac{4L^2}{q} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{4L^2}{q} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2. \end{aligned}$$

380 *Proof.* We have

$$\begin{aligned} &\sum_{i=1}^m \mathbb{E} \left[\left\| \frac{\xi_i}{q} (\nabla f_i(x_i^t) - \nabla f_i(w_i^t)) \right\|^2 \right] \\ &= \frac{1}{q} \sum_{i=1}^m \|\nabla f_i(x_i^t) - \nabla f_i(w_i^t)\|^2 \\ &\leq \frac{4}{q} \sum_{i=1}^m \left[\|\nabla f_i(x_i^t) - \nabla f_i(\bar{x}^t)\|^2 + \|\nabla f_i(\bar{x}^t) - \nabla f_i(x^*)\|^2 + \|\nabla f_i(x^*) - \nabla f_i(\bar{w}^t)\|^2 + \|\nabla f_i(\bar{w}^t) - \nabla f_i(w_i^t)\|^2 \right] \\ &\stackrel{\text{Asm. 1}}{\leq} \frac{4L^2}{q} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{8mL}{q} (f(\bar{x}^t) - f(x^*)) + \frac{8mL}{q} (f(\bar{w}^t) - f(x^*)) + \frac{4L^2}{q} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \\ &\stackrel{\text{Asm. 2}}{\leq} \frac{8mL}{q} (\theta_1(f(\bar{z}^t) - f(x^*)) + \theta_2(f(\bar{w}^t) - f(x^*)) + (1-\theta_1-\theta_2)(f(\bar{z}^t) - f(x^*))) \\ &\quad + \frac{4L^2}{q} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{8mL}{q} (f(\bar{w}^t) - f(x^*)) + \frac{4L^2}{q} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \\ &\stackrel{\text{Asm. 1}}{\leq} \frac{4mL^2\theta_1}{q} \|\bar{z}^t - x^*\|^2 + \frac{8mL}{q} (1+\theta_2)(f(\bar{w}^t) - f(x^*)) + \frac{8mL}{q} (1-\theta_1-\theta_2)(f(\bar{y}^t) - f(x^*)) \\ &\quad + \frac{4L^2}{q} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{4L^2}{q} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \end{aligned}$$

$$\begin{aligned}
&= \frac{8\eta mL\theta_1}{(1+\eta\sigma)q} \mathcal{Z}^t + \frac{8mLp\gamma\theta_1}{q\theta_2} (1+\theta_2)\mathcal{W}^t + \frac{8mL\theta_1}{q} (1-\theta_1-\theta_2)\mathcal{Y}^t \\
&\quad + \frac{4L^2}{q} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{4L^2}{q} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2.
\end{aligned}$$

381

□

382 Now we are ready to bound $\|\mathbf{v}^t - \mathbf{1}\bar{v}^t\|^2$.

383 **Lemma 14.** *Under the settings of Lemma 2, it holds that*

$$\begin{aligned}
&\mathbb{E} \left[\|\mathbf{v}^{t+1} - \mathbf{v}^t\|^2 \right] \\
&\leq \frac{4Lmp\gamma\theta_1}{\theta_2} \left(\frac{4(1+\theta_2)}{q} + 3 \right) (\mathcal{W}^{t+1} + \mathcal{W}^t) + \frac{16\eta mL\theta_1}{(1+\eta\sigma)q} (\mathcal{Z}^{t+1} + \mathcal{Z}^t) \\
&\quad + \frac{16mL\theta_1}{q} (1-\theta_1-\theta_2)(\mathcal{Y}^{t+1} + \mathcal{Y}^t) \\
&\quad + 6 \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2 + \frac{8L^2}{q} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + 2L^2 \left(\frac{4}{q} + 3 \right) \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \\
&\quad + 6 \|\mathbf{u}^{t+1} - \mathbf{1}\bar{u}^{t+1}\|^2 + \frac{8L^2}{q} \|\mathbf{x}^{t+1} - \mathbf{1}\bar{x}^{t+1}\|^2 + 2L^2 \left(\frac{4}{q} + 3 \right) \|\mathbf{w}^{t+1} - \mathbf{1}\bar{w}^{t+1}\|^2.
\end{aligned}$$

384 *Proof.* It holds that

$$\begin{aligned}
\mathbb{E} \left[\|\mathbf{v}^t\|^2 \right] &= \sum_{i=1}^m \mathbb{E} \left[\left\| u_i^t + \frac{\xi_i}{q} (\nabla f_i(x_i^t) - \nabla f_i(w_i^t)) \right\|^2 \right] \\
&= \sum_{i=1}^m \mathbb{E} \left[\|u_i^t\|^2 + \left\| \frac{\xi_i}{q} (\nabla f_i(x_i^t) - \nabla f_i(w_i^t)) \right\|^2 \right] \\
&\stackrel{(12),(13)}{\leq} \frac{2Lmp\gamma\theta_1}{\theta_2} \left(\frac{4(1+\theta_2)}{q} + 3 \right) \mathcal{W}^t + \frac{8\eta mL\theta_1}{(1+\eta\sigma)q} \mathcal{Z}^t + \frac{8mL\theta_1}{q} (1-\theta_1-\theta_2)\mathcal{Y}^t \\
&\quad + 3 \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2 + \frac{4L^2}{q} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + L^2 \left(\frac{4}{q} + 3 \right) \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2.
\end{aligned}$$

385 Then we use the fact that $\|\mathbf{v}^{t+1} - \mathbf{v}^t\|^2 \leq 2\|\mathbf{v}^{t+1}\|^2 + 2\|\mathbf{v}^t\|^2$, we can obtain the result. □

386 Substituting the result of Lemma 11 and Lemma 14 into Lemma 9, we obtain Lemma 2. Here, we
387 rewrite the result of Lemma 2 by taking $q = b/m$, which contains the detailed expressions of \mathbf{A} , \mathbf{B}
388 and e^t

Lemma 15 (The complete version of Lemma 2). *Let*

$$r^t = \frac{L}{m} [\|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2, \frac{\eta^2}{L^2} \|\mathbf{u}^t - \mathbf{1}\bar{u}^t\|^2, \frac{\eta^2}{L^2} \|\mathbf{s}^t - \mathbf{1}\bar{s}^t\|^2, \|\mathbf{z}^t - \mathbf{1}\bar{z}^t\|^2, \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2, \|\mathbf{y}^t - \mathbf{1}\bar{y}^t\|^2]^\top.$$

389 *Under the settings of Lemma 1, we run Algorithm 2 by taking*

$$K = \left\lceil \frac{\log(1/\rho)}{\sqrt{1 - \lambda_2(W)}} \right\rceil$$

390 *with*

$$\rho \leq \frac{q\theta_1^3}{9} \min \left\{ \frac{1}{2}\eta^2\sigma^2, \frac{1}{2} \left(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma} \right), \frac{1}{2}p(1-\gamma) \right\},$$

391 *then we have*

$$\mathbb{E} [r^{t+1}] \leq \rho^2 \cdot \left(\mathbf{B} \cdot \mathbf{A} \cdot r^t + e^t \right),$$

392 where we define matrix \mathbf{A} , elementary matrix \mathbf{B} and vector e^t as

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 8\eta^2(1 + \rho^2(1 - p)) & 8\eta^2\rho^2p \\ a_{31} & 12(1 + 2\rho^2) & 2 & a_{34} & a_{35} & a_{36} \\ \frac{3\eta^2\sigma^2}{(1+\eta\sigma)^2} & 0 & \frac{3}{(1+\eta\sigma)^2} & \frac{3}{(1+\eta\sigma)^2} & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 - p & p \\ \frac{9\rho^2\theta_1^2\eta^2\sigma^2}{(1+\eta\sigma)^2} & 0 & \frac{9\rho^2\theta_1^2}{(1+\eta\sigma)^2} & 3\theta_1^2 + \frac{9\rho^2\theta_1^2}{(1+\eta\sigma)^2} & 0 & 0 \end{bmatrix},$$

393

$$\mathbf{B} = \begin{bmatrix} 1 & 0 & 0 & 3\rho^2\theta_1^2 & 3\rho^2\theta_2^2 & 3\rho^2(1 - \theta_1 - \theta_2)^2 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

394 and

$$e^t = \begin{bmatrix} 0 \\ \frac{16\eta^2p\gamma\theta_1}{\theta_2}(\mathcal{W}^{t+1} + \mathcal{W}^t) \\ \frac{32\eta^3\theta_1m}{(1+\eta\sigma)b}(\mathcal{Z}^t + \mathcal{Z}^{t+1}) + \frac{8\eta^2p\gamma\theta_1}{\theta_2}\left(\frac{4m(1+\theta_2)}{b}\right) + 3 + 48\rho^2(\mathcal{W}^t + \mathcal{W}^{t+1}) + \frac{32\eta^3\theta_1m}{b}(1 - \theta_1 - \theta_2)(\mathcal{Y}^t + \mathcal{Y}^{t+1}) \\ 0 \\ 0 \\ 0 \end{bmatrix},$$

395 where

$$\begin{aligned} a_{31} &= \frac{16\eta^2m}{b} \left(1 + 9\rho^4 \left[\frac{\eta^2\sigma^2\theta_1^2}{(1+\eta\sigma)^2} + (1 - \theta_1 - \theta_2)^2 \right] \right), \\ a_{34} &= \frac{144\eta^2\rho^4m}{b} \left(\frac{\theta_1^2}{(1+\eta\sigma)^2} + \theta_1^2(1 - \theta_1 - \theta_2)^2 \right), \\ a_{35} &= \frac{16\eta^2m}{b} [1 + \rho^2(1 - p) + 9\rho^4\theta_2^2(1 - p)] + 12\eta^2(1 + 8\rho^2)(1 + \rho^2(1 - p)), \\ a_{36} &= 4\eta^2\rho^2p \left[\frac{4m}{b}(1 + 3\rho^2\theta_2^2) + 3(1 + 8\rho^2) \right]. \end{aligned}$$

396 to simplify our equation. Then we can do a simple estimation with substituting $\eta = 1/(13\theta_1)$ to
397 obtain that

$$\|\mathbf{A}\| < \frac{4m}{b\theta_1^2}, \quad \|\mathbf{B}\| < 2 \quad \text{and} \quad \|e^t\| < \frac{2m}{3b\theta_1^2}(V^{t+1} + V^t).$$

398 **C Proof of Theorem 1**

399 Combining Lemma 1 and 2, we have got what we need to prove our main theorem.

400 *Proof for Theorem 1.* We prove the theorem by induction. By Lemma 1, the theorem holds for $t = 1$.
401 Now, we assume that

$$\mathbb{E}[V^t] \leq \left(\underbrace{\max \left(1 - \eta\sigma, 1 - \frac{1}{2} \left(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma} \right), 1 - \frac{1}{2}p(1 - \gamma) \right)}_{\alpha} \right)^t (V^0 + \|r^0\|). \quad (\text{C.1})$$

402 holds when $t \leq k$ and are going to prove it holds for $t = k + 1$. We use the notation \mathbf{A} , \mathbf{B} and e^t by
 403 following Lemma 2 (Section B). By the definition of V^t and Lemma 15, we can obtain that

$$\begin{aligned}
 & \mathbb{E} [\|r^k\|] \\
 & \leq \frac{2\rho^2 m}{3b\theta_1^2} \cdot \sum_{i=1}^k (2\rho^2 \|\mathbf{A}\|)^{k-i} (V^i + V^{i-1}) + (2\rho^2 \|\mathbf{A}\|)^k \|r^0\| \\
 & \leq \frac{2\rho^2 m}{3b\theta_1^2} \cdot \sum_{i=1}^k (2\rho^2 \|\mathbf{A}\|)^{k-i} \alpha^{i-1} (\alpha + 1) (V^0 + \|r^0\|) + (2\rho^2 \|\mathbf{A}\|)^k \|r^0\| \\
 & \leq \frac{2\rho^2 m (\alpha + 1)}{3b\theta_1^2 \alpha} \cdot \left(\frac{\alpha}{2}\right)^k \sum_{i=1}^k \left(\frac{\alpha}{2}\right)^{-i} \alpha^i (V^0 + \|r^0\|) + (2\rho^2 \|\mathbf{A}\|)^k \|r^0\| \\
 & = \frac{2\rho^2 m (\alpha + 1)}{3b\theta_1^2 \alpha} \cdot \left(\frac{\alpha}{2}\right)^k \cdot (2^{k+1} - 2) (V^0 + \|r^0\|) + (2\rho^2 \|\mathbf{A}\|)^k \|r^0\| \\
 & \leq \left(\frac{8\rho^2 m}{3b\theta_1^2} \alpha^{k-1} + (2\rho^2 \|\mathbf{A}\|)^k \right) \cdot (V^0 + \|r^0\|),
 \end{aligned} \tag{C.2}$$

404 where the second inequality is because the condition for ρ in Theorem 1 implies $\rho^2 \leq \alpha/(4\|\mathbf{A}\|)$
 405 when $\eta = 1/(13\theta_1)$ and we assume that for all $i \leq k$, $\mathbb{E}[V^i] \leq \alpha^i (V^0 + \|r^0\|)$. Furthermore, we
 406 can obtain

$$\begin{aligned}
 & \mathbb{E} \left[\frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \right] \\
 & \leq \frac{1}{3b\theta_1} \mathbb{E} [\|r^k\|] \\
 & \leq \frac{1}{3b\theta_1} \left(\frac{8\rho^2 m}{3b\theta_1^2} \alpha^{k-1} + (2\rho^2 \|\mathbf{A}\|)^k \right) \cdot (V^0 + \|r^0\|)
 \end{aligned} \tag{C.3}$$

407 and

$$\begin{aligned}
 & \mathbb{E} \left[\sqrt{\frac{2\eta LV^t}{(1 + \eta\sigma)m}} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\| \right] \\
 & \leq \sqrt{2\eta V^t} \sqrt{\mathbb{E} \|r^k\|} \\
 & \leq \sqrt{2\eta V^t} \sqrt{\left(\frac{8\rho^2 m}{3b\theta_1^2} \alpha^{k-1} + (2\rho^2 \|\mathbf{A}\|)^k \right) \cdot (V^0 + \|r^0\|)}.
 \end{aligned} \tag{C.4}$$

408 Furthermore, if we denote that

$$\beta = \max \left(\frac{1}{1 + \eta\sigma}, 1 - \left(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma} \right), 1 - p(1 - \gamma) \right),$$

409 we have

$$\begin{aligned}
 & \mathbb{E} [V^{k+1}] \\
 & \leq \beta V^k + \sqrt{\frac{2\eta LV^t}{(1 + \eta\sigma)m}} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\| + \frac{L}{3mb\theta_1} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 + \frac{L}{4mb\theta_1} \|\mathbf{w}^t - \mathbf{1}\bar{w}^t\|^2 \\
 & \leq \alpha^k (V^0 + \|r^0\|) \left(\beta + \frac{1}{3b\theta_1} \left(\frac{8\rho^2 m}{3b\theta_1^2} \alpha^{-1} + (2\rho^2 \|\mathbf{A}\| \alpha^{-1})^k \right) \right) \\
 & \quad + \sqrt{\frac{2\eta \alpha^{-k}}{3b\theta_1} \left(\frac{8\rho^2 m}{3b\theta_1^2} \alpha^{k-1} + (2\rho^2 \|\mathbf{A}\|)^k \right)} \\
 & \leq \alpha^k (V^0 + \|r^0\|) \left(\beta + \frac{1}{3b\theta_1} \left(\frac{16\rho^2 m}{3b\theta_1^2} + \left(\frac{16\rho^2 m}{b\theta_1^2} \right)^k \right) \right)
 \end{aligned} \tag{C.5}$$

$$\begin{aligned}
& + \sqrt{\frac{2\eta}{3b\theta_1}} \left(\sqrt{\frac{8\rho^2 m}{3b\theta_1^2}} \alpha^{-1/2} + \left(\frac{16\rho^2 m}{b\theta_1^2} \right)^{\frac{k}{2}} \right) \tag{C.6} \\
& \leq \alpha^k (V^0 + \|r^0\|) \left(\beta + \frac{1}{3b\theta_1} \left(\frac{16m}{3b\theta_1^2} + \frac{16m}{b\theta_1^2} \right) \rho^2 + \sqrt{\frac{2\eta}{3b\theta_1}} \left(2\sqrt{\frac{8m}{3b\theta_1^2}} + \sqrt{\frac{16m}{b\theta_1^2}} \right) \rho \right) \\
& = \alpha^k (V^0 + \|r^0\|) \left(\beta + \frac{64m}{9b^2\theta_1^3} \rho^2 + 4 \left(1 + \sqrt{\frac{2}{3}} \right) \sqrt{\frac{2}{39}} \frac{\sqrt{m}}{b\theta_1^2} \rho \right) \\
& \leq \alpha^k (V^0 + \|r^0\|) \left(\max \left(\frac{1}{1 + \eta\sigma}, 1 - \left(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma} \right), 1 - p(1 - \gamma) \right) + \frac{9m}{b\theta_1^3} \rho \right) \\
& \leq \alpha^{k+1} (V^0 + \|r^0\|),
\end{aligned}$$

410 where the last inequality is because of the condition of ρ in Theorem 1. The first inequality is because
411 of Lemma 1. Inequality (C.5) is because of the equation (C.1),(C.3) and (C.4) and inequality (C.6)
412 is because of $\alpha \geq 1/2$ and $\|\mathbf{A}\| < (4m)/(b\theta_1^2)$. Thus, the theorem also holds for $t = k + 1$ and we
413 complete the proof by induction. Furthermore, Equation C.2 and condition of ρ in Theorem 1 imply
414 that $\rho^2 \leq \alpha/(4\|\mathbf{A}\|)$. Then we obtain

$$\begin{aligned}
& \mathbb{E} \left[\frac{L}{m} \|\mathbf{x}^t - \mathbf{1}\bar{x}^t\|^2 \right] \\
& \leq \left(\frac{8}{243} + 2^{-t} \right) \max \left(1 - \eta\sigma, 1 - \frac{1}{2} \left(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma} \right), 1 - \frac{1}{2} p(1 - \gamma) \right)^t \cdot (V^0 + \|\mathbf{r}^0\|).
\end{aligned}$$

415

□

416 D Proof of Corollary 2

417 *Proof.* We first prove that KNOT (Algorithm 2) can find an ϵ -suboptimal solution in expectation.

418 We run KNOT with the setting of Theorem 1 and let

$$T = \mathcal{O} \left(\left(\frac{1}{\eta\sigma} + \frac{2}{(\theta_1 + \theta_2 - \frac{\theta_2}{\gamma})} + \frac{2}{p(1 - \gamma)} \right) \log \frac{1}{\epsilon} \right).$$

419 Then Theorem 1 means

$$\frac{L}{2} \mathbb{E} \left[\|\bar{z}^T - x^*\|^2 \right] \leq \frac{\epsilon}{3}, \quad \mathbb{E}[f(\bar{y}^T) - f(x^*)] \leq \frac{\epsilon}{3}, \quad \mathbb{E}[f(\bar{w}^T) - f(x^*)] \leq \frac{\epsilon}{3}$$

420 and $\mathbb{E} \left[\|x_i^T - \bar{x}^T\|^2 \right] \leq m\epsilon/L$ for each $i = 1, \dots, m$. From Assumption 1 and 3 we obtain that

$$\begin{aligned}
& \mathbb{E} [f(\bar{x}^T) - f(x^*)] \\
& \leq \theta_1 \mathbb{E} [f(\bar{z}^T) - f(x^*)] + \theta_2 \mathbb{E} [f(\bar{w}^T) - f(x^*)] + (1 - \theta_1 - \theta_2) \mathbb{E} [f(\bar{y}^T) - f(x^*)] \\
& \leq \frac{\theta_1 L}{2} \mathbb{E} \left[\|\bar{z}^T - x^*\|^2 \right] + \theta_2 \mathbb{E} [f(\bar{w}^T) - f(x^*)] + (1 - \theta_1 - \theta_2) \mathbb{E} [f(\bar{y}^T) - f(x^*)] \tag{D.1} \\
& \leq (\theta_1 + \theta_2 + (1 - \theta_1 - \theta_2)) \epsilon = \epsilon.
\end{aligned}$$

421 Moreover, Proposition 1 means step $\mathbf{x}_{\text{out}} = \text{AccGossip}(\mathbf{x}_T, K_{\text{out}})$ with $K_{\text{out}} = \mathcal{O}(\sqrt{1/\alpha} \log m)$
422 (line 23 of Algorithm 2) leads to $\bar{x}^{\text{out}} = \bar{x}^T$ and

$$L \mathbb{E} \left[\|x_i^{\text{out}} - \bar{x}^{\text{out}}\|^2 \right] \leq \frac{\epsilon}{3} \tag{D.2}$$

423 for each $i = 1, \dots, m$. Applying the smoothness of f (Assumption 1) and Young's inequality, we
 424 have

$$\begin{aligned} f(x_i^{\text{out}}) - f(\bar{x}^{\text{out}}) &\leq \langle \nabla f(\bar{x}^{\text{out}}), x_i^{\text{out}} - \bar{x}^{\text{out}} \rangle + \frac{L}{2} \|x_i^{\text{out}} - \bar{x}^{\text{out}}\|^2 \\ &\leq \frac{c}{2} \|\nabla f(\bar{x}^{\text{out}})\|^2 + \frac{1}{2c} \|x_i^{\text{out}} - \bar{x}^{\text{out}}\|^2 + \frac{L}{2} \|x_i^{\text{out}} - \bar{x}^{\text{out}}\|^2. \end{aligned}$$

425 for any $i = 1, \dots, m$ and any $c > 0$. Additionally, Lemma 3 implies

$$f(\bar{x}^{\text{out}}) - f(x^*) \geq \frac{1}{2L} \|\nabla f(\bar{x}^{\text{out}})\|^2.$$

426 Combing above two results with $c = 1/(2L)$, we have

$$\begin{aligned} \mathbb{E} [f(x_i^{\text{out}}) - f(x^*)] &\leq cL \mathbb{E} [f(\bar{x}^{\text{out}}) - f(x^*)] + \left(\frac{1}{2c} + \frac{L}{2} \right) \mathbb{E} [\|x_i^{\text{out}} - \bar{x}^{\text{out}}\|^2] \\ &\leq \frac{\epsilon}{2} + \frac{3}{2} \cdot \frac{\epsilon}{3} = \epsilon \end{aligned}$$

427 for any $i = 1, \dots, m$. This implies the output x^{out} is an ϵ -suboptimal solution in expectation.

428 Then we analyze the complexity of KNOT by following the parameter settings of Theorem 1.

429 **Case 1:** In the case of $m < \kappa$, we choose $b = \sqrt{m}$ and $p = 1/\sqrt{m}$. Thus,

$$\theta_1 = \min \left\{ \sqrt{\frac{b}{\kappa p}} \theta_2, \theta_2 \right\} = \sqrt{\frac{b}{\kappa p}} \theta_2 = \frac{1}{2\sqrt{\kappa}}.$$

430 By choosing $\gamma = 1 - \frac{1}{3}\sqrt{\frac{m}{\kappa}} \in (2/3, 1)$, we have

$$\frac{2}{\theta_1 + \theta_2 - \frac{\theta_2}{\gamma}} = \frac{2}{\theta_1(1 - \frac{1}{3\gamma})} \leq \frac{4}{\theta_1} = 2\sqrt{\kappa},$$

431 which means

$$\frac{1}{\eta\sigma} + \frac{2}{\theta_1 + \theta_2 - \frac{\theta_2}{\gamma}} + \frac{2}{p(1 - \gamma)} \leq \left(\frac{13}{2} + 4 + \frac{4}{1 - \gamma} \right) \frac{1}{\theta_1} = \mathcal{O}(\sqrt{\kappa}).$$

432 **Case 2:** In the case of $m \geq \kappa$, we choose $b = \sqrt{\kappa}$ and $p = 1/\sqrt{\kappa}$. Thus

$$\theta_1 = \min \left\{ \sqrt{\frac{b}{\kappa p}} \theta_2, \theta_2 \right\} = \theta_2 = \frac{1}{2\sqrt{\kappa}}.$$

433 By choosing $\gamma \in (2/3, 1)$, we have

$$\frac{1}{\eta\sigma} + \frac{2}{\theta_1 + \theta_2 - \frac{\theta_2}{\gamma}} + \frac{2}{p(1 - \gamma)} \leq \left(\frac{13}{2} + 4 + \frac{4}{1 - \gamma} \right) \frac{1}{\theta_1} = \mathcal{O}(\sqrt{\kappa}).$$

434 Therefore, the number of iterations for KNOT to achieve an expected ϵ -suboptimal solution is

$$T = \mathcal{O} \left(\sqrt{\kappa} \log \frac{1}{\epsilon} \right).$$

435 Thus, the expected first-order oracle complexity is

$$T \cdot (b + mp) = \mathcal{O} \left((m + \sqrt{m\kappa}) \log \frac{1}{\epsilon} \right),$$

436 and the expected communication complexity is

$$T \cdot K = \mathcal{O} \left(\frac{\sqrt{\kappa} \log(m\kappa)}{\sqrt{1 - \lambda_2(W)}} \log \frac{1}{\epsilon} \right).$$

437

□