# Stochastic Distributed Optimization under Average Second-order Similarity: Algorithms and Analysis

**Dachao Lin** [*†]     **Yuze Han** [*‡]     **Haishan Ye** [§]     **Zhihua Zhang** [¶]

## Abstract

We study finite-sum distributed optimization problems involving a master node and $n-1$ local nodes under the popular $\delta$-similarity and $\mu$-strong convexity conditions. We propose two new algorithms, SVRS and AccSVRS, motivated by previous works. The non-accelerated SVRS method combines the techniques of gradient sliding and variance reduction and achieves a better communication complexity of $\tilde{\mathcal{O}}(n+\sqrt{n}\delta/\mu)$ compared to existing non-accelerated algorithms. Applying the framework proposed in Katyusha X [6], we also develop a directly accelerated version named AccSVRS with the $\tilde{\mathcal{O}}(n+n^{3/4}\sqrt{\delta/\mu})$ communication complexity. In contrast to existing results, our complexity bounds are entirely smoothness-free and exhibit superiority in ill-conditioned cases. Furthermore, we establish a nearly matched lower bound to verify the tightness of our AccSVRS method.

## 1 Introduction

We have witnessed the development of distributed optimization in recent years. Distributed optimization aims to cooperatively solve a learning task over a predefined social network by exchanging information exclusively with immediate neighbors. This class of problems has found extensive applications in various fields, including machine learning, healthcare, network information processing, telecommunications, manufacturing, natural language processing tasks, and multi-agent control [54, 30, 48, 45, 60, 8]. In this paper, we focus on the following classical finite-sum optimization problem in a centralized setting:

$$\min_{\boldsymbol{x}\in\mathbb{R}^d} f(\boldsymbol{x}) := \frac{1}{n}\sum_{i=1}^n f_i(\boldsymbol{x}), \tag{1}$$

where each $f_i$ is differentiable and corresponds to a client or node, and the target objective is their average function $f$. Without loss of generality, we assume $f_1$ is the master node and the others are local nodes. In each round, every local node can communicate with the master node certain information, such as the local parameter $\boldsymbol{x}$, local gradient $\nabla f_i(\boldsymbol{x})$, and some global information gathered at the master node. Such a scheme can also be viewed as decentralized optimization over a star network [55].

Following the wisdom of statistical similarity residing in the data at different nodes, many previous works study scenarios where the individual functions exhibit relationships or, more specifically, certain homogeneity shared among the local $f_i$'s and $f$. The most common one is under the $\delta$-second-order

---

[*]Equal Contribution.

[†]Academy for Advanced Interdisciplinary Studies; Peking University; `lindachao@pku.edu.cn`;

[‡]School of Mathematical Sciences; Peking University; `hanyuze97@pku.edu.cn`;

[§]Corresponding Author; School of Management; Xi'an Jiaotong University; SGIT AI Lab, State Grid Corporation of China; `yehaishan@xjtu.edu.cn`;

[¶]School of Mathematical Sciences; Peking University; `zhzhang@math.pku.edu.cn`.

similarity assumption [33, 35], that is,

$$\left\| \nabla^2 f_i(\boldsymbol{x}) - \nabla^2 f(\boldsymbol{x}) \right\| \leq \delta, \forall \boldsymbol{x} \in \mathbb{R}^d, i \in \{1, \ldots, n\}.$$

Such an assumption also has different names in the literature, such as $\delta$-related assumption, bounded Hessian dissimilarity, or function similarity [7, 32, 51, 54, 62]. The rigorous definitions are deferred to Section 2. Moreover, the second-order similarity assumption can hold with a relatively small $\delta$ compared to the smoothness coefficient of $f_i$'s in many practical settings, such as statistical learning. More insights on this can be found in the discussion presented in [54, Section 2]. The similarity assumption indicates that the data across different clients share common information on the second-order derivative, potentially leading to a reduction in communication among clients. Meanwhile, the cost of communication is often much higher than that of local computation in distributed optimization settings [9, 44, 30]. Hence, researchers are motivated to develop efficient algorithms characterized by low communication complexity, which is the primary objective of this paper as well.

Furthermore, we need to emphasize that prior research [25, 56, 63, 19, 43, 5] has shown tightly matched lower and upper bounds on computation complexity for the finite-sum objective in Eq. (1). These works focus on gradient complexity under (average) smoothness [63] instead of communication complexity under similarity. Indeed, we will also discuss and compare the gradient complexity as shown in [35], to explore the trade-off between communication and gradient complexity.

Although the development of distributed optimization with similarity has lasted for years, the optimal complexity under full participation was only recently achieved by Kovalev et al. [35]. They employed gradient-sliding [37] and obtained the optimal communication complexity $\tilde{\mathcal{O}}(n\sqrt{\delta/\mu})$ for $\mu$-strongly convex $f$ and $\delta$-related $f_i$'s in Eq. (1). However, the full participation model requires the calculation of the whole gradient $\nabla f(\cdot)$, which incurs a communication cost of $n-1$ in each round. In contrast, partial participation could reduce the communication burden and yield improved complexity. Hence, Khaled and Jin [33] introduced client sampling, a technique that selects one client for updating in each round. They developed a non-accelerated algorithm SVRP, which achieves the communication complexity of $\tilde{\mathcal{O}}(n+\delta^2/\mu^2)$. Additionally, they proposed a Catalyzed version of SVRP with the complexity $\tilde{\mathcal{O}}(n+n^{3/4}\sqrt{\delta/\mu})$, which is better than the rates obtained in the full participation setting.

We believe there are several potential avenues for improvement inspired by [33]. 1) Khaled and Jin [33] introduced the requirement that each individual function is strongly convex (see [33, Assumption 2]). However, this constraint is absent in prior works. Notably, in the context of full participation, even non-convexity is deemed acceptable[6]. A prominent example is the shift-and-invert approach to solving PCA [52, 23], where each component is smooth and non-convex, but the average function remains convex. Thus we doubt the necessity of requiring strong convexity for individual components. 2) In hindsight, it seems that the directly accelerated SVRP could only achieve a bound of $\tilde{\mathcal{O}}(n + \sqrt{n} \cdot \delta/\mu)$ based on the current analysis, which is far from being satisfactory compared to its Catalyzed version. Consequently, there might be room for the development of a more effective algorithm for direct acceleration. 3) It is essential to note that the Catalyst framework introduces an additional log term in the overall complexity, along with the challenge of parameter tuning. This aspect is discussed in detail in [6, Section 1.2]. Therefore, we intend to address the aforementioned concerns, particularly on designing directly accelerated methods under the second-order similarity assumption.

## 1.1 Main Contributions

In this paper, we address the above concerns under the average similarity condition. Our contributions are presented in detail below and we provide a comparison with previous works in Table 1:

- First, we combine gradient sliding and client sampling techniques to develop an improved non-accelerated algorithm named SVRS (Algorithms 1). SVRS achieves a communication complexity of $\tilde{\mathcal{O}}(n + \sqrt{n} \cdot \delta/\mu)$, surpassing SVRP in ill-conditioned cases. Notably, this rate does not need component strong convexity and applies to the function value gap instead of the parameter distance.

- Second, building on SVRS, we employ a classical interpolation framework motivated by Katyusha X [6] to introduce the directly accelerated SVRS (AccSVRS, Algorithm 2).

---

[6]Readers can check that the proof of [35] only requires $f_1(\cdot) + \frac{1}{2\theta} \|\cdot\|^2$ is strongly convex, which can be guaranteed by $\delta$-second-order similarity since $f$ is $\mu$-strongly convex and $\theta = 1/(2\delta)$ therein.

Table 1: Comparison of communication under similarity for the strongly convex objective.

| | Method/Reference | Communication complexity | Assumptions |
|---|---|---|---|
| No Sampling | AccExtragradient [35] | $\mathcal{O}\left(n\sqrt{\frac{\delta}{\mu}}\log\frac{1}{\varepsilon}\right)$ | SS only for $f_1$ |
| | Lower bound [7] | $\Omega\left(n\sqrt{\frac{\delta}{\mu}}\log\frac{1}{\varepsilon}\right)$ | SS for $f_i$'s |
| Client Sampling | SVRP [33] | $\mathcal{O}\left(\left(n+\frac{\delta^2}{\mu^2}\right)\log\frac{1}{\varepsilon}\right)^{(1)}$ | SC for $f_i$'s, AveSS |
| | Catalyzed SVRP [33] | $\mathcal{O}\left(\left(n+n^{3/4}\sqrt{\frac{\delta}{\mu}}\right)\log\frac{1}{\varepsilon}\log\frac{L}{\mu}\right)^{(2)}$ | SC for $f_i$'s, AveSS |
| | SVRS (Thm 3.3) | $\mathcal{O}\left(\left(n+\sqrt{n}\cdot\frac{\delta}{\mu}\right)\log\frac{1}{\varepsilon}\right)$ | AveSS |
| | AccSVRS (Thm 3.6) | $\mathcal{O}\left(\left(n+n^{3/4}\sqrt{\frac{\delta}{\mu}}\right)\log\frac{1}{\varepsilon}\right)$ | AveSS |
| | Lower bound (Thm 4.4) | $\Omega\left(n+n^{3/4}\sqrt{\frac{\delta}{\mu}}\log\frac{1}{\varepsilon}\right)^{(3)}$ | AveSS |

(1) The rate only applies to $\mathbb{E}\left\|\boldsymbol{x}_k-\boldsymbol{x}_*\right\|^2$, otherwise it would introduce $L$ in the log term; (2) The term $\log(L/\mu)$ comes from the Catalyst framework. See Appendix C for the detail. (2, 3) Here we only list the rates of the common ill-conditioned case: $\mu=\mathcal{O}(\delta/\sqrt{n})$. See Appendices for the remaining case. *Notation:* $\delta$=similarity parameter (both for SS and AveSS), $L$=smoothness constant of $f$, $\mu$=strong convexity constant of $f$(or $f_i$'s), $\varepsilon$=error of the solution for $\mathbb{E}f(\boldsymbol{x}_k)-f(\boldsymbol{x}_*)$. Here $L\geq\delta\geq\mu\gg\epsilon>0$. *Abbreviation:* SC=strong convexity, SS=second-order similarity, AveSS=average SS.

AccSVRS achieves the same communication bound of $\tilde{\mathcal{O}}(n+n^{3/4}\sqrt{\delta/\mu})$ as Catalyzed SVRP. Specifically, our bound is entirely smoothness-free and slightly outperforms Catalyzed SVRP, featuring a log improvement and not requiring component strong convexity.

- Third, by considering the proximal incremental first-order oracle in the centralized distributed framework, we establish a lower bound, which nearly matches the upper bound of AccSVRS in ill-conditioned cases.

## 1.2 Related Work

**Gradient sliding/Oracle Complexity Separation.** For optimization problems with a separated structure or multiple building blocks, such as Eq. (1), there are scenarios where computing the gradients/values of some parts (or the whole) is more expensive than the others (or a partial one). In response to this challenge, techniques such as the gradient-sliding method [37] and the concept of oracle complexity separation [28] have emerged. These methods advocate for the infrequent use of more expensive oracles compared to their less resource-intensive counterparts. This strategy has found applications in zero-order [12, 21, 28, 53], first-order [37–39, 31] and high-order methods [31, 24, 3], as well as in addressing saddle point problems [4, 13]. Our algorithms can be viewed as a variance-reduced version of gradient sliding tailored to leverage the similarity assumption.

**Distributed optimization under similarity.** Distributed optimization has a long history with a plethora of existing works and surveys. To streamline our discussion, we only list the most relevant references, particularly under the similarity and strong convexity assumptions. In the full participation setting, which involves deterministic methods, the first algorithm credits to DANE [51], though its analysis is limited to quadratic objectives. Subsequently, AIDE [50], DANE-LS and DANE-HB [58] improved the rates for quadratic objective; Disco [62] SPAG [27], ACN [1] and DiRegINA [18] improved the rates for self-concordant objectives. As for general strongly convex objectives, Sun et al. [54] introduced the SONATA algorithm, and Tian et al. [55] proposed accelerated SONATA. However, their complexity bounds include additional log factors. These factors have recently been removed by Accelerated Extragradient [35], whose complexity bound perfectly matches the lower bound in [7]. We highly recommend the comparison of rates in [35, Table 1] for a comprehensive overview. Once the discussion of deterministic methods is concluded, Khaled and Jin [33] shifted their focus to stochastic methods using client sampling. They proposed SVRP and its Catalyzed version, both of which exhibited superior rates compared to deterministic methods.

## 2 Preliminaries

**Notation.** We denote vectors by lowercase bold letters (e.g., $\boldsymbol{w}, \boldsymbol{x}$), and matrices by capital bold letters (e.g., $\boldsymbol{A}, \boldsymbol{B}$). We let $\|\cdot\|$ be the $\ell_2$-norm for vectors, or induced $\ell_2$-norm for a given matrix: $\|\boldsymbol{A}\| = \sup_{\boldsymbol{u} \neq 0} \|\boldsymbol{A}\boldsymbol{u}\| / \|\boldsymbol{u}\|$. We abbreviate $[n] = \{1, \ldots, n\}$ and $\boldsymbol{I}_d \in \mathbb{R}^{d \times d}$ is the identity matrix. We use $\boldsymbol{0}$ for the all-zero vector/matrix, whose size will be specified by a subscript, if necessary, and otherwise is clear from the context. We denote $\text{Unif}(\mathcal{S})$ as the uniform distribution over set $\mathcal{S}$. We say $T \sim \text{Geom}(p)$ for $p \in (0, 1]$ if $\mathbb{P}(T = k) = (1-p)^{k-1}p, \forall k \in \{1, 2, \ldots\}$, i.e., $T$ obeys a geometric distribution. We adopt $\mathbb{E}_k$ as the expectation for all randomness appeared in step $k$, and $\mathbb{1}_A$ as the indicator function on event $A$, i.e., $\mathbb{1}_A = 1$ if event $A$ holds, and $0$ otherwise. We use $\mathcal{O}(\cdot), \Omega(\cdot), \Theta(\cdot)$ and $\tilde{\mathcal{O}}(\cdot)$ notation to hide universal constants and log-factors. We define the Bregman divergence induced by a differentiable (convex) function $h \colon \mathbb{R}^d \to \mathbb{R}$ as $D_h(\boldsymbol{x}, \boldsymbol{y}) := h(\boldsymbol{x}) - h(\boldsymbol{y}) - \langle \nabla h(\boldsymbol{y}), \boldsymbol{x} - \boldsymbol{y} \rangle$.

**Definitions.** We present the following common definitions used in this paper.

**Definition 2.1** *A differentiable function $g \colon \mathbb{R}^d \to \mathbb{R}$ is $\mu$-strongly convex (SC) if*

$$g(\boldsymbol{y}) \geq g(\boldsymbol{x}) + \langle \nabla g(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \frac{\mu}{2} \|\boldsymbol{y} - \boldsymbol{x}\|^2, \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d. \tag{2}$$

*Particularly, if $\mu = 0$, we say that $g$ is convex.*

**Definition 2.2** *A differentiable function $g \colon \mathbb{R}^d \to \mathbb{R}$ is L-smooth if*

$$g(\boldsymbol{y}) \leq g(\boldsymbol{x}) + \langle \nabla g(\boldsymbol{x}), \boldsymbol{y} - \boldsymbol{x} \rangle + \frac{L}{2} \|\boldsymbol{y} - \boldsymbol{x}\|^2, \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d. \tag{3}$$

There are many basic inequalities involving strong convexity and smoothness, see [22, Appendix A.1] for an introduction. Next, we present the definition of second-order similarity in distributed optimization.

**Definition 2.3** *The differentiable functions $f_i$'s satisfy $\delta$-average second-order similarity (AveSS) if the following inequality holds for $f_i$'s and $f = \frac{1}{n} \sum_{i=1}^n f_i$:*

$$\text{(AveSS)} \quad \frac{1}{n} \sum_{i=1}^n \|[\nabla[f_i - f](\boldsymbol{x}) - \nabla[f_i - f](\boldsymbol{y})]\|^2 \leq \delta^2 \|\boldsymbol{x} - \boldsymbol{y}\|^2, \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d. \tag{4}$$

**Definition 2.4** *The differentiable functions $f_i$'s satisfy $\delta$-component second-order similarity (SS) if the following inequality holds for $f_i$'s and $f = \frac{1}{n} \sum_{i=1}^n f_i$:*

$$\text{(SS)} \quad \|[\nabla[f_i - f](\boldsymbol{x}) - \nabla[f_i - f](\boldsymbol{y})]\|^2 \leq \delta^2 \|\boldsymbol{x} - \boldsymbol{y}\|^2, \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d, i \in [n]. \tag{5}$$

Definitions 2.3 and 2.4 first appear in [33], which is an analogy to (average) smoothness in prior literature [63]. Particularly, $f_i$'s satisfy $\delta$-AveSS implies that $(f - f_i)$'s satisfy $\delta$-average smoothness, while $f_i$'s satisfy $\delta$-SS implies that $(f - f_i)$'s satisfy $\delta$-smoothness. Additionally, many researchers [32, 7, 51, 62, 54, 35] use the equivalent one defined by Hessian similarity (HS) if assuming that $f_i$'s are twice differentiable. Thus we also list them below and leave the derivation in Appendix B.

$$\text{(AveHS)} \left\| \frac{1}{n} \sum_{i=1}^n \left[ \nabla^2 f_i(\boldsymbol{x}) - \nabla^2 f(\boldsymbol{x}) \right]^2 \right\| \leq \delta^2; \text{(HS)} \left\| \nabla^2 f_i(\boldsymbol{x}) - \nabla^2 f(\boldsymbol{x}) \right\| \leq \delta, \forall i \in [n]. \tag{6}$$

Since our algorithm is a first-order method, we adopt the gradient description of similarity (Definitions 2.3 and 2.4) without assuming twice differentiability for brevity.

As mentioned in [7, 54], if $f_i$'s satisfy $\delta$-AveSS (or SS), and $f$ is $\mu$-strongly convex and $L$-smooth, then generally $L \gg \delta \gg \mu > 0$ for large datasets in practice. Therefore, researchers aim to develop algorithms that achieve communication complexity solely related to $\delta, \mu$ (or log terms of $L$). This is also our objective. To finish this section, we will clarify several straightforward yet essential propositions, and the proofs are deferred to Appendix A.

**Proposition 2.5** *We have the following properties among SS, AveSS, and SC: 1) $\delta$-SS implies $\delta$-AveSS, but $\delta$-AveSS only implies $\sqrt{n}\delta$-SS. 2) If $f_i$'s satisfy $\delta$-SS and $f$ is $\mu$-strongly convex, then for all $i \in [n]$, $f_i(\cdot) + \frac{\delta - \mu}{2} \|\cdot\|^2$ is convex, i.e., $f_i$ is $(\delta - \mu)$-almost convex [14].*

## 3 Algorithm and Theory

In this section, we introduce our main algorithms, which are developed to solve the distributed optimization problem in Eq. (1) under Assumption 1 below:

**Assumption 1** *We assume that $f_i$'s satisfy $\delta$-AveSS, and $f$ is $\mu$-strongly convex with $\delta \geq \mu > 0$.*

Assumption 1 does not need each $f_i$ to be $\mu$-strongly convex. In fact, it is acceptable that $f_i$'s are non-convex, since by Proposition 2.5, $f_i$'s are $(\sqrt{n}\delta - \mu)$-almost convex [14]. In the following, we first propose our new algorithm SVRS, which combines the techniques of gradient sliding and variance reduction, resulting in improved rates. Then we establish the directly accelerated method motivated by [6].

### 3.1 No Acceleration Version: SVRS

We first show the one-epoch Stochastic Variance-Reduced Sliding (SVRS[1ep]) method in Algorithm 1. Before delving into the theoretical analysis, we present some key insights into our method. These insights aim to enhance comprehension and facilitate connections with other algorithms.

**Variance Reduction.** Our algorithm can be viewed as adding variance reduction from [35]. Besides the acceleration step, the main difference lies in the proximal step, where Kovalev et al. [35] solved:

$$\boldsymbol{x}_{t+1} \approx \underset{\boldsymbol{x} \in \mathbb{R}^d}{\arg\min} \, B_\theta^t(\boldsymbol{x}) := \langle \nabla f(\boldsymbol{x}_t) - \nabla f_1(\boldsymbol{x}_t), \boldsymbol{x} - \boldsymbol{x}_t \rangle + \frac{1}{2\theta} \|\boldsymbol{x} - \boldsymbol{x}_t\|^2 + f_1(\boldsymbol{x}).$$

To save the heavy communication burden of calculating $\nabla f(\boldsymbol{x}_t)$, we apply client sampling by selecting a random $\nabla f_{i_t}(\boldsymbol{x}_t)$ in the $t$-th step. However, this substitution introduces significant noise. To mitigate this, we incorporate a correction term $\boldsymbol{g}_t = \nabla f_{i_t}(\boldsymbol{w}_0) - \nabla f(\boldsymbol{w}_0)$ from previous wisdom [29] to reduce the variance.

**Gradient sliding.** Our algorithm can be viewed as adding gradient sliding from SVRP [33]. The main difference also lies in the proximal point problem, where Khaled and Jin [33] solved:

$$\boldsymbol{x}_{t+1} \approx \underset{\boldsymbol{x} \in \mathbb{R}^d}{\arg\min} \, C_\theta^t(\boldsymbol{x}) := \langle -\boldsymbol{g}_t, \boldsymbol{x} - \boldsymbol{x}_t \rangle + \frac{1}{2\theta} \|\boldsymbol{x} - \boldsymbol{x}_t\|^2 + f_{i_t}(\boldsymbol{x}).$$

Here we adopt a fixed proximal function $f_1$ instead of $f_{i_t}$, which can be viewed as approximating $f_{i_t}(\boldsymbol{x}) \approx f_1(\boldsymbol{x}) + [f_{i_t} - f_1](\boldsymbol{x}) + \langle \nabla [f_{i_t} - f_1](\boldsymbol{x}_t), \boldsymbol{x} - \boldsymbol{x}_t \rangle + \frac{1}{2\theta'} \|\boldsymbol{x} - \boldsymbol{x}_t\|^2$ with a properly chosen $\theta' > 0$. Such a modification is motivated by [35], where they reformulated the objective as $f(\boldsymbol{x}) = [f(\boldsymbol{x}) - f_1(\boldsymbol{x})] + f_1(\boldsymbol{x})$. Thus they could employ gradient sliding to skip heavy computations of $\nabla [f - f_1](\boldsymbol{x})$ by utilizing the easy computations of $\nabla f_1(\boldsymbol{x})$ more times. Fixing the proximal function $f_1$ leads to the same metric space owned by $f_1$ in each step, which could benefit the analysis and alleviate the requirements on $f_i$'s compared to SVRP. Indeed, in our setting $f_1$ can be replaced by any other **fixed** client $f_b, b \in [n]$. In this case, the master node would be $f_b$ instead of $f_1$.

**Bregman-SVRG.** Our algorithm can be viewed as the classical Bregman-SVRG [20] with the reference function $f_1(\cdot) + \frac{1}{2\theta} \|\cdot\|^2$ after introducing the Bregman divergence:

$$\boldsymbol{x}_{t+1} \approx \underset{\boldsymbol{x} \in \mathbb{R}^d}{\arg\min} \, A_\theta^t(\boldsymbol{x}) \overset{(7)}{=} \underset{\boldsymbol{x} \in \mathbb{R}^d}{\arg\min} \, \langle \nabla f_{i_t}(\boldsymbol{x}_t) - \nabla [f_{i_t} - f](\boldsymbol{w}_0), \boldsymbol{x} - \boldsymbol{x}_t \rangle + D_{f_1(\cdot) + \frac{1}{2\theta} \|\cdot\|^2}(\boldsymbol{x}, \boldsymbol{x}_t).$$

We need to emphasize that the proof of Bregman-SVRG requires additional structural assumptions [20, Assumption 3], which is not directly applicable in our setting. Hence, the rigorous proof of Bregman-SVRG under our similarity assumption is still meaningful as far as we are concerned.

#### 3.1.1 Communication Complexity under Distributed Settings

When applied to the distributed system, the communication complexity of SVRS[1ep] can be described as follows: At the beginning of each epoch, the master (corresponding to $f_1$) sends $\boldsymbol{w}_0$ to all clients. Each client computes $\nabla f_i(\boldsymbol{w}_0)$ from its local data and sends it back to the master. The master then builds $\nabla f(\boldsymbol{w}_0)$ after collecting all $\nabla f_i(\boldsymbol{w}_0)$'s. The communication complexity is $2(n-1)$ in this

---

**Algorithm 1** $\text{SVRS}^{1\text{ep}}(f, \boldsymbol{w}_0, \theta, p)$

---

1: **Input:** $\boldsymbol{w}_0 \in \mathbb{R}^d$, $p \in (0,1)$, $\theta > 0$
2: Initialize $\boldsymbol{x}_0 = \boldsymbol{w}_0$, compute $\nabla f(\boldsymbol{w}_0)$, and set $T \sim \text{Geom}(p)$
3: **for** $t = 0, 1, 2, \ldots, T-1$ **do**
4:     Sample $i_t \sim \text{Unif}([n])$ and compute $\boldsymbol{g}_t = \nabla f_{i_t}(\boldsymbol{w}_0) - \nabla f(\boldsymbol{w}_0)$
5:     Approximately solve the local proximal point problem:

$$\boldsymbol{x}_{t+1} \approx \underset{\boldsymbol{x} \in \mathbb{R}^d}{\arg\min} \, A_\theta^t(\boldsymbol{x}) := \langle \nabla f_{i_t}(\boldsymbol{x}_t) - \nabla f_1(\boldsymbol{x}_t) - \boldsymbol{g}_t, \boldsymbol{x} - \boldsymbol{x}_t \rangle + \frac{1}{2\theta} \|\boldsymbol{x} - \boldsymbol{x}_t\|^2 + f_1(\boldsymbol{x}) \tag{7}$$

6: **end for**
7: **Output:** $\boldsymbol{x}_T$

---

case. Next, the algorithm enters into the loop iterations. In each iteration, the master only sends current $\boldsymbol{x}_t$ to the chosen client $i_t$. The $i_t$-th client computes $\nabla f_{i_t}(\boldsymbol{x}_t)$ and sends it to the master (the first client). Then the master solves (inexactly) the local problem (Line 5 in Algorithm 1) to get an inexact solution $\boldsymbol{x}_{t+1}$. The communication complexity is 2 in this case. Thus, the total communication complexity of $\text{SVRS}^{1\text{ep}}$ is $2(n-1) + 2T$. Note that $\mathbb{E}T = 1/p$ and generally $p = 1/n$. We obtain that one epoch communication complexity is $4n - 2$ in expectation.

We would like to emphasize that our setup differs from that in [41, 46], where the authors assume the nodes can perform calculations and transmit vectors in parallel. We recognize the significance of both setups. However, there are situations where communication is more expensive than computation. For instance, in a business network or communication network the communication between any two nodes can result in charges and the risk of information leakage. To mitigate these costs, we should reduce the frequency of communication. Thus, we focus on the nonparallel setting.

### 3.1.2 Convergence Analysis of SVRS

Based on the one-epoch method $\text{SVRS}^{1\text{ep}}$, we could introduce our non-accelerated algorithm SVRS, which starts from $\boldsymbol{w}_0 \in \mathbb{R}^d$ and repeatedly performs the update[7]

$$\boldsymbol{w}_{k+1} = \text{SVRS}^{1\text{ep}}(f, \boldsymbol{w}_k, \theta, p), \ \forall k \geq 0.$$

Now we derive the convergence rate of SVRS[8]. The main technique we apply is replacing the Euclidean distance with the Bergman divergence. Denote the reference function

$$h(\boldsymbol{x}) := f_1(\boldsymbol{x}) + \frac{1}{2\theta} \|\boldsymbol{x}\|^2 - f(\boldsymbol{x}). \tag{8}$$

By Assumption 1 and 1) in Proposition 2.5, we see that $f_i$'s are $\sqrt{n}\delta$-SS. i.e., $[f_1 - f](\cdot)$ is $(\sqrt{n}\delta)$-smooth. Thus, $h(\cdot)$ is $(\frac{1}{\theta} - \sqrt{n}\delta)$-strongly convex and $(\frac{1}{\theta} + \sqrt{n}\delta)$-smooth if $\theta < \frac{1}{\sqrt{n}\delta}$, that is,

$$0 \leq \frac{1 - \sqrt{n}\theta\delta}{2\theta} \|\boldsymbol{x} - \boldsymbol{y}\|^2 \overset{(2)}{\leq} D_h(\boldsymbol{x}, \boldsymbol{y}) \overset{(3)}{\leq} \frac{1 + \sqrt{n}\theta\delta}{2\theta} \|\boldsymbol{x} - \boldsymbol{y}\|^2. \tag{9}$$

Hence, if $\sqrt{n}\theta\delta = \Theta(1)$, $h(\cdot)$ is nearly a rescaled Euclidean norm since its condition number related to $\|\cdot\|$ is $\frac{1+\sqrt{n}\theta\delta}{1-\sqrt{n}\theta\delta} = \Theta(1)$. Next, we employ the properties of the Bregman divergence $D_h(\cdot, \cdot)$ to build the one-epoch progress of $\text{SVRS}^{1\text{ep}}$ as shown below:

**Lemma 3.1** *Suppose Assumption 1 holds. Let* $\boldsymbol{w}^+ = \text{SVRS}^{1\text{ep}}(f, \boldsymbol{w}_0, \theta, p)$ *with* $\theta = 1/(4\sqrt{n}\delta)$, *and the approximated solution* $\boldsymbol{x}_{t+1}$ *satisfies*

$$\left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2 \leq \frac{\mu}{20\theta} \left\|\boldsymbol{x}_t - \underset{\boldsymbol{x} \in \mathbb{R}^d}{\arg\min} \, A_\theta^t(\boldsymbol{x})\right\|^2, \forall t \geq 0. \tag{10}$$

---

[7]See Algorithm 3 in Appendix D for the details.
[8]Similar results for the popular loopless version [34] can also be derived, see Appendix D.5 for the detail.

---

**Algorithm 2** Accelerated SVRS (AccSVRS)

---

1: **Input:** $z_0 = y_0 \in \mathbb{R}^d, p, \tau \in (0, 1), \alpha, \theta > 0, K \in \{1, 2, \dots\}$
2: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
3:     $x_{k+1} = \tau z_k + (1 - \tau) y_k$
4:     $y_{k+1} = \mathrm{SVRS}^{1\mathrm{ep}}(f, x_{k+1}, \theta, p)$
5:     $\mathcal{G}_{k+1} = p\left(\nabla[f_1 - f_{j_k}](x_{k+1}) - \nabla[f_1 - f_{j_k}](y_{k+1}) + \frac{1}{\theta}(x_{k+1} - y_{k+1})\right), j_k \sim \mathrm{Unif}([n])$
6:     $z_{k+1} = \arg\min_{z \in \mathbb{R}^d} \frac{1}{2\alpha} \|z - z_k\|^2 + \langle \mathcal{G}_{k+1}, z \rangle + \frac{3\mu}{20} \|z - y_{k+1}\|^2 = \frac{z_k + 0.3\mu\alpha y_{k+1} - \alpha\mathcal{G}_{k+1}}{1 + 0.3\mu\alpha}$
7: **end for**
8: **Output:** $y_K$

---

Then for all $x \in \mathbb{R}^d$ that is independent of the indices $i_1, i_2, \dots, i_T$ in $\mathrm{SVRS}^{1\mathrm{ep}}(f, w_0, \theta, p)$, we have

$$\mathbb{E}f(w^+) - f(x) \leq \mathbb{E}\, p\langle x - w_0, \nabla h(w^+) - \nabla h(w_0)\rangle - \left(p - \frac{2}{9n}\right)D_h(w_0, w^+) - \frac{2\mu\theta}{5}D_h(x, w^+). \tag{11}$$

**Remark 3.2** *We note that some papers [10, 11] assume the smoothness and convexity of component functions, and adopt local updates for solving the proximal step. However, we replace these assumptions with a proximal approximately solvable assumption (10), which could even cover some nonsmooth and non-convex but proximal trackable component functions. We regard our assumption as more essential since the local updates can be viewed as partially solving this proximal step.*

The proof of Lemma 3.1 is left in Appendix D.1. From Lemma 3.1, we find a well-behaved proximal operator is sufficient to ensure favorable progress. Finally, we establish the convergence rate and communication complexity of the SVRS method, and the proof is deferred to Appendix D.2.

**Theorem 3.3** *Suppose Assumption 1 holds. If in* $\mathrm{SVRS}^{1\mathrm{ep}}$*(Algorithm 1), the hyperparameters are set as* $\theta = 1/(4\sqrt{n}\delta), p = 1/n$*, and the approximate solution* $x_{t+1}$ *in each proximal step satisfies Eq.* (10)*. Then for any error* $\varepsilon > 0$*, when*

$$k \geq K_1 := \max\left\{2, \frac{5\delta}{\mu\sqrt{n}}\right\} \log \frac{3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right)[f(w_0) - f(x_*)]}{\varepsilon},$$

*i.e., after* $\tilde{\mathcal{O}}(n + \sqrt{n}\delta/\mu)$ *communications in expectation, we obtain that* $\mathbb{E}f(w_k) - f(x_*) \leq \varepsilon$.

**Remark 3.4** *Our results enjoy the following advantages over SVRP [33]: The convergence of SVRP ([33, Theorem 2]) only applied to* $\mathbb{E}\|w_k - x_*\|^2$*, which can also be derived by our results from strong convexity:* $f(w_k) - f(x_*) \geq \frac{\mu}{2}\|w_k - x_*\|^2$*. However, the reverse is not applicable since we do not assume the smoothness of* $f$*, or indeed the smoothness coefficient is very large. Moreover, for ill-conditioned problems (e.g.,* $\delta/\mu \gg \sqrt{n}$*), our step size* $1/(4\sqrt{n}\delta)$ *is much larger than* $\mu/(2\delta^2)$ *used in SVRP, and the convergence rate is also faster than SVRP:* $\tilde{\mathcal{O}}(n + \sqrt{n}\delta/\mu)$ *vs.* $\tilde{\mathcal{O}}(n + \delta^2/\mu^2)$*. Finally, we do not need the strong convexity assumption of component functions.*

## 3.2 Acceleration Version: AccSVRS

Now we apply the classical interpolation technique motivated by Katyusha X [6] to establish accelerated SVRS (AccSVRS, Algorithm 2). The main difference between AccSVRS and Katyusha X is due to the different choices of distance spaces. Specifically, we adopt $D_h(\cdot, \cdot)$ instead of the Euclidean distance used in Katyusha X. Thus, the gradient mapping step (corresponding to Step 2 in [6, (4.1)]) should be built on the reference function $h(\cdot)$ defined in Eq. (8), i.e., $\nabla h(x_{k+1}) - \nabla h(y_{k+1})$ instead of $(x_{k+1} - y_{k+1})/\theta$. Moreover, noting that $\nabla h(\cdot)$ could involve the heavy gradient computing part $\nabla f(\cdot)$, we further employ its stochastic version (Step 5 in Algorithm 2) to reduce the overall communication complexity.

Next, we delve into the convergence analysis. We first give the core lemma for AccSVRS, which is also motivated by the framework of Katyusha X [6]. The proof is deferred to Appendix D.3.

**Lemma 3.5** *Suppose Assumption 1 holds, and $\theta = 1/(4\sqrt{n}\delta), p = 1/n, \alpha \leq n\theta/(2\tau)$ in Algorithm 2, where $\mathrm{SVRS}^{\mathrm{1ep}}(f, \boldsymbol{x}_{k+1}, \theta, p)$ satisfies Eq. (10) in each iteration. Then for all $\boldsymbol{x} \in \mathbb{R}^d$ that is independent of the random indices $i_1^{(k)}, i_2^{(k)}, \ldots, i_T^{(k)}$ in $\mathrm{SVRS}^{\mathrm{1ep}}(f, \boldsymbol{x}_{k+1}, \theta, p)$, we have that*

$$\mathbb{E}_k \frac{\alpha}{\tau} [f(\boldsymbol{y}_{k+1}) - f(\boldsymbol{x})] \leq \mathbb{E}_k (1-\tau) \cdot \frac{\alpha}{\tau} [f(\boldsymbol{y}_k) - f(\boldsymbol{x})] + \frac{\|\boldsymbol{x} - \boldsymbol{z}_k\|^2}{2} - \frac{1 + 0.3\mu\alpha}{2} \|\boldsymbol{x} - \boldsymbol{z}_{k+1}\|^2 . \tag{12}$$

Finally, we present the convergence rate and communication complexity of AccSVRS based on Lemma 3.5, and the proof is left in Appendix D.4.

**Theorem 3.6** *Suppose Assumption 1 holds. Consider AccSVRS with the following hyperparameters*

$$\theta = \frac{1}{4\sqrt{n}\delta}, p = \frac{1}{n}, \tau = \frac{1}{4} \min \left\{ 1, \frac{n^{1/4}}{2} \sqrt{\frac{\mu}{\delta}} \right\}, \alpha = \frac{\sqrt{n}}{8\delta\tau},$$

*and Eq. (10) is satisfied in each iteration of $\mathrm{SVRS}^{\mathrm{1ep}}(f, \boldsymbol{x}_{k+1}, \theta, p)$. Then for any $\varepsilon > 0$, when*

$$k \geq K_2 := \max \left\{ 4, 8n^{-1/4}\sqrt{\delta/\mu} \right\} \log \frac{2[f(\boldsymbol{y}_0) - f(\boldsymbol{x}_*)]}{\varepsilon},$$

*i.e., after $\tilde{\mathcal{O}}\left(n + n^{3/4}\sqrt{\delta/\mu}\right)$ communications in expectation, we obtain that $\mathbb{E}f(\boldsymbol{y}_k) - f(\boldsymbol{x}_*) \leq \varepsilon$.*

**Remark 3.7** *Although roughly the same as the communication complexity obtained by Catalyzed SVRP in [33, Theorem 3], our results have the following advantages.*

***Fewer assumptions.*** *Except for the strong convexity of $f$ and AveSS of $f_i$'s, we do not need to assume component strong convexity appearing in [33, Assumption 2].*

***Inexact proximal step.*** *Khaled and Jin [33, Theorem 3] require exact evaluations of the proximal operator, though they mention that this is only for the convenience of analysis. Our framework allows approximated solutions in each proximal step, and the approximation criterion (10) is error-independent, i.e., irrelevant to the final error $\varepsilon$. Since the local proximal function is strongly convex, we could solve the problem in a few steps if additionally assuming the smoothness of $f_1$.*

***Smoothness-free bound.*** *As shown in [33, Appendix G.1] or Appendix C, even if an exact proximal step is allowed, a dependence on the smoothness coefficient would be introduced in the total communication iterations of Catalyzed SVRP, though only in a log scale. Our directly accelerated method has no dependence on the smoothness coefficient.*

### 3.3 Gradient Complexity under Smooth Assumption

Due to the importance of total computation in the machine learning and optimization community, we consider a more common setup by **additionally assuming** that $f_1$ is $L$-smooth with $L \geq \delta \geq \mu > 0$, which together with Assumption 1 facilitates the quantification of Eq. (10). Then we can compute the total gradient complexity for AccSVRS as shown below. By Proposition 2.5 and our assumptions, $A_\theta^t(\boldsymbol{x})$ is $(\frac{1}{\theta} - \sqrt{n}\delta)$-strongly convex and $(\frac{1}{\theta} + L)$-smooth. Using accelerated methods starting from $\boldsymbol{x}_t$, we can guarantee that Eq. (10) holds after $T_{\mathrm{app}} = \tilde{\mathcal{O}}\left(\sqrt{\frac{1+\theta L}{1-\sqrt{n}\theta\delta}}\right) = \tilde{\mathcal{O}}\left(1 + n^{-1/4}\sqrt{L/\delta}\right)$ iterations with the choice of $\theta$ in Theorem 3.6. Hence, the total gradient calls in expectation are

$$\mathcal{O}(n T_{\mathrm{app}} \cdot K_2) = \tilde{\mathcal{O}}\left(n + n^{3/4}\left(\sqrt{\delta/\mu} + \sqrt{L/\delta}\right) + \sqrt{nL/\mu}\right).$$

Since $\delta \in [\mu, L]$, we recover the optimal gradient complexity $\tilde{\mathcal{O}}(n + n^{3/4}\sqrt{L/\mu})$ for the average smooth setting [63, Table 1] if neglecting log factors. Particularly, when $\delta = \Theta(\sqrt{\mu L})$, we even obtain the nearly optimal gradient complexity $\tilde{\mathcal{O}}(n + \sqrt{nL/\mu})$ for the component smooth setting [25, 26, 56]. We leave the details in Appendix E. Although the gradient complexity is not the primary focus of our work, we have demonstrated that the gradient complexity bound of AccSVRS is nearly optimal for certain values of $\delta$ in specific cases.

## 4 Lower Bound

In this section, we establish the lower bound of the communication complexity, which nearly matches the upper bound of AccSVRS.

## 4.1 Definition of Algorithms

In this subsection, we specify the class of algorithms to which our lower bound can apply. We first introduce the Proximal Incremental First-order Oracle (PIFO) [56, 25], which is defined as $h_{f_i}^{\mathrm{P}}(\boldsymbol{x}, \gamma) = [f_i(\boldsymbol{x}), \nabla f_i(\boldsymbol{x}), \mathrm{prox}_{f_i}^\gamma(\boldsymbol{x})]$ with $\gamma > 0$. Here the proximal operator is defined as $\mathrm{prox}_{f_i}^\gamma(\boldsymbol{x}) := \arg\min_{\boldsymbol{u}}\{f_i(\boldsymbol{u}) + \frac{1}{2\gamma}\|\boldsymbol{x} - \boldsymbol{u}\|^2\}$. In addition to the local zero-order and first-order information of $f_i$ at $\boldsymbol{x}$, the PIFO $h_{f_i}^{\mathrm{P}}(\boldsymbol{x}, \gamma)$ also provides some global information through the proximal operator[9]. Then we assume the algorithm has access to the PIFO and the definition of algorithms is presented as follows.

**Definition 4.1** *Consider a randomized algorithm $\mathcal{A}$ to solve problem* (1). *Suppose the number of communication rounds is $T$. At the initialization stage, the master node 1 communicates with all the others. In round $t$ $(0 \leq t \leq T-1)$, the algorithm samples a node $i_t \sim \mathrm{Unif}([n])$, and node 1 communicates with node $i_t$. Then the algorithm samples a Bernoulli random variable $a_t$ with constant expectation $c_0/n$. If $a_t = 1$, node 1 communicates with all the others. Define the information set $\mathcal{I}_{t+1}$ as the set of all the possible points $\mathcal{A}$ can obtain after round $t$. The algorithm updates $\mathcal{I}_{t+1}$ based on the linear-span operation and PIFO, and finally outputs a certain point in $\mathcal{I}_T$.*

At the initialization stage, the communication cost is $2(n-1)$. In each communication round, the Bernoulli random variable $a_t$ determines whether the master node communicates with all the others, i.e., whether to calculate the full gradient. Since $\mathbb{E}a_t = c_0/n$, the expected communication cost of each round is of the order $\Theta(1)$. Thus the total communication cost is of the order $\Theta(n+T)$ and we can use $T$ to measure the communication complexity. Moreover, one can check Algorithm 2 satisfies Definition 4.1. The formal definition and detailed analysis are deferred to Appendix F.1.

## 4.2 The Construction and Results

In this section, we construct a hard instance of problem (1) and then use it to establish the lower bound. Due to space limitations, we only present several key properties. The complete framework of construction is deferred to Appendix F.2.

Inspired by [25], we consider the class of matrices $\boldsymbol{B}(m, \zeta) = \begin{bmatrix} 1 & -1 & & \\ & \ddots & \ddots & \\ & & 1 & -1 \\ & & & \zeta \end{bmatrix} \in \mathbb{R}^{m \times m}$. This class

of matrices is widely used to establish lower bounds for minimax optimization problems [61, 49, 59], and $\boldsymbol{A}(m, \zeta) := \boldsymbol{B}(m, \zeta)^\top \boldsymbol{B}(m, \zeta)$ is the well-known tridiagonal matrix in the analysis of lower bounds for convex optimization [47, 40, 15]. Denote the $l$-th row of $\boldsymbol{B}(m, \zeta)$ as $\boldsymbol{b}_l(m, \zeta)^\top$. We partition the row vectors of $\boldsymbol{B}(m, \zeta)$ according to the index sets $\mathcal{L}_i = \{l : 1 \leq l \leq m, l \equiv i - 1 \,(\mathrm{mod}\,(n-1))\}$ for $2 \leq i \leq n$ and $\mathcal{L}_1 = \varnothing$[10]. These sets are mutually exclusive and their union is $[m]$. Then we consider the following problem

$$\min_{\boldsymbol{x} \in \mathbb{R}^m} r(\boldsymbol{x}; m, \zeta, c) = \frac{1}{n}\sum_{i=1}^n \left[ r_i(\boldsymbol{x}; m, \zeta, c) := \begin{cases} \frac{c}{2}\|\boldsymbol{x}\|^2 - n\langle\boldsymbol{e}_1, \boldsymbol{x}\rangle & \text{for } i = 1, \\ \frac{c}{2}\|\boldsymbol{x}\|^2 + \frac{n}{2}\sum_{l \in \mathcal{L}_i}\|\boldsymbol{b}_l(m, \zeta)^\top \boldsymbol{x}\|^2 & \text{for } i \neq 1. \end{cases} \right] \quad (13)$$

Here $\boldsymbol{e}_i \in \mathbb{R}^m$ denotes the unit vector with the $i$-th element equal to 1 and others equal to 0. Then one can check $r(\boldsymbol{x}; m, \zeta, c) = \frac{1}{2}\boldsymbol{x}^\top \boldsymbol{A}(m, \zeta)\boldsymbol{x} + \frac{c}{2}\|\boldsymbol{x}\|^2 - \langle\boldsymbol{e}_1, \boldsymbol{x}\rangle$. Clearly, $r$ is $c$-strongly convex. We can also determine the AveSS parameter as follows. The proof is deferred to Appendix F.3.

**Proposition 4.2** *Suppose that $0 < \zeta \leq \sqrt{2}$, $n \geq 3$ and $m \geq 3$. Then $r_i$'s satisfy $\sqrt{8n+4}$-AveSS.*

Define the subspaces $\{\mathcal{F}_k\}_{k=0}^m$ as $\mathcal{F}_0 = \{\boldsymbol{0}\}$ and $\mathcal{F}_k = \mathrm{span}\{\boldsymbol{e}_1, \boldsymbol{e}_2, \ldots, \boldsymbol{e}_k\}$ for $1 \leq k \leq m$. The next lemma is fundamental to our analysis. The proof is deferred to Appendix F.5.

**Lemma 4.3** *Suppose the algorithm $\mathcal{A}$ satisfies Definition 4.1 and apply it to solve problem* (13) *with $n \geq 3$ and $m \geq 4$. We have (i) $\mathcal{I}_0 = \mathcal{F}_1$. (ii) Suppose $\mathcal{I}_t \subseteq \mathcal{F}_k$ $(1 \leq k \leq m-3)$. If $i_t$ satisfies $k \in \mathcal{L}_{i_t}$ or $a_t = 1$, then $\mathcal{I}_{t+1} \subseteq \mathcal{F}_{k+3}$; otherwise, $\mathcal{I}_{t+1} \subseteq \mathcal{F}_k$.*

---

[9]If we let $\gamma \to \infty$, $\mathrm{prox}_{f_i}^\gamma(\boldsymbol{x})$ converges to the exact minimizer of $f_i$, irrelevant to the choice of $\boldsymbol{x}$.

[10]Such a way of partitioning is also inspired by [25] and similar to that in [36]. However, our setting is different from theirs.
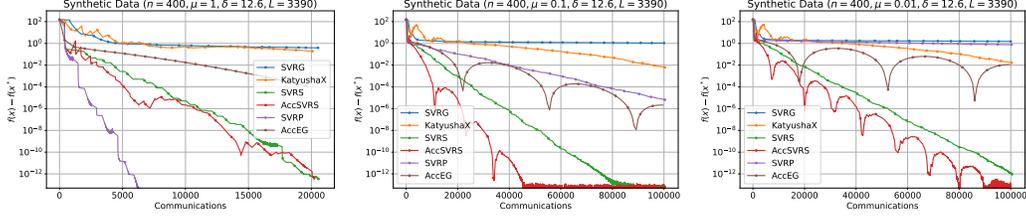
Figure 1: Numerical experiments on synthetic data. The corresponding coefficients are shown in the title of each graph. We plot the function gap on a log scale versus the number of communication steps, where one exchange of vectors counts as a communication step.

Lemma 4.3 guarantees that in each round, only when a specific component is sampled or the full gradient is calculated, can we expand the information set by at most three dimensions. For problem (13), we could never obtain an approximate solution unless we expand the information set to the whole space (see Proposition F.6 in Appendix F.2), while Lemma 4.3 implies that the process of expanding is very slow. Then we can establish the following lower bound.

**Theorem 4.4** *For any $n \geq 3$, $\delta, \mu > 0$, algorithm $\mathcal{A}$ satisfying Definition 4.1 and sufficiently small $\epsilon > 0$, there exists a rescaled version of problem (13) such that (i) Assumption 1 holds; (ii) In order to find an $\epsilon$-suboptimal solution $\hat{\boldsymbol{x}}$ such that $\mathbb{E}r(\hat{\boldsymbol{x}}) - \min_{\boldsymbol{x}} r(\boldsymbol{x}) < \epsilon$ by $\mathcal{A}$, the communication complexity in expectation is $\tilde{\Omega}(n + n^{3/4}\sqrt{\delta/\mu})$.*

This lower bound nearly matches the upper bound in Theorem 3.6 up to log factors, implying Algorithm 2 is nearly optimal in terms of communication complexity. The detailed statement and proof are deferred to Appendices F.2 and F.9.

## 5 Experiments

To demonstrate the advantages of our algorithms, we conduct the same numerical experiments as those in [35, 33]. We focus on the linear ridge regression problem with $\ell_2$ regularization, where the average loss $f$ has the formulation: $f(\boldsymbol{x}) = \frac{1}{n}\sum_{i=1}^{n}\left[f_i(\boldsymbol{x}) := \frac{1}{m}\sum_{j=1}^{m}\left(\boldsymbol{z}_{i,j}^\top\boldsymbol{x} - y_{i,j}\right)^2 + \frac{\mu}{2}\|\boldsymbol{x}\|^2\right]$. Here $\boldsymbol{z}_{i,j} \in \mathbb{R}^d$ and $y_{i,j} \in \mathbb{R}, \forall i \in [n], j \in [m]$ serve as the feature and label respectively, and $m$ can be viewed as data size in each local client. We consider a synthetic dataset generated by adding a small random noise matrix to the center matrix, ensuring a small $\delta$. To capture the differences in convergence rates between our methods and SVRP caused by different magnitudes of $\mu$, we vary $\mu = 10^{-i}, i \in \{0, 1, 2\}$. We compare our methods (SVRS and AccSVRS) against SVRG, KatyushaX, SVRP (Catalyzed SVRP is somehow hard to tune so we omit it), and Accelerated Extragradient (AccEG) using their theoretical step sizes, except that we scale the interpolation parameter $\tau$ in KatyushaX and AccSVRS for producing practical performance (see Appendix G for detail). From Figure 1, we can observe that for a large $\mu$, SVRP outperforms existing algorithms due to its high-order dependence on $\mu$. However, when the problem becomes ill-conditioned with a small $\mu$, AccSVRS exhibits significant improvements compared to other algorithms.

## 6 Conclusion

In this paper, we have introduced two new algorithms, SVRS and its directly accelerated version AccSVRS, and established improved communication complexity bounds for distributed optimization under the similarity assumption. Our rates are entirely smoothness-free and only require strong convexity of the objective, average similarity, and proximal friendliness of components. Moreover, our methods also have nearly optimal gradient complexity (leaving out the log term) when applied to smooth components in specific cases. It would be interesting to remove additional log terms to achieve both optimal communication and local gradient calls as [35], as well as investigating the complexity under other similarity assumptions (such as SS instead of AveSS) in future research.

## Acknowledgments and Disclosure of Funding

## References

[1] Artem Agafonov, Pavel Dvurechensky, Gesualdo Scutari, Alexander Gasnikov, Dmitry Kamzolov, Aleksandr Lukashevich, and Amir Daneshmand. An accelerated second-order method for distributed stochastic optimization. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pages 2407–2413. IEEE, 2021.

[2] Alekh Agarwal and Leon Bottou. A lower bound for the optimization of finite sums. In *ICML*, 2015.

[3] Masoud Ahookhosh and Yurii Nesterov. High-order methods beyond the classical complexity bounds, ii: inexact high-order proximal-point methods with segment search. *arXiv preprint arXiv:2109.12303*, 2021.

[4] Mohammad S Alkousa, Alexander Vladimirovich Gasnikov, Darina Mikhailovna Dvinskikh, Dmitry A Kovalev, and Fedor Sergeevich Stonyakin. Accelerated methods for saddle-point problem. *Computational Mathematics and Mathematical Physics*, 60:1787–1809, 2020.

[5] Zeyuan Allen-Zhu. Katyusha: the first direct acceleration of stochastic gradient methods. *The Journal of Machine Learning Research*, 18(1):8194–8244, 2017.

[6] Zeyuan Allen-Zhu. Katyusha x: Practical momentum method for stochastic sum-of-nonconvex optimization. *arXiv preprint arXiv:1802.03866*, 2018.

[7] Yossi Arjevani and Ohad Shamir. Communication complexity of distributed convex learning and optimization. *Advances in neural information processing systems*, 28, 2015.

[8] Syreen Banabilah, Moayad Aloqaily, Eitaa Alsayed, Nida Malik, and Yaser Jararweh. Federated learning review: Fundamentals, enabling technologies, and future applications. *Information processing & management*, 59(6):103061, 2022.

[9] Ron Bekkerman, Mikhail Bilenko, and John Langford. *Scaling up machine learning: Parallel and distributed approaches*. Cambridge University Press, 2011.

[10] Aleksandr Beznosikov and Alexander Gasnikov. Compression and data similarity: Combination of two techniques for communication-efficient solving of distributed variational inequalities. In *International Conference on Optimization and Applications*, pages 151–162. Springer, 2022.

[11] Aleksandr Beznosikov and Alexander Gasnikov. Similarity, compression and local steps: Three pillars of efficient communications for distributed variational inequalities. *arXiv preprint arXiv:2302.07615*, 2023.

[12] Aleksandr Beznosikov, Eduard Gorbunov, and Alexander Gasnikov. Derivative-free method for composite optimization with applications to decentralized distributed optimization. *IFAC-PapersOnLine*, 53(2):4038–4043, 2020.

[13] Aleksandr Beznosikov, Gesualdo Scutari, Alexander Rogozin, and Alexander Gasnikov. Distributed saddle-point problems under data similarity. *Advances in Neural Information Processing Systems*, 34:8172–8184, 2021.

[14] Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Accelerated methods for nonconvex optimization. *SIAM Journal on Optimization*, 28(2):1751–1772, 2018.

[15] Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary points ii: first-order methods. *Mathematical Programming*, 185(1-2):315–355, 2021.

[16] Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.

[17] Gong Chen and Marc Teboulle. Convergence analysis of a proximal-like minimization algorithm using bregman functions. *SIAM Journal on Optimization*, 3(3):538–543, 1993.

[18] Amir Daneshmand, Gesualdo Scutari, Pavel Dvurechensky, and Alexander Gasnikov. Newton method over networks is fast up to the statistical precision. In *International Conference on Machine Learning*, pages 2398–2409. PMLR, 2021.

[19] Aaron Defazio. A simple practical accelerated method for finite sums. *Advances in neural information processing systems*, 29, 2016.

[20] Radu Alexandru Dragomir, Mathieu Even, and Hadrien Hendrikx. Fast stochastic bregman gradient methods: Sharp analysis and variance reduction. In *International Conference on Machine Learning*, pages 2815–2825. PMLR, 2021.

[21] Darina Mikhailovna Dvinskikh, Sergey Sergeevich Omelchenko, Alexander Vladimirovich Gasnikov, and AI Tyurin. Accelerated gradient sliding for minimizing a sum of functions. In *Doklady Mathematics*, volume 101, pages 244–246. Springer, 2020.

[22] Alexandre d'Aspremont, Damien Scieur, Adrien Taylor, et al. Acceleration methods. *Foundations and Trends® in Optimization*, 5(1-2):1–245, 2021.

[23] Dan Garber, Elad Hazan, Chi Jin, Cameron Musco, Praneeth Netrapalli, Aaron Sidford, et al. Faster eigenvector computation via shift-and-invert preconditioning. In *International Conference on Machine Learning*, pages 2626–2634. PMLR, 2016.

[24] Geovani Nunes Grapiglia and Yurii Nesterov. Adaptive third-order methods for composite convex optimization. *arXiv preprint arXiv:2202.12730*, 2022.

[25] Yuze Han, Guangzeng Xie, and Zhihua Zhang. Lower complexity bounds of finite-sum optimization problems: The results and construction. *arXiv preprint arXiv:2103.08280*, 2021.

[26] Robert Hannah, Yanli Liu, Daniel O'Connor, and Wotao Yin. Breaking the span assumption yields fast finite-sum minimization. *Advances in Neural Information Processing Systems*, 31, 2018.

[27] Hadrien Hendrikx, Lin Xiao, Sebastien Bubeck, Francis Bach, and Laurent Massoulie. Statistically preconditioned accelerated gradient method for distributed optimization. In *International conference on machine learning*, pages 4203–4227. PMLR, 2020.

[28] Anastasiya Ivanova, Pavel Dvurechensky, Evgeniya Vorontsova, Dmitry Pasechnyuk, Alexander Gasnikov, Darina Dvinskikh, and Alexander Tyurin. Oracle complexity separation in convex optimization. *Journal of Optimization Theory and Applications*, 193(1-3):462–490, 2022.

[29] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. *Advances in neural information processing systems*, 26, 2013.

[30] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.

[31] Dmitry Kamzolov, Alexander Gasnikov, and Pavel Dvurechensky. Optimal combination of tensor optimization methods. In *Optimization and Applications: 11th International Conference, OPTIMA 2020, Moscow, Russia, September 28–October 2, 2020, Proceedings 11*, pages 166–183. Springer, 2020.

[32] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.

[33] Ahmed Khaled and Chi Jin. Faster federated optimization under second-order similarity. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=ElC6LYO4MfD.

[34] Dmitry Kovalev, Samuel Horváth, and Peter Richtárik. Don't jump through hoops and remove those loops: Svrg and katyusha are better without the outer loop. In *Algorithmic Learning Theory*, pages 451–467. PMLR, 2020.

[35] Dmitry Kovalev, Aleksandr Beznosikov, Ekaterina Dmitrievna Borodich, Alexander Gasnikov, and Gesualdo Scutari. Optimal gradient sliding and its application to optimal distributed optimization under similarity. In *Advances in Neural Information Processing Systems*, 2022.

[36] Dmitry Kovalev, Aleksandr Beznosikov, Abdurakhmon Sadiev, Michael Persiianov, Peter Richtárik, and Alexander Gasnikov. Optimal algorithms for decentralized stochastic variational inequalities. *Advances in Neural Information Processing Systems*, 35:31073–31088, 2022.

[37] Guanghui Lan. Gradient sliding for composite optimization. *Mathematical Programming*, 159:201–235, 2016.

[38] Guanghui Lan and Yuyuan Ouyang. Accelerated gradient sliding for structured convex optimization. *Computational Optimization and Applications*, 82(2):361–394, 2022.

[39] Guanghui Lan and Yi Zhou. Conditional gradient sliding for convex optimization. *SIAM Journal on Optimization*, 26(2):1379–1409, 2016.

[40] Guanghui Lan and Yi Zhou. An optimal randomized incremental gradient method. *Mathematical programming*, 171:167–215, 2018.

[41] Kfir Y Levy. Slowcal-sgd: Slow query points improve local-sgd for stochastic convex optimization. *arXiv preprint arXiv:2304.04169*, 2023.

[42] Bingcong Li, Meng Ma, and Georgios B Giannakis. On the convergence of sarah and beyond. In *International Conference on Artificial Intelligence and Statistics*, pages 223–233. PMLR, 2020.

[43] Zhize Li. Anita: An optimal loopless accelerated variance-reduced gradient method. *arXiv preprint arXiv:2103.11333*, 2021.

[44] 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. *Advances in neural information processing systems*, 30, 2017.

[45] Ming Liu, Stella Ho, Mengqi Wang, Longxiang Gao, Yuan Jin, and He Zhang. Federated learning meets natural language processing: A survey. *arXiv preprint arXiv:2107.12603*, 2021.

[46] Konstantin Mishchenko, Grigory Malinovsky, Sebastian Stich, and Peter Richtárik. Proxskip: Yes! local gradient steps provably lead to communication acceleration! finally! In *International Conference on Machine Learning*, pages 15750–15769. PMLR, 2022.

[47] Yurii Nesterov et al. *Lectures on convex optimization*, volume 137. Springer, 2018.

[48] Dinh C Nguyen, Quoc-Viet Pham, Pubudu N Pathirana, Ming Ding, Aruna Seneviratne, Zihuai Lin, Octavia Dobre, and Won-Joo Hwang. Federated learning for smart healthcare: A survey. *ACM Computing Surveys (CSUR)*, 55(3):1–37, 2022.

[49] Yuyuan Ouyang and Yangyang Xu. Lower complexity bounds of first-order methods for convex-concave bilinear saddle-point problems. *Mathematical Programming*, 185(1-2):1–35, 2021.

[50] Sashank J Reddi, Jakub Konečnỳ, Peter Richtárik, Barnabás Póczós, and Alex Smola. Aide: Fast and communication efficient distributed optimization. *arXiv preprint arXiv:1608.06879*, 2016.

[51] Ohad Shamir, Nati Srebro, and Tong Zhang. Communication-efficient distributed optimization using an approximate newton-type method. In *International conference on machine learning*, pages 1000–1008. PMLR, 2014.

[52] Danny C Sorensen. Numerical methods for large eigenvalue problems. *Acta Numerica*, 11: 519–584, 2002.

[53] Ivan Stepanov, Artyom Voronov, Aleksandr Beznosikov, and Alexander Gasnikov. One-point gradient-free methods for composite optimization with applications to distributed optimization. *arXiv preprint arXiv:2107.05951*, 2021.

[54] Ying Sun, Gesualdo Scutari, and Amir Daneshmand. Distributed optimization based on gradient tracking revisited: Enhancing convergence rate via surrogation. *SIAM Journal on Optimization*, 32(2):354–385, 2022.

[55] Ye Tian, Gesualdo Scutari, Tianyu Cao, and Alexander Gasnikov. Acceleration in distributed optimization under similarity. In *International Conference on Artificial Intelligence and Statistics*, pages 5721–5756. PMLR, 2022.

[56] Blake E Woodworth and Nati Srebro. Tight complexity bounds for optimizing composite objectives. *Advances in neural information processing systems*, 29, 2016.

[57] Lin Xiao and Tong Zhang. A proximal stochastic gradient method with progressive variance reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.

[58] Xiao-Tong Yuan and Ping Li. On convergence of distributed approximate newton methods: Globalization, sharper bounds and beyond. *The Journal of Machine Learning Research*, 21(1): 8502–8552, 2020.

[59] Junyu Zhang, Mingyi Hong, and Shuzhong Zhang. On lower iteration complexity bounds for the convex concave saddle point problems. *Mathematical Programming*, 194(1-2):901–935, 2022.

[60] Sai Qian Zhang, Jieyu Lin, and Qi Zhang. A multi-agent reinforcement learning approach for efficient client selection in federated learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 9091–9099, 2022.

[61] Siqi Zhang, Junchi Yang, Cristóbal Guzmán, Negar Kiyavash, and Niao He. The complexity of nonconvex-strongly-concave minimax optimization. In *Uncertainty in Artificial Intelligence*, pages 482–492. PMLR, 2021.

[62] Yuchen Zhang and Xiao Lin. Disco: Distributed optimization for self-concordant empirical loss. In *International conference on machine learning*, pages 362–370. PMLR, 2015.

[63] Dongruo Zhou and Quanquan Gu. Lower bounds for smooth nonconvex finite-sum optimization. In *International Conference on Machine Learning*, pages 7574–7583. PMLR, 2019.

# A  Auxiliary Results

**Proposition A.1 (Three-point identity [17, Lemma3.1])** *Given a differentiable function $h\colon \mathbb{R}^d \to \mathbb{R}$, we have the following equality:*

$$\langle \boldsymbol{x} - \boldsymbol{y}, \nabla h(\boldsymbol{y}) - \nabla h(\boldsymbol{z})\rangle = D_h(\boldsymbol{x}, \boldsymbol{z}) - D_h(\boldsymbol{x}, \boldsymbol{y}) - D_h(\boldsymbol{y}, \boldsymbol{z}), \forall \boldsymbol{x}, \boldsymbol{y}, \boldsymbol{z} \in \mathbb{R}^d. \qquad (14)$$

**Proposition A.2** *Denote $\forall i \in \mathbb{N}, X_i = \begin{cases} 1 & \text{with probability} \quad p \\ 0 & \text{with probability} \quad 1 - p \end{cases}$, and $X_1, X_2, \ldots$ are independent and identically distributed random variables. Then $Y := \inf_i\{i : X_i = 1\} \sim \mathrm{Geom}(p)$.*

Proof: We direct verify the probability distribution:

$$\mathbb{P}(Y = k) = \prod_{i=1}^{k-1} \mathbb{P}(X_i = 0)\mathbb{P}(X_k = 1) = (1 - p)^{k-1}p, \quad k \in \{1, 2, \ldots\}.$$

Hence, we see that $Y \sim \mathrm{Geom}(p)$. $\qquad\qquad \square$

**Proposition A.3 (Proposition 2.5 in the main text)** *We have the following properties among SS, AveSS, and SC: 1) The $\delta$-SS can deduce $\delta$-AveSS, but $\delta$-AveSS can only deduce $\sqrt{n}\delta$-SS. 2) If $f_i$'s satisfy $\delta$-SS and $f$ is $\mu$-strongly convex, then for all $i \in [n]$, $f_i(\cdot) + \frac{\delta-\mu}{2}\|\cdot\|^2$ is convex, i.e., $f_i$ is $(\delta - \mu)$-almost convex [14].*

Proof: 1) The first part "$\delta$-SS $\Rightarrow$ $\delta$-AveSS" is trivial. The second part is because for all $i \in [n]$,

$$\|[\nabla[f_i - f](\boldsymbol{x}) - \nabla[f_i - f](\boldsymbol{y})]\|^2 \leq \sum_{j=1}^{n} \|[\nabla[f_j - f](\boldsymbol{x}) - \nabla[f_j - f](\boldsymbol{y})]\|^2 \overset{(4)}{\leq} n\delta^2 \|\boldsymbol{x} - \boldsymbol{y}\|^2.$$

Thus Eq. (5) holds with parameter $\sqrt{n}\delta$.

2) Since $f_i$'s satisfy $\delta$-SS, we get $\forall i \in [n], f - f_i$ is $\delta$-smooth, thus $\frac{\delta}{2}\|\boldsymbol{x}\|^2 - [f(\boldsymbol{x}) - f_i(\boldsymbol{x})]$ is convex (e.g., [22, Theorem A.1]). Moreover, we also have $f(\boldsymbol{x}) - \frac{\mu}{2}\|\boldsymbol{x}\|^2$ is a convex function since $f$ is $\mu$-strongly convex (e.g., [22, Theorem A.2]). Therefore, we obtain that

$$f_i(\boldsymbol{x}) + \frac{\delta - \mu}{2}\|\boldsymbol{x}\|^2 = \left(\frac{\delta}{2}\|\boldsymbol{x}\|^2 - [f(\boldsymbol{x}) - f_i(\boldsymbol{x})]\right) + \left(f(\boldsymbol{x}) - \frac{\mu}{2}\|\boldsymbol{x}\|^2\right)$$

is also convex. The proof is finished. $\qquad\qquad \square$

**Lemma A.4 (Allen-Zhu [6, Fact 2.3])** *Given sequence $D_0, D_1, \ldots$ of reals, if $N \sim \mathrm{Geom}(p)$, then*

$$\mathbb{E}_N[D_{N-1} - D_N] = p\mathbb{E}[D_0 - D_N], \mathbb{E}_N[D_{N-1}] = (1 - p)\mathbb{E}[D_N] + pD_0 \qquad (15)$$

**Lemma A.5 (Allen-Zhu [6, Lemma 2.4])** *If $g(\cdot)$ is proper convex and $\sigma$-strongly convex and $\boldsymbol{z}_{k+1} = \arg\min_{\boldsymbol{z} \in \mathbb{R}^d} \frac{1}{2\alpha}\|\boldsymbol{z} - \boldsymbol{z}_k\|^2 + \langle \boldsymbol{\xi}, \boldsymbol{z}\rangle + g(\boldsymbol{z})$, then for every $\boldsymbol{x} \in \mathbb{R}^d$, we have*

$$\langle \boldsymbol{\xi}, \boldsymbol{z}_k - \boldsymbol{x}\rangle + g(\boldsymbol{z}_{k+1}) - g(\boldsymbol{x}) \leq \frac{\alpha}{2}\|\boldsymbol{\xi}\|^2 + \frac{\|\boldsymbol{x} - \boldsymbol{z}_k\|^2}{2\alpha} - \frac{(1 + \sigma\alpha)\|\boldsymbol{x} - \boldsymbol{z}_{k+1}\|^2}{2\alpha}. \qquad (16)$$

**Lemma A.6 (Han et al. [25, Lemma 2.10])** *Let $\{Y_i\}_{i=1}^m$ be independent random variables such that $Y_i \sim \mathrm{Geom}(p_i)$ with $p_i > 0$. Then for $m \geq 2$, we have*

$$\mathbb{P}\left(\sum_{i=1}^m Y_i > \frac{m^2}{4(\sum_{i=1}^m p_i)}\right) \geq \frac{1}{9}.$$

# B  Hessian Similarity

In this section, we show that AveHS (HS) defined in Eq. (6) is equivalent to AveSS (SS).

**Proposition B.1** *For twice differentiability $f_i$'s and $f$, AveSS $\Leftrightarrow$ AveHS, SS $\Leftrightarrow$ HS.*

Proof: Indeed, we only need to prove the following results for twice differentiability $g$:

$$\frac{1}{n}\sum_{i=1}^{n}\|\nabla g_i(\boldsymbol{x}) - \nabla g_i(\boldsymbol{y})\|^2 \le \delta^2 \|\boldsymbol{y} - \boldsymbol{x}\|^2 \Leftrightarrow \left\|\frac{1}{n}\sum_{i=1}^{n}\left(\nabla^2 g_i(\boldsymbol{x})\right)^2\right\| \le \delta^2, \forall \boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^d. \quad (17)$$

"⇒": Taking $\boldsymbol{y} = \boldsymbol{x} + t\boldsymbol{v}, t \in \mathbb{R}\backslash\{0\}, \boldsymbol{v} \in \mathbb{R}^d, \|\boldsymbol{v}\| = 1$ and letting $t \to 0$, we get

$$
\begin{aligned}
\delta^2 &= \lim_{t\to 0}\frac{\delta^2\|\boldsymbol{x}-\boldsymbol{y}\|^2}{t^2} \ge \lim_{t\to 0}\frac{1}{n}\sum_{i=1}^{n}\left\|\frac{[\nabla g_i(\boldsymbol{x}) - \nabla g_i(\boldsymbol{x}+t\boldsymbol{v})]}{t}\right\|^2\\
&= \frac{1}{n}\sum_{i=1}^{n}\left\|\nabla^2 g_i(\boldsymbol{x})\boldsymbol{v}\right\|^2 = \boldsymbol{v}^\top\left[\frac{1}{n}\sum_{i=1}^{n}\left(\nabla^2 g_i(\boldsymbol{x})\right)^2\right]\boldsymbol{v}.
\end{aligned}
$$

The final equality uses the fact that $\nabla^2 g_i(\boldsymbol{x})$ is a symmetric matrix. Now by the arbitrary of $\boldsymbol{v} \in \mathbb{R}^d$ with $\|\boldsymbol{v}\| = 1$, we get $\left\|\frac{1}{n}\sum_{i=1}^{n}\left(\nabla^2 g_i(\boldsymbol{x})\right)^2\right\| \le \delta^2$.

"⇐": We use the integral formulation:

$$
\begin{aligned}
\frac{1}{n}\sum_{i=1}^{n}\|\nabla g_i(\boldsymbol{x}) - \nabla g_i(\boldsymbol{y})\|^2 &= \frac{1}{n}\sum_{i=1}^{n}\left\|\int_0^1\nabla^2 g_i(\boldsymbol{x}+t(\boldsymbol{y}-\boldsymbol{x}))(\boldsymbol{y}-\boldsymbol{x})\,dt\right\|^2\\
&= (\boldsymbol{y}-\boldsymbol{x})\left[\frac{1}{n}\sum_{i=1}^{n}\int_0^1\nabla^2 g_i(\boldsymbol{x}+s(\boldsymbol{y}-\boldsymbol{x}))\nabla^2 g_i(\boldsymbol{x}+t(\boldsymbol{y}-\boldsymbol{x}))ds dt\right](\boldsymbol{y}-\boldsymbol{x})\\
&\overset{(i)}{\le} (\boldsymbol{y}-\boldsymbol{x})\left[\frac{1}{n}\sum_{i=1}^{n}\int_0^1\left(\nabla^2 g_i(\boldsymbol{x}+t(\boldsymbol{y}-\boldsymbol{x}))\right)^2 dt\right](\boldsymbol{y}-\boldsymbol{x})\\
&\le \int_0^1\left\|\frac{1}{n}\sum_{i=1}^{n}\left(\nabla^2 g_i(\boldsymbol{x}+t(\boldsymbol{y}-\boldsymbol{x}))\right)^2\right\|\cdot\|\boldsymbol{y}-\boldsymbol{x}\|^2 dt \le \delta^2\|\boldsymbol{y}-\boldsymbol{x}\|^2,
\end{aligned}
$$

where $(i)$ uses the inequality $\boldsymbol{A}_t^2 + \boldsymbol{A}_s^2 \succeq \boldsymbol{A}_s\boldsymbol{A}_t + \boldsymbol{A}_t\boldsymbol{A}_s$ for symmetric matrices $\boldsymbol{A}_s, \forall s \in [0,1]$ since $(\boldsymbol{A}_t - \boldsymbol{A}_s)^2 \succeq 0$, and the final inequality uses the assumption.

Hence, Eq. (17) is proved. Now choosing $g_i = f_i - f, \forall i \in [n]$, we obtain "AveSS ⇔ AveHS". Additionally, letting $n = 1$ and noting that $\left\|\left(\nabla^2 g_i(\boldsymbol{x})\right)^2\right\| = \left\|\nabla^2 g_i(\boldsymbol{x})\right\|^2$, we obtain "SS ⇔ HS". The proof is finished. □

## C   Concrete Complexity of Catalyst SVRP

Inherited from the computation of [33, Appendix G.1], we see that the total iterations of Catalyst SVRP is

$$
\begin{aligned}
\mathbb{E}T_{\text{iter}}^{\text{total}} &= 8\sqrt{\frac{\mu+\gamma}{\mu}}\max\left\{\frac{\delta^2}{(\gamma+\mu)^2}, n\right\}\log\left(\frac{f(\boldsymbol{x}_0)-f(\boldsymbol{x}_*)}{\varepsilon}\cdot\frac{32(\mu+\gamma)}{\mu}\right)\log\iota,\\
\iota &:= A\left(\frac{2}{1-\rho}+\frac{2592\gamma}{\mu(1-\rho)^2(\sqrt{q}-\rho)^2}\right),
\end{aligned}
$$

where $\rho = \sqrt{q}/2 = \frac{\sqrt{\mu/(\mu+\gamma)}}{2} \in (0, \frac{1}{2})$, $A = \frac{L+\gamma}{\mu+\gamma}\left(1+\frac{(\gamma+\mu)^2 n}{\delta^2}\right)$. Letting $\gamma = \max\left\{\frac{\delta}{\sqrt{n}}-\mu, 0\right\}$, we recover the complexity:

$$
\begin{aligned}
\mathbb{E}T_{\text{iter}}^{\text{total}} &= 8\max\left\{n, n^{3/4}\sqrt{\frac{\delta}{\mu}}\right\}\log\left(\max\left\{32, \frac{32\delta}{\mu\sqrt{n}}\right\}\cdot\frac{f(\boldsymbol{x}_0)-f(\boldsymbol{x}_*)}{\varepsilon}\right)\log\iota,\\
\iota &= A\left(\frac{2}{1-\rho}+\frac{2592\gamma}{\mu(1-\rho)^2(\sqrt{q}-\rho)^2}\right) = \Theta\left(A\left(1+\frac{\gamma(\mu+\gamma)}{\mu^2}\right)\right) = \Theta\left(\frac{A(\mu+\gamma)^2}{\mu^2}\right).
\end{aligned}
$$

When $\delta/\mu \le \sqrt{n}$, leading to $\gamma = 0$, then we get

$$\frac{A(\mu+\gamma)^2}{\mu^2} = \frac{L}{\mu}\left(1+\frac{\mu^2 n}{\delta^2}\right) = \Theta\left(\frac{L\mu n}{\delta^2}\right).$$

---

**Algorithm 3** Stochastic Variance-Reduced Sliding (SVRS)

---

1: **Input:** $\boldsymbol{w}_0 \in \mathbb{R}^d, p \in (0,1), \theta > 0, K \in \{1,2,\dots\}$
2: **for** $k = 0, 1, 2, \dots, K-1$ **do**
3: $\quad \boldsymbol{w}_{k+1} = \text{SVRS}^{1\text{ep}}(f, \boldsymbol{w}_k, \theta, p)$
4: **end for**
5: **Output:** $\boldsymbol{w}_K$

---

Thus, $\mathbb{E}T_{\text{iter}}^{\text{total}} = \mathcal{O}\left( \left( n + n^{3/4}\sqrt{\frac{\delta}{\mu}} \right) \log \frac{f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*)}{\varepsilon} \log \frac{L\mu n}{\delta^2} \right)$.

When $\delta/\mu \geq \sqrt{n}$, i.e., $\max\left\{ n, n^{3/4}\sqrt{\frac{\delta}{\mu}} \right\} = n^{3/4}\sqrt{\frac{\delta}{\mu}}$, we get $\gamma = \frac{\delta}{\sqrt{n}} - \mu \leq L - \mu$ (note that $L \geq \delta \geq \mu, n \geq 1$ by assumption), leading to

$$\frac{2L}{\mu} \leq \frac{A(\mu + \gamma)^2}{\mu^2} = \frac{2(L+\gamma)(\mu+\gamma)}{\mu^2} \leq \frac{4L^2}{\mu^2}.$$

Thus, $\mathbb{E}T_{\text{iter}}^{\text{total}} = \mathcal{O}\left( \left( n + n^{3/4}\sqrt{\frac{\delta}{\mu}} \right) \log \frac{f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*)}{\varepsilon} \log \frac{L}{\mu} \right)$ (for small enough error $\varepsilon$).

## D  Proofs for Section 3

The complete procedure of SVRS is presented in Algorithm 3. Before giving the omit proofs, we need the following one-step lemma.

**Lemma D.1** *Suppose Assumption 1 holds. If the step size $\theta \leq 1/(2\sqrt{n}\delta)$ in $\text{SVRS}^{1\text{ep}}$(Algorithm 1), then the following inequality holds for all $\boldsymbol{x} \in \mathbb{R}^d$ that is independent to the index $i_k$:*

$$\mathbb{E}_t[f(\boldsymbol{x}_{t+1}) - f(\boldsymbol{x})] \leq \mathbb{E}_t D_h(\boldsymbol{x}, \boldsymbol{x}_t) - \left( 1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta} \right) D_h(\boldsymbol{x}, \boldsymbol{x}_{t+1}) + \frac{2\theta^2\delta^2}{(1 - \sqrt{n}\theta\delta)^2} D_h(\boldsymbol{w}_0, \boldsymbol{x}_t)$$

$$+ \frac{2 + \mu\theta}{2\mu} \left[ \left\| \nabla A_\theta^t(\boldsymbol{x}_{t+1}) \right\|^2 - \frac{\mu}{20\theta} \left\| \boldsymbol{x}_t - \arg\min_{\boldsymbol{x} \in \mathbb{R}^d} A_\theta^t(\boldsymbol{x}) \right\|^2 \right]. \qquad (18)$$

Proof: First, note that

$$\nabla A_\theta^t(\boldsymbol{x}) = \frac{\boldsymbol{x} - \boldsymbol{x}_t}{\theta} + \nabla f_1(\boldsymbol{x}) - \nabla f_1(\boldsymbol{x}_t) + \nabla f_{i_t}(\boldsymbol{x}_t) - [\nabla f_{i_t}(\boldsymbol{w}_0) - \nabla f(\boldsymbol{w}_0)]$$
$$= \nabla f(\boldsymbol{x}) + \nabla h(\boldsymbol{x}) - \nabla h(\boldsymbol{x}_t) + \nabla(f_{i_t} - f)(\boldsymbol{x}_t) - \nabla(f_{i_t} - f)(\boldsymbol{w}_0). \qquad (19)$$

Now we begin from the strong convexity of function $f$ in Assumption 1,

$$\mathbb{E}_t[f(\boldsymbol{x}_{t+1}) - f(\boldsymbol{x})] \overset{(2)}{\leq} \mathbb{E}_t \langle \boldsymbol{x} - \boldsymbol{x}_{t+1}, -\nabla f(\boldsymbol{x}_{t+1}) \rangle - \frac{\mu}{2} \| \boldsymbol{x}_{t+1} - \boldsymbol{x} \|^2$$

$$\overset{(19)}{=} \mathbb{E}_t \langle \boldsymbol{x} - \boldsymbol{x}_{t+1}, \nabla h(\boldsymbol{x}_{t+1}) - \nabla h(\boldsymbol{x}_t) \rangle + \langle \boldsymbol{x} - \boldsymbol{x}_{t+1}, \nabla(f_{i_t} - f)(\boldsymbol{x}_t) - \nabla(f_{i_t} - f)(\boldsymbol{w}_0) \rangle$$
$$- \langle \boldsymbol{x} - \boldsymbol{x}_{t+1}, \nabla A_\theta^t(\boldsymbol{x}_{t+1}) \rangle - \frac{\mu}{2} \| \boldsymbol{x} - \boldsymbol{x}_{t+1} \|^2$$

$$\overset{(i)}{=} \mathbb{E}_t D_h(\boldsymbol{x}, \boldsymbol{x}_t) - D_h(\boldsymbol{x}, \boldsymbol{x}_{t+1}) - D_h(\boldsymbol{x}_{t+1}, \boldsymbol{x}_t) + \langle \boldsymbol{x}_t - \boldsymbol{x}_{t+1}, \nabla(f_{i_t} - f)(\boldsymbol{x}_t) - \nabla(f_{i_t} - f)(\boldsymbol{w}_0) \rangle$$
$$- \langle \boldsymbol{x} - \boldsymbol{x}_{t+1}, \nabla A_\theta^t(\boldsymbol{x}_{t+1}) \rangle - \frac{\mu}{2} \| \boldsymbol{x}_{t+1} - \boldsymbol{x} \|^2$$

$$\leq \mathbb{E}_t D_h(\boldsymbol{x}, \boldsymbol{x}_t) - D_h(\boldsymbol{x}, \boldsymbol{x}_{t+1}) - D_h(\boldsymbol{x}_{t+1}, \boldsymbol{x}_t)$$
$$+ \frac{1 - \sqrt{n}\theta\delta}{4\theta} \| \boldsymbol{x}_{t+1} - \boldsymbol{x}_t \|^2 + \frac{\theta}{1 - \sqrt{n}\theta\delta} \| \nabla(f_{i_t} - f)(\boldsymbol{x}_t) - \nabla(f_{i_t} - f)(\boldsymbol{w}_0) \|^2$$
$$+ \left[ \frac{\mu}{4} \| \boldsymbol{x}_{t+1} - \boldsymbol{x} \|^2 + \frac{1}{\mu} \| \nabla A_\theta^t(\boldsymbol{x}_{t+1}) \|^2 \right] - \frac{\mu}{2} \| \boldsymbol{x}_{t+1} - \boldsymbol{x} \|^2 \qquad (20)$$

where $(i)$ uses Eq. (14) and $\mathbb{E}_{i_t} \langle \boldsymbol{x}_t - \boldsymbol{x}, \nabla(f_{i_t} - f)(\boldsymbol{x}_t) - \nabla(f_{i_t} - f)(\boldsymbol{w}_0) \rangle = 0$ since $\boldsymbol{x}_t - \boldsymbol{x}$ is independent to $i_t$ and the final inequality uses $\langle \boldsymbol{a}, \boldsymbol{b} \rangle \leq t^2 \| \boldsymbol{a} \|^2 + \frac{\|\boldsymbol{b}\|^2}{4t^2}$ twice. Next, we continue

17

using Eq. (9) to convert $\|\cdot\|$ with $D_h(\cdot, \cdot)$ by assumption $\theta \le 1/(2\sqrt{n}\delta)$:

$$\mathbb{E}_t[f(\boldsymbol{x}_{t+1}) - f(\boldsymbol{x})]$$

$$\stackrel{(20)(9)}{\le} \mathbb{E}_t D_h(\boldsymbol{x}, \boldsymbol{x}_t) - \left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{x}, \boldsymbol{x}_{t+1}) - \frac{1 - \sqrt{n}\theta\delta}{4\theta} \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_t\|^2$$

$$+ \frac{\theta}{1 - \sqrt{n}\theta\delta} \|\nabla(f_{i_t} - f)(\boldsymbol{x}_t) - \nabla(f_{i_t} - f)(\boldsymbol{w}_0)\|^2 + \frac{1}{\mu} \left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2$$

$$\stackrel{(4)}{\le} \mathbb{E}_t D_h(\boldsymbol{x}, \boldsymbol{x}_t) - \left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{x}, \boldsymbol{x}_{t+1}) + \frac{\theta\delta^2}{1 - \sqrt{n}\theta\delta} \|\boldsymbol{x}_t - \boldsymbol{w}_0\|^2$$

$$- \frac{1 - \sqrt{n}\theta\delta}{4\theta} \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_t\|^2 + \frac{1}{\mu} \left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2$$

$$\stackrel{(9)}{\le} \mathbb{E}_t D_h(\boldsymbol{x}, \boldsymbol{x}_t) - \left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{x}, \boldsymbol{x}_{t+1}) + \frac{2\theta^2\delta^2}{(1 - \sqrt{n}\theta\delta)^2} D_h(\boldsymbol{w}_0, \boldsymbol{x}_t)$$

$$- \frac{1 - \sqrt{n}\theta\delta}{4\theta} \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_t\|^2 + \frac{1}{\mu} \left\|\nabla A_\theta^t(x_{t+1})\right\|^2 .$$

Finally, we show the error analysis if an approximate solution, i.e., $\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\| \ne 0$ is allowed. Using Proposition 2.5, we see that $f_1(\boldsymbol{x}) + \frac{\sqrt{n}\delta - \mu}{2}\|\boldsymbol{x}\|^2$ is a convex function, leading to $A_\theta^t(\boldsymbol{x})$ is $\left(\frac{1}{\theta} - \sqrt{n}\delta + \mu\right)$-strongly convex function. Let $\hat{\boldsymbol{x}}_{k+1} \in \arg\min_{\boldsymbol{x} \in \mathbb{R}^d} A_\theta^t(\boldsymbol{x})$, i.e., $\nabla A_\theta^t(\hat{\boldsymbol{x}}_{k+1}) = 0$. Since $\theta \le 1/(2\sqrt{n}\delta)$, we can further bound the last two terms:

$$-\frac{1 - \sqrt{n}\theta\delta}{4\theta} \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_t\|^2 + \frac{1}{\mu} \left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2 \le \frac{1}{\mu} \left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2 - \frac{1}{8\theta} \|\boldsymbol{x}_{t+1} - \boldsymbol{x}_t\|^2$$

$$\le \frac{1}{\mu} \left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2 + \frac{1}{8\theta} \|\boldsymbol{x}_{t+1} - \hat{\boldsymbol{x}}_{t+1}\|^2 - \frac{1}{16\theta} \|\hat{\boldsymbol{x}}_{t+1} - \boldsymbol{x}_t\|^2$$

$$\le \frac{1}{\mu} \left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2 + \frac{\theta}{8(1 - (\sqrt{n}\delta - \mu)\theta)^2} \left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1}) - \nabla A_\theta^t(\hat{\boldsymbol{x}}_{t+1})\right\|^2 - \frac{1}{16\theta} \|\hat{\boldsymbol{x}}_{t+1} - \boldsymbol{x}_t\|^2$$

$$\le \frac{1}{\mu} \left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2 + \frac{\theta}{2} \left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2 - \frac{1}{16\theta} \|\boldsymbol{x}_t - \hat{\boldsymbol{x}}_{t+1}\|^2$$

$$= \frac{2 + \mu\theta}{2\mu} \left[\left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2 - \frac{\mu}{8\theta(2 + \mu\theta)} \left\|\boldsymbol{x}_t - \arg\min_{\boldsymbol{x} \in \mathbb{R}^d} A_\theta^t(\boldsymbol{x})\right\|^2\right]$$

$$\le \frac{2 + \mu\theta}{2\mu} \left[\left\|\nabla A_\theta^t(\boldsymbol{x}_{t+1})\right\|^2 - \frac{\mu}{20\theta} \left\|\boldsymbol{x}_t - \arg\min_{\boldsymbol{x} \in \mathbb{R}^d} A_\theta^t(\boldsymbol{x})\right\|^2\right].$$

Therefore, Eq. (18) is proved. $\qquad\square$

## D.1 Proof of Lemma 3.1

Proof: Since $\theta = 1/(4\sqrt{n}\delta)$ satisfies the condition required in Lemma D.1, we get

$$\mathbb{E}_t[f(\boldsymbol{x}_{t+1}) - f(\boldsymbol{x})] \stackrel{(18)}{\le} \mathbb{E}_t D_h(\boldsymbol{x}, \boldsymbol{x}_t) - \left(1 + \frac{2\mu\theta}{5}\right) D_h(\boldsymbol{x}, \boldsymbol{x}_{t+1}) + \frac{2}{9n} D_h(\boldsymbol{w}_0, \boldsymbol{x}_t).$$

Taking $t = T - 1$ with $T \sim \text{Geom}(p)$ and noting that $\boldsymbol{w}^+ = \boldsymbol{x}_T, \boldsymbol{w}_0 = \boldsymbol{x}_0$, by Lemma A.4, we get

$$\mathbb{E}[f(\boldsymbol{w}^+) - f(\boldsymbol{x})] = \mathbb{E}[f(\boldsymbol{x}_T) - f(\boldsymbol{x})]$$

$$\le \mathbb{E} D_h(\boldsymbol{x}, \boldsymbol{x}_{T-1}) - D_h(\boldsymbol{x}, \boldsymbol{x}_T) - \frac{2\mu\theta}{5} D_h(\boldsymbol{x}, \boldsymbol{x}_T) + \frac{2}{9n} D_h(\boldsymbol{w}_0, \boldsymbol{x}_{T-1})$$

$$\stackrel{(15)}{=} \mathbb{E}\, p D_h(\boldsymbol{x}, \boldsymbol{x}_0) - p D_h(\boldsymbol{x}, \boldsymbol{x}_T) - \frac{2\mu\theta}{5} D_h(\boldsymbol{x}, \boldsymbol{x}_T)$$

$$+ \frac{2}{9n} [(1 - p) D_h(\boldsymbol{w}_0, \boldsymbol{x}_T) + p D_h(\boldsymbol{w}_0, \boldsymbol{x}_0)]$$

$$\le \mathbb{E}\, p D_h(\boldsymbol{x}, \boldsymbol{w}_0) - p D_h(\boldsymbol{x}, \boldsymbol{w}^+) - \frac{2\mu\theta}{5} D_h(\boldsymbol{x}, \boldsymbol{w}^+) + \frac{2}{9n} D_h(\boldsymbol{w}_0, \boldsymbol{w}^+) \qquad (21)$$

$$\stackrel{(14)}{=} \mathbb{E}\, p \langle \boldsymbol{x} - \boldsymbol{w}_0, \nabla h(\boldsymbol{w}^+) - \nabla h(\boldsymbol{w}_0)\rangle - \frac{9pn - 2}{9n} D_h(\boldsymbol{w}_0, \boldsymbol{w}^+) - \frac{2\mu\theta}{5} D_h(\boldsymbol{x}, \boldsymbol{w}^+).$$

Thus, Eq. (11) is proved. $\qquad\square$

## D.2 Proof of Theorem 3.3

Proof: Choosing $\boldsymbol{x} = \boldsymbol{x}_*$ and $\boldsymbol{x} = \boldsymbol{w}_k$ in Eq. (21), which are all independent to indices $i_1, i_2 \ldots, i_T$ in $\mathrm{SVRS}^{\mathrm{1ep}}(f, \boldsymbol{w}_k, \theta, p)$, then we get

$$\mathbb{E}_k[f(\boldsymbol{w}_{k+1}) - f(\boldsymbol{x}_*)] \le \mathbb{E}_k p D_h(\boldsymbol{x}_*, \boldsymbol{w}_k) - p D_h(\boldsymbol{x}_*, \boldsymbol{w}_{k+1}) - \frac{2\mu\theta}{5} D_h(\boldsymbol{x}_*, \boldsymbol{w}_{k+1}) + \frac{2}{9n} D_h(\boldsymbol{w}_k, \boldsymbol{w}_{k+1}).$$

$$\mathbb{E}_k[f(\boldsymbol{w}_{k+1}) - f(\boldsymbol{w}_k)] \le \mathbb{E}_k p D_h(\boldsymbol{w}_k, \boldsymbol{w}_k) - p D_h(\boldsymbol{w}_k, \boldsymbol{w}_{k+1}) - \frac{2\mu\theta}{5} D_h(\boldsymbol{w}_k, \boldsymbol{w}_{k+1}) + \frac{2}{9n} D_h(\boldsymbol{w}_k, \boldsymbol{w}_{k+1}).$$

Adding both inequalities together, we could obtain

$$\mathbb{E}\left[2f(\boldsymbol{w}_{k+1}) - f(\boldsymbol{w}_k) - f(\boldsymbol{x}_*)\right] \le \mathbb{E} p D_h(\boldsymbol{x}_*, \boldsymbol{w}_k) - \left(p + \frac{2\mu\theta}{5}\right) D_h(\boldsymbol{x}_*, \boldsymbol{w}_{k+1})$$
$$- \left(p + \frac{2\mu\theta}{5} - \frac{4}{9n}\right) D_h(\boldsymbol{w}_k, \boldsymbol{w}_{k+1}).$$

Noting that $p = 1/n$, thus $p + \frac{2\mu\theta}{5} - \frac{4}{9n} > 0$. Based on Eq. (9), after rearranging the terms, we get

$$\mathbb{E}[f(\boldsymbol{w}_{k+1}) - f(\boldsymbol{x}_*)] + \frac{1}{2}\left(p + \frac{2\mu\theta}{5}\right) D_h(\boldsymbol{x}_*, \boldsymbol{w}_{k+1}) \le \mathbb{E}\frac{1}{2}[f(\boldsymbol{w}_k) - f(\boldsymbol{x}_*)] + \frac{p}{2} D_h(\boldsymbol{x}_*, \boldsymbol{w}_k).$$

Now we denote the potential function as

$$\Phi_k = \mathbb{E}[f(\boldsymbol{w}_k) - f(\boldsymbol{x}_*)] + \frac{1}{2}\left(p + \frac{2\mu\theta}{5}\right) D_h(\boldsymbol{x}_*, \boldsymbol{w}_k). \tag{22}$$

By $\theta = 1/(4\sqrt{n}\delta)$, we obtain

$$\mathbb{E}\Phi_{k+1} \le \max\left\{1 - \frac{1}{2}, \left(1 + \frac{2\mu\theta}{5p}\right)^{-1}\right\} \mathbb{E}\Phi_k = \max\left\{1 - \frac{1}{2}, \left(1 + \frac{2\mu\sqrt{n}}{5\delta}\right)^{-1}\right\} \mathbb{E}\Phi_k.$$

When $\frac{2\mu\sqrt{n}}{5\delta} \ge 1$, we get $\mathbb{E}\Phi_{k+1} \le \frac{1}{2}\mathbb{E}\Phi_k$. Otherwise, $\frac{2\mu\sqrt{n}}{5\delta} < 1$, by inequality $\frac{1}{1+x} \le 1 - \frac{x}{2}, \forall 0 \le x \le 1$, we get $\left(1 + \frac{2\mu\sqrt{n}}{5\delta}\right)^{-1} \le 1 - \frac{\mu\sqrt{n}}{5\delta}$. Hence, $\mathbb{E}\Phi_{k+1} \le \left(1 - \frac{\mu\sqrt{n}}{5\delta}\right)\mathbb{E}\Phi_k$. Therefore, we obtain $\mathbb{E}\Phi_{k+1} \le \max\left\{1 - \frac{1}{2}, 1 - \frac{\mu\sqrt{n}}{5\delta}\right\}\mathbb{E}\Phi_k$. Moreover, the initial term

$$\Phi_0 \stackrel{(22)}{=} f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*) + \frac{1}{2}\left(p + \frac{2\mu\theta}{5}\right) D_h(\boldsymbol{x}_*, \boldsymbol{w}_0)$$
$$\stackrel{(9)}{\le} f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*) + \frac{1}{2}\left(p + \frac{2\mu\theta}{5}\right) \frac{1 + \sqrt{n}\theta\delta}{2\theta} \|\boldsymbol{x}_* - \boldsymbol{w}_0\|^2$$
$$\stackrel{(i)}{\le} f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*) + \frac{1}{2}\left(\frac{5\delta}{2\sqrt{n}} + \frac{\mu}{4}\right) \|\boldsymbol{w}_0 - \boldsymbol{x}_*\|^2 \stackrel{(ii)}{\le} \left[1 + \frac{1}{2}\left(\frac{5\delta}{\mu\sqrt{n}} + \frac{1}{2}\right)\right][f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*)]$$
$$\le 3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right)[f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*)],$$

where $(i)$ uses $\theta = \sqrt{p}/(4\delta)$ and $(ii)$ uses $\frac{\mu}{2}\|\boldsymbol{w}_0 - \boldsymbol{x}_*\|^2 \stackrel{(2)}{\le} f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*)$. Then we finally get

$$\mathbb{E} f(\boldsymbol{w}_k) - f(\boldsymbol{x}_*) \stackrel{(22)}{=} \mathbb{E}\Phi_k \le \left(\max\left\{1 - \frac{1}{2}, 1 - \frac{\mu\sqrt{n}}{5\delta}\right\}\right)^k \cdot 3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right)[f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*)].$$

In order to make $\mathbb{E}\Phi_k \le \varepsilon$, we need

$$\exp\left\{-\frac{k}{\max\left\{2, \frac{5\delta}{\mu\sqrt{n}}\right\}}\right\} \cdot 3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right)[f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*)] \le \varepsilon,$$

which leads to $k \ge K_1 := \max\left\{2, \frac{5\delta}{\mu\sqrt{n}}\right\} \log \frac{3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right)[f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*)]}{\varepsilon}$.

Noting that one-epoch communication complexity in $\mathrm{SVRS}^{\mathrm{1ep}}$ is $\Theta(n)$ in expectation when $p = 1/n$ (shown in Section 3.1.1), we get total communication complexity is $\tilde{\mathcal{O}}(n + \sqrt{n}\delta/\mu)$. $\qquad\square$

## D.3 Proof of Lemma 3.5

Proof: Based on Lemma 3.1 and noting that $\boldsymbol{y}_{k+1} = \text{SVRS}^{1\text{ep}}(f, \boldsymbol{x}_{k+1}, \theta, p)$, we get

$$\mathbb{E}_k[f(\boldsymbol{y}_{k+1}) - f(\boldsymbol{x})] \overset{(11)}{\leq} \mathbb{E}_k p \langle \boldsymbol{x} - \boldsymbol{x}_{k+1}, \nabla h(\boldsymbol{y}_{k+1}) - \nabla h(\boldsymbol{x}_{k+1}) \rangle$$
$$- \frac{7p}{9} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1}) - \frac{2\mu\theta}{5} D_h(\boldsymbol{x}, \boldsymbol{y}_{k+1}). \tag{23}$$

Here $\boldsymbol{x} \in \mathbb{R}^d$ should be independent to random indices $i_1^{(k)}, i_2^{(k)}, \dots, i_T^{(k)}$ in $\text{SVRS}^{1\text{ep}}(f, \boldsymbol{x}_{k+1}, \theta, p)$. Then we can apply interpolation $\boldsymbol{z}_k$ to derive

$$\mathbb{E}_k[f(\boldsymbol{y}_{k+1}) - f(\boldsymbol{x})] \overset{(23)}{\leq} \mathbb{E}_k \, p \langle \boldsymbol{z}_k - \boldsymbol{x}_{k+1}, \nabla h(\boldsymbol{y}_{k+1}) - \nabla h(\boldsymbol{x}_{k+1}) \rangle - \frac{7p}{9} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1})$$

$$+ p \langle \boldsymbol{x} - \boldsymbol{z}_k, \nabla h(\boldsymbol{y}_{k+1}) - \nabla h(\boldsymbol{x}_{k+1}) \rangle - \frac{2\mu\theta}{5} D_h(\boldsymbol{x}, \boldsymbol{y}_{k+1})$$

$$= \mathbb{E}_k \frac{1-\tau}{\tau} \cdot p \langle \boldsymbol{x}_{k+1} - \boldsymbol{y}_k, \nabla h(\boldsymbol{y}_{k+1}) - \nabla h(\boldsymbol{x}_{k+1}) \rangle - \frac{7p}{9} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1})$$

$$+ p \langle \boldsymbol{x} - \boldsymbol{z}_k, \nabla h(\boldsymbol{y}_{k+1}) - \nabla h(\boldsymbol{x}_{k+1}) \rangle - \frac{2\mu\theta}{5} D_h(\boldsymbol{x}, \boldsymbol{y}_{k+1})$$

$$\overset{(i)}{\leq} \mathbb{E}_k \frac{1-\tau}{\tau} \left[ f(\boldsymbol{y}_k) - f(\boldsymbol{y}_{k+1}) - \frac{7p}{9} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1}) \right] - \frac{7p}{9} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1})$$

$$+ p \langle \boldsymbol{z}_k - \boldsymbol{x}, \nabla h(\boldsymbol{x}_{k+1}) - \nabla h(\boldsymbol{y}_{k+1}) \rangle - \frac{2\mu\theta}{5} D_h(\boldsymbol{x}, \boldsymbol{y}_{k+1}) \tag{24}$$

where $(i)$ uses Eq. (23) with $\boldsymbol{x} = \boldsymbol{y}_k$, which is independent to indices in $\text{SVRS}^{1\text{ep}}(f, \boldsymbol{x}_{k+1}, \theta, p)$. We continue obtaining

$$\mathbb{E}_k[f(\boldsymbol{y}_{k+1}) - f(\boldsymbol{x})] \overset{(24)(9)}{\leq} \mathbb{E}_k \frac{1-\tau}{\tau} [f(\boldsymbol{y}_k) - f(\boldsymbol{y}_{k+1})] - \frac{7p}{9\tau} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1})$$

$$+ p \langle \boldsymbol{z}_k - \boldsymbol{x}, \nabla h(\boldsymbol{x}_{k+1}) - \nabla h(\boldsymbol{y}_{k+1}) \rangle - \frac{\mu(1 - \sqrt{n}\theta\delta)}{5} \|x - y_{k+1}\|^2$$

$$\overset{(i)}{\leq} \mathbb{E}_k \frac{1-\tau}{\tau} [f(\boldsymbol{y}_k) - f(\boldsymbol{y}_{k+1})] - \frac{7p}{9\tau} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1})$$

$$+ \mathbb{E}_k \left[ \mathbb{E}_{j_k} \langle \boldsymbol{z}_k - \boldsymbol{x}, \boldsymbol{\mathcal{G}}_{k+1} \rangle - \frac{3\mu}{20} \|\boldsymbol{x} - \boldsymbol{y}_{k+1}\|^2 \right]$$

$$\overset{(16)}{\leq} \mathbb{E}_k \frac{1-\tau}{\tau} [f(\boldsymbol{y}_k) - f(\boldsymbol{y}_{k+1})] - \frac{7p}{9\tau} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1})$$

$$+ \mathbb{E}_k \left[ \mathbb{E}_{j_k} \frac{\alpha}{2} \|\boldsymbol{\mathcal{G}}_{k+1}\|^2 + \frac{\|\boldsymbol{x} - \boldsymbol{z}_k\|^2}{2\alpha} - \frac{1 + 0.3\mu\alpha}{2\alpha} \|\boldsymbol{x} - \boldsymbol{z}_{k+1}\|^2 \right],$$

where $(i)$ uses $\mathbb{E}_{j_k} \boldsymbol{\mathcal{G}}_{k+1} = p \left[ \nabla h(\boldsymbol{x}_{k+1}) - \nabla h(\boldsymbol{y}_{k+1}) \right]$ and $\sqrt{n}\theta\delta = 1/4$.

Furthermore, we can estimate

$$\mathbb{E}_{j_k} \|\boldsymbol{\mathcal{G}}_{k+1}\|^2 = p^2 \mathbb{E}_{j_k} \|\nabla h(\boldsymbol{x}_{k+1}) - \nabla h(\boldsymbol{y}_{k+1}) + \nabla[f - f_{j_k}](\boldsymbol{x}_{k+1}) - \nabla[f - f_{j_k}](\boldsymbol{y}_{k+1})\|^2$$

$$\overset{(i)}{=} p^2 \mathbb{E}_{j_k} \|\nabla h(\boldsymbol{x}_{k+1}) - \nabla h(\boldsymbol{y}_{k+1})\|^2 + \|\nabla[f - f_{j_k}](\boldsymbol{x}_{k+1}) - \nabla[f - f_{j_k}](\boldsymbol{y}_{k+1})\|^2$$

$$\overset{(4)}{\leq} p^2 \mathbb{E}_{j_k} \|\nabla h(\boldsymbol{x}_{k+1}) - \nabla h(\boldsymbol{y}_{k+1})\|^2 + p^2 \delta^2 \|\boldsymbol{x}_{k+1} - \boldsymbol{y}_{k+1}\|^2$$

$$\overset{(ii)}{\leq} \frac{2(1 + \sqrt{n}\theta\delta)p^2}{\theta} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1}) + \frac{2\theta p^2 \delta^2}{1 - \sqrt{n}\theta\delta} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1})$$

$$= \frac{5p^2}{2\theta} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1}) + \frac{p^2}{6n\theta} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1}) \leq \frac{8p^2}{3\theta} D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1})$$

where $(i)$ uses $\mathbb{E}_{j_k} \nabla[f - f_{j_k}](\boldsymbol{x}_{k+1}) - \nabla[f - f_{j_k}](\boldsymbol{y}_{k+1}) = \boldsymbol{0}$, $(ii)$ uses the convexity and smoothness of $h$ (e.g., [22, Theorem A.1 (iii)]) and Eq. (9). After rearrangement, we get

$$\mathbb{E}_k \frac{\alpha}{\tau}[f(\boldsymbol{y}_{k+1}) - f(\boldsymbol{x})] \leq \mathbb{E}_k (1-\tau) \cdot \frac{\alpha}{\tau}[f(\boldsymbol{y}_k) - f(\boldsymbol{x})] + \frac{\|\boldsymbol{x} - \boldsymbol{z}_k\|^2}{2} - \frac{1 + 0.3\mu\alpha}{2}\|\boldsymbol{x} - \boldsymbol{z}_{k+1}\|^2$$
$$+ \alpha \left(\frac{4\alpha p^2}{3\theta} - \frac{7p}{9\tau}\right) D_h(\boldsymbol{x}_{k+1}, \boldsymbol{y}_{k+1}).$$

Hence, we see that once $2\tau\alpha p \leq \theta$, Eq. (12) holds. $\qquad\square$

### D.4   Proof of Theorem 3.6

Proof: Taking $\boldsymbol{x} = \boldsymbol{x}_*$ in Eq. (12), which is independent of any index during the process, we get

$$\mathbb{E} \frac{\alpha}{\tau}[f(\boldsymbol{y}_{k+1}) - f(\boldsymbol{x}_*)] + \frac{(1 + 0.3\mu\alpha)\|\boldsymbol{x}_* - \boldsymbol{z}_{k+1}\|^2}{2} \leq \mathbb{E}(1-\tau) \cdot \frac{\alpha}{\tau}[f(\boldsymbol{y}_k) - f(\boldsymbol{x}_*)] + \frac{\|\boldsymbol{x}_* - \boldsymbol{z}_k\|^2}{2}.$$

Denote the potential function as

$$\Phi_k = [f(\boldsymbol{y}_k) - f(\boldsymbol{x}_*)] + \frac{\tau(1 + 0.3\mu\alpha)}{2\alpha}\|\boldsymbol{x}_* - \boldsymbol{z}_k\|^2.$$

We obtain

$$\mathbb{E}\Phi_{k+1} \leq \max\left\{1 - \tau, \left(1 + \frac{\mu\sqrt{n}}{\delta} \cdot \frac{3}{80\tau}\right)^{-1}\right\}\mathbb{E}\Phi_k.$$

When $\tau = \frac{1}{4} \leq \frac{1}{8}n^{1/4}\sqrt{\frac{\mu}{\delta}}$, then we have that

$$(1-\tau)\left(1 + \frac{\mu\sqrt{n}}{\delta} \cdot \frac{3}{80\tau}\right) \geq (1-\tau)\left(1 + \frac{3}{20\tau}\right) \geq 1 \Rightarrow \mathbb{E}\Phi_{k+1} \leq \left(1 - \frac{1}{4}\right)\mathbb{E}\Phi_k.$$

When $\tau = \frac{n^{1/4}}{8}\sqrt{\frac{\mu}{\delta}} \leq \frac{1}{4}$, we get

$$t := \frac{\mu\sqrt{n}}{\delta} \cdot \frac{3}{80\tau} = \frac{3n^{1/4}}{10}\sqrt{\frac{\mu}{\delta}} \leq \frac{3}{5} \Rightarrow \frac{1}{1+t} \leq 1 - \frac{5t}{8} \Rightarrow \mathbb{E}\Phi_{k+1} \leq \left(1 - \frac{n^{1/4}}{8}\sqrt{\frac{\mu}{\delta}}\right)\mathbb{E}\Phi_k.$$

Therefore, we finally obtain

$$\mathbb{E}\Phi_{k+1} \leq \max\left\{1 - \frac{1}{4}, 1 - \frac{n^{1/4}}{8}\sqrt{\frac{\mu}{\delta}}\right\}\mathbb{E}\Phi_k.$$

By the strong convexity of $f$ in Assumption 1 and the choice of $\tau$ and $\alpha$, the initial term

$$\begin{aligned}
\Phi_0 &= [f(\boldsymbol{y}_0) - f(\boldsymbol{x}_*)] + \frac{\tau(1 + 0.3\mu\alpha)}{2\alpha}\|\boldsymbol{x}_* - \boldsymbol{y}_0\|^2 \\
&= [f(\boldsymbol{y}_0) - f(\boldsymbol{x}_*)] + \left(\frac{8\delta\tau^2}{\sqrt{n}\mu} + 0.3\tau\right)\frac{\mu}{2}\|\boldsymbol{x}_* - \boldsymbol{y}_0\|^2 \\
&\leq \left(1 + \frac{1}{8} + \frac{0.3}{4}\right)[f(\boldsymbol{y}_0) - f(\boldsymbol{x}_*)] \leq 2[f(\boldsymbol{y}_0) - f(\boldsymbol{x}_*)].
\end{aligned}$$

To obtain $\varepsilon$-error solution, we need

$$k \geq K_2 = \max\left\{4, 8n^{-1/4}\sqrt{\frac{\delta}{\mu}}\right\}\log\frac{2[f(\boldsymbol{y}_0) - f(\boldsymbol{x}_*)]}{\varepsilon}.$$

Note that every call of Algorithm SVRS$^{\text{1ep}}$ requires $4n$ communication in expectation (shown in Section 3.1.1). The remaining communication in one iteration of AccSVRS need 4 communication (the master sends $\boldsymbol{x}_{k+1}$ and $\boldsymbol{y}_{k+1}$ to the client $j_k$, and then receives $\nabla f_{j_k}(\boldsymbol{x}_{k+1})$ and $\nabla f_{j_k}(\boldsymbol{y}_{k+1})$). Thus one iteration of AccSVRS is $\Theta(n)$ in expectation, leading to the total communication complexity for $\varepsilon$-error solution is $\tilde{\mathcal{O}}\left(n + n^{3/4}\sqrt{\frac{\delta}{\mu}}\right)$. $\qquad\square$

**Algorithm 4** Loopless Stochastic Variance-Reduced Sliding (SVRS)

---

1: **Input:** $\boldsymbol{w}_0 \in \mathbb{R}^d, p \in (0,1), \theta > 0, K \in \{1, 2, \dots\}$
2: Initialize $\boldsymbol{x}_0 = \boldsymbol{w}_0$ and compute $\nabla f(\boldsymbol{w}_0)$
3: **for** $k = 0, 1, 2, \dots, K-1$ **do**
4:     Sample $i_k \sim \text{Unif}([n])$ and compute $\boldsymbol{g}_k = \nabla f_{i_k}(\boldsymbol{w}_k) - \nabla f(\boldsymbol{w}_k)$
5:     Approximately solve the local proximal point problem:

$$\boldsymbol{x}_{k+1} \approx \underset{\boldsymbol{x} \in \mathbb{R}^d}{\arg\min}\, A_\theta^k(\boldsymbol{x}) := \langle \nabla f_{i_k}(\boldsymbol{x}_k) - \boldsymbol{g}_k - \nabla f_1(\boldsymbol{x}_k), \boldsymbol{x} - \boldsymbol{x}_k \rangle + \frac{1}{2\theta} \|\boldsymbol{x} - \boldsymbol{x}_k\|^2 + f_1(\boldsymbol{x})$$

6:     $\boldsymbol{w}_{k+1} = \begin{cases} \boldsymbol{x}_{k+1} & \text{with probability} \quad p \\ \boldsymbol{w}_k & \text{with probability} \quad 1-p \end{cases}$
7: **end for**
8: **Output:** $\boldsymbol{w}_K$

---

### D.5 Loopless SVRS

In this section, we describe the loopless SVRS (Algorithm 4). By simple facts shown in Proposition A.2, $\text{SVRS}^{\text{1ep}}(f, \boldsymbol{w}_k, \theta, p)$ can be viewed as the inter iteration until $\boldsymbol{w}_k$ in loopless SVRS is updated. Thus, the one-step variation in Lemma D.1 still holds. Hence, we can derive a similar convergence rate and communication complexity for loopless SVRS.

**Theorem D.2** *Suppose Assumption 1 holds. If in loopless SVRS (Algorithm 4), the hyperparameters are set as $\theta = 1/(4\sqrt{n}\delta), p = 1/n$, and the approximate solution in each proximal step satisfies Eq. (10), then for any error $\varepsilon > 0$, when*

$$k \geq K_1 := \max\left\{2n, \frac{11\sqrt{n}\delta}{\mu}\right\} \log \frac{3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right) [f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*)]}{\varepsilon},$$

*i.e., after $\tilde{\mathcal{O}}(n + \sqrt{n}\delta/\mu)$ communications, we can guarantee that $\mathbb{E}f(\boldsymbol{w}_k) - f(\boldsymbol{x}_*) \leq \varepsilon$.*

Proof: Noting that in each step of loopless SVRS, the anchor point is $\boldsymbol{w}_k$ instead of $\boldsymbol{w}_0$, thus Eq. (18) holds after replacing $\boldsymbol{w}_0$ to $\boldsymbol{w}_k$. Now choosing $\boldsymbol{x} = \boldsymbol{x}_*$ and $\boldsymbol{x} = \boldsymbol{w}_k$ in Eq. (18), which are all independent to index $i_k$, we get

$$\mathbb{E}_k[f(\boldsymbol{x}_{k+1}) - f(\boldsymbol{x}_*)] \leq \mathbb{E}_k D_h(\boldsymbol{x}_*, \boldsymbol{x}_k) - \left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{x}_*, \boldsymbol{x}_{k+1}) + \frac{2\theta^2\delta^2}{(1 - \sqrt{n}\theta\delta)^2} D_h(\boldsymbol{w}_k, \boldsymbol{x}_k).$$

$$\mathbb{E}_k[f(\boldsymbol{x}_{k+1}) - f(\boldsymbol{w}_k)] \leq \mathbb{E}_k D_h(\boldsymbol{w}_k, \boldsymbol{x}_k) - \left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{w}_k, \boldsymbol{x}_{k+1}) + \frac{2\theta^2\delta^2}{(1 - \sqrt{n}\theta\delta)^2} D_h(\boldsymbol{w}_k, \boldsymbol{x}_k).$$

Adding both inequalities together and noting that

$$\mathbb{E}_k D_h(\boldsymbol{w}_{k+1}, \boldsymbol{x}_{k+1}) = \mathbb{E}_k(1-p)D_h(\boldsymbol{w}_k, \boldsymbol{x}_{k+1}) + pD_h(\boldsymbol{x}_{k+1}, \boldsymbol{x}_{k+1}) = (1-p)\mathbb{E}_k D_h(\boldsymbol{w}_k, \boldsymbol{x}_{k+1}),$$

as well as

$$\mathbb{E}_k f(\boldsymbol{w}_{k+1}) = \mathbb{E}_k(1-p)f(\boldsymbol{w}_k) + pf(\boldsymbol{x}_{k+1}),$$

we could obtain

$$\mathbb{E}\frac{2}{p}[f(\boldsymbol{w}_{k+1}) - (1-p)f(\boldsymbol{w}_k)] - f(\boldsymbol{w}_k) - f(\boldsymbol{x}_*)$$

$$\leq \mathbb{E}D_h(\boldsymbol{x}_*, \boldsymbol{x}_k) - \left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{x}_*, \boldsymbol{x}_{k+1}) + \left(1 + \frac{4\theta^2\delta^2}{(1 - \sqrt{n}\theta\delta)^2}\right) D_h(\boldsymbol{w}_k, \boldsymbol{x}_k) - D_h(\boldsymbol{w}_k, \boldsymbol{x}_{k+1})$$

$$= \mathbb{E}D_h(\boldsymbol{x}_*, \boldsymbol{x}_k) - \left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{x}_*, \boldsymbol{x}_{k+1}) + \left(1 + \frac{4\theta^2\delta^2}{(1 - \sqrt{n}\theta\delta)^2}\right) D_h(\boldsymbol{w}_k, \boldsymbol{x}_k) - \frac{D_h(\boldsymbol{w}_{k+1}, \boldsymbol{x}_{k+1})}{1-p}.$$

Rearranging the terms, we get

$$\mathbb{E}[f(\boldsymbol{w}_{k+1}) - f(\boldsymbol{x}_*)] + \frac{p}{2}\left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{x}_*, \boldsymbol{x}_{k+1}) + \frac{p}{2(1-p)} D_h(\boldsymbol{w}_{k+1}, \boldsymbol{x}_{k+1})$$

$$\leq \quad \mathbb{E}(1 - \frac{p}{2})[f(\boldsymbol{w}_k) - f(\boldsymbol{x}_*)] + \frac{p}{2}D_h(\boldsymbol{x}_*, \boldsymbol{x}_k) + \frac{p}{2}\left(1 + \frac{4\theta^2\delta^2}{(1 - \sqrt{n}\theta\delta)^2}\right) D_h(\boldsymbol{w}_k, \boldsymbol{x}_k).$$

Now we denote the potential function as

$$\Phi_k = \mathbb{E}[f(\boldsymbol{w}_k) - f(\boldsymbol{x}_*)] + \frac{p}{2}\left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{x}_*, \boldsymbol{x}_k) + \frac{p}{2(1-p)} D_h(\boldsymbol{w}_k, \boldsymbol{x}_k).$$

Then we obtain

$$\mathbb{E}\Phi_{k+1} \le \max\left\{1 - \frac{p}{2}, \left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right)^{-1}, \left(1 + \frac{4\theta^2\delta^2}{(1 - \sqrt{n}\theta\delta)^2}\right)(1-p)\right\}\mathbb{E}\Phi_k.$$

Since we choose $\theta = 1/(4\sqrt{n}\delta)$, we get $\theta\mu \le \sqrt{n}\theta\delta \le 1/4$ by Assumption 1, which shows that

$$\left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right)^{-1} = 1 - \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta + \mu\theta/2} = 1 - \frac{4\mu\theta}{11} = 1 - \frac{\mu}{11\delta\sqrt{n}}.$$

Additionally, by $p = 1/n$ and $\theta = 1/(4\delta\sqrt{n})$, we also have that

$$\left(1 + \frac{4\theta^2\delta^2}{(1 - \sqrt{n}\theta\delta)^2}\right)(1-p) = \left(1 + \left(\frac{\frac{1}{2\sqrt{n}}}{1 - \frac{1}{4}}\right)^2\right)(1-p) = \left(1 + \frac{4p}{9}\right)(1-p) \le 1 - \frac{5p}{9}.$$

Therefore, we obtain the ratio between $\mathbb{E}\Phi_{k+1}$ and $\mathbb{E}\Phi_k$:

$$\mathbb{E}\Phi_{k+1} \le \max\left\{1 - \frac{p}{2}, 1 - \frac{\mu}{11\delta\sqrt{n}}, 1 - \frac{5p}{9}\right\}\mathbb{E}\Phi_k \le \max\left\{1 - \frac{p}{2}, 1 - \frac{\mu}{11\delta\sqrt{n}}\right\}\mathbb{E}\Phi_k.$$

Moreover, the initial term

$$\Phi_0 = f(\boldsymbol{w}_0) - f(\boldsymbol{x}_*) + \frac{p}{2}\left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right) D_h(\boldsymbol{x}_*, \boldsymbol{x}_0)$$

$$\overset{(9)}{\le} f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*) + \frac{p}{2}\left(1 + \frac{\mu\theta/2}{1 + \sqrt{n}\theta\delta}\right)\frac{1 + \sqrt{n}\theta\delta}{2\theta}\|\boldsymbol{x}_* - \boldsymbol{x}_0\|^2$$

$$\overset{(i)}{\le} f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*) + \frac{p}{2}\left(\frac{2.5\delta}{\sqrt{p}} + \frac{\mu}{4}\right)\|\boldsymbol{x}_0 - \boldsymbol{x}_*\|^2 \overset{(ii)}{\le} \left[1 + \frac{p}{2}\left(\frac{5\delta}{\mu\sqrt{p}} + \frac{1}{2}\right)\right][f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*)]$$

$$\le 3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right)[f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*)],$$

where $(i)$ uses $\theta = \sqrt{p}/(4\delta)$ and $(ii)$ uses $f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*) \overset{(2)}{\ge} \frac{\mu}{2}\|\boldsymbol{x}_0 - \boldsymbol{x}_*\|^2$. Then we finally get

$$\mathbb{E}f(\boldsymbol{w}_k) - f(\boldsymbol{x}_*) \le \mathbb{E}\Phi_k \le \left(\max\left\{1 - \frac{1}{2n}, 1 - \frac{\mu}{11\delta\sqrt{n}}\right\}\right)^k \cdot 3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right)[f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*)].$$

In order to make $\mathbb{E}\Phi_k \le \varepsilon$, we need

$$\exp\left\{-\frac{k}{\max\left\{2n, \frac{11\sqrt{n}\delta}{\mu}\right\}}\right\} \cdot 3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right)[f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*)] \le \varepsilon,$$

which leads to

$$k \ge K_1 := \max\left\{2n, \frac{11\sqrt{n}\delta}{\mu}\right\}\log\frac{3\left(1 + \frac{\delta}{\mu\sqrt{n}}\right)[f(\boldsymbol{x}_0) - f(\boldsymbol{x}_*)]}{\varepsilon},$$

Noting that communication complexity in each iteration is $2p(n-1) + 2$ in expectation (by similar analysis in Section 3.1.1), we get communication complexity is 4 in each iteration in expectation. Therefore, the total communication complexity is $\tilde{\mathcal{O}}(n + \frac{\sqrt{n}\delta}{\mu})$ in expectation. $\qquad\square$

23

# E    Computation of Gradient Complexity

We show the detail omitted in Section 3.3. Let $\boldsymbol{x}_{t,*} = \arg\min_{\boldsymbol{x}\in\mathbb{R}^d} A_\theta^t(\boldsymbol{x})$. Noting that by [47, Theorem 2.2.2], we could obtain

$$\left\| \nabla A_\theta^t(\boldsymbol{x}_{t,s}) \right\|^2 \leq 2L' \left( A_\theta^t(\boldsymbol{x}_{t,s}) - A_\theta^t(\boldsymbol{x}_{t,*}) \right) \leq \frac{20\mu'L' \left\| \boldsymbol{x}_t - \boldsymbol{x}_{t,*} \right\|^2}{3} \left[ e^{(s+1)/\sqrt{\kappa'}} - 1 \right]^{-1}$$

if we start from $\boldsymbol{x}_t$ in the proximal step for optimizing $A_\theta^t(\boldsymbol{x})$, where $\kappa' = L'/\mu', L' = L+1/\theta, \mu' = -\sqrt{n}\delta + 1/\theta$ based on assumptions. Then Eq. (10) could be satisfied after $T_{\mathrm{app}}$ iterations when

$$\frac{20\mu'L' \left\| \boldsymbol{x}_t - \boldsymbol{x}_{t,*} \right\|^2}{3} \left[ e^{(T_{\mathrm{app}}+1)/\sqrt{\kappa'}} - 1 \right]^{-1} \leq \frac{\mu}{20\theta} \left\| \boldsymbol{x}_t - \boldsymbol{x}_{t,*} \right\|^2 .$$

Note that $\theta = 1/(4\sqrt{n}\delta)$, which leads to

$$T_{\mathrm{app}} = \mathcal{O}\left( \sqrt{\frac{1+\theta L}{1-\sqrt{n}\theta\delta}} \cdot \log\left( \frac{(1+\theta L)(1-\sqrt{n}\theta\delta)}{\mu\theta} \right) \right) = \mathcal{O}\left( \left( 1 + n^{-1/4}\sqrt{\delta/\mu} \right) \log\frac{\sqrt{n}\delta + L}{\mu} \right).$$

Hence, the total number of gradient calls in expectation is

$$\mathcal{O}(nT_{\mathrm{app}} \cdot K_2) = \tilde{\mathcal{O}}\left[ \left( n + n^{3/4}\sqrt{\frac{\delta}{\mu}} \right)\left( 1 + \frac{1}{n^{1/4}}\sqrt{\frac{L}{\delta}} \right) \right] = \tilde{\mathcal{O}}\left( n + n^{3/4}\left( \sqrt{\frac{\delta}{\mu}} + \sqrt{\frac{L}{\delta}} \right) + \sqrt{\frac{nL}{\mu}} \right).$$

Since $\delta \in [\mu, L]$, we obtain $\sqrt{\frac{\delta}{\mu}} + \sqrt{\frac{L}{\delta}} \leq \sqrt{\frac{L}{\mu}} + 1$, leading to

$$n + n^{3/4}\left( \sqrt{\frac{\delta}{\mu}} + \sqrt{\frac{L}{\delta}} \right) + \sqrt{\frac{nL}{\mu}} \leq 2\left( n + n^{3/4}\sqrt{\frac{L}{\mu}} \right).$$

Thus, the gradient complexity is $\tilde{\mathcal{O}}\left( n + n^{3/4}\sqrt{L/\mu} \right)$. Moreover, when $\delta = \Theta(\sqrt{\mu L})$, we obtain

$$n + n^{3/4}\left( \sqrt{\frac{\delta}{\mu}} + \sqrt{\frac{L}{\delta}} \right) = n + \Theta\left( n^{3/4}\left( \frac{L}{\mu} \right)^{1/4} \right) + \sqrt{\frac{nL}{\mu}} = \Theta\left( n + \sqrt{\frac{nL}{\mu}} \right).$$

Thus, the gradient complexity is $\tilde{\mathcal{O}}\left( n + \sqrt{nL/\mu} \right)$ in this time.

Note that Assumption 1 and smoothness of $f_1$ could only guarantee

$$\frac{1}{n}\sum_{i=1}^n \left\| \nabla f_i(\boldsymbol{x}) - \nabla f_i(\boldsymbol{y}) \right\|^2 \overset{(4)}{\leq} \delta^2 + \left\| \nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y}) \right\|^2$$

$$\leq \delta^2 + 2\left\| \nabla[f-f_1](\boldsymbol{x}) - \nabla[f-f_1](\boldsymbol{y}) \right\|^2 + 2\left\| \nabla f_1(\boldsymbol{x}) - \nabla f_1(\boldsymbol{y}) \right\|^2 \leq \left[ (2n+1)\delta^2 + 2L^2 \right] \left\| \boldsymbol{x} - \boldsymbol{y} \right\|^2 ,$$

that is, $f_i$'s are $(2L + 2\sqrt{n}\delta)$-average smooth. Hence, the tightness of our gradient complexity holds for the average smooth setting only when $\sqrt{n}\delta = \mathcal{O}(L)$. Moreover, we can also compute

$$\left\| \nabla f_i(\boldsymbol{x}) - \nabla f_i(\boldsymbol{y}) \right\|^2 \leq 2\left\| \nabla[f-f_i](\boldsymbol{x}) - \nabla[f-f_i](\boldsymbol{y}) \right\|^2 + 2\left\| \nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{y}) \right\|^2$$

$$\overset{(4)}{\leq} 2n\delta^2 + 4\left\| \nabla[f-f_1](\boldsymbol{x}) - \nabla[f-f_1](\boldsymbol{y}) \right\|^2 + 4\left\| \nabla f_1(\boldsymbol{x}) - \nabla f_1(\boldsymbol{y}) \right\|^2 \leq \left[ 6n\delta^2 + 4L^2 \right] \left\| \boldsymbol{x} - \boldsymbol{y} \right\|^2 ,$$

that is, $f_i$'s are $(2L + 3\sqrt{n}\delta)$-smooth. Hence, the tightness of our gradient complexity holds for the component smooth setting only when $\delta = \Theta(\sqrt{\mu L})$ and $n\mu = \mathcal{O}(L)$.

# F    Omitted Details of Section 4

In this section, we give the omitted details of Section 4 as well as their proofs.

## F.1    Formal Statement of Definition 4.1 and Discussion

In this subsection, we give the formal statement of Definition 4.1 and show that Algorithm 2 satisfies our definition.

We first introduce the two oracles: the incremental first-order oracle (IFO) [2, 63] and the Proximal Incremental First-order Oracle (PIFO)[11] [56, 25], which are defined as $h_{f_i}^{\mathrm{I}}(\boldsymbol{x}) = [f_i(\boldsymbol{x}), \nabla f_i(\boldsymbol{x})]$

---

[11]Although we have defined PIFO in Section 4.1, we restate it here for completement.

and $h_{f_i}^{\mathrm{P}}(\boldsymbol{x}, \gamma) = [f_i(\boldsymbol{x}), \nabla f_i(\boldsymbol{x}), \mathrm{prox}_{f_i}^{\gamma}(\boldsymbol{x})]$ with $\gamma > 0$ respectively. Here the proximal operator is

$$\mathrm{prox}_{f_i}^{\gamma}(\boldsymbol{x}) := \arg\min_{\boldsymbol{u}} \left\{ f_i(\boldsymbol{u}) + \frac{1}{2\gamma} \|\boldsymbol{x} - \boldsymbol{u}\|^2 \right\} = \arg\min_{\boldsymbol{u}} \left\{ \gamma f_i(\boldsymbol{u}) + \frac{1}{2} \|\boldsymbol{x} - \boldsymbol{u}\|^2 \right\}.$$

The IFO $h_{f_i}^{\mathrm{I}}(\boldsymbol{x})$ takes a point $\boldsymbol{x}$ and a component $f_i$ as input and returns the zero-order and first-order information of the component at $\boldsymbol{x}$. The PIFO $h_{f_i}^{\mathrm{P}}(\boldsymbol{x}, \gamma)$ has an additional input $\gamma > 0$, which can be viewed as the step size of the proximal operator. Besides the local zero-order and first-order information returned by $h_{f_i}^{\mathrm{I}}(\boldsymbol{x})$, $h_{f_i}^{\mathrm{P}}(\boldsymbol{x}, \gamma)$ also provides some global information of $f_i$ by means of the proximal operator. To see this, if we let $\gamma \to +\infty$, $\mathrm{prox}_{f_i}^{\gamma}(\boldsymbol{x})$ converges to the exact minimizer of $f_i$, irrelevant to the choice of $\boldsymbol{x}$. In practice, it could be hard to compute $\mathrm{prox}_{f_i}^{\gamma}(\boldsymbol{x})$ precisely. Nevertheless, since we only focus on communication complexity, it makes no difference to distinguish between the IFO and the PIFO[12]. Thus we assume the algorithm has access to the PIFO and the definition is as follows.

**Definition F.1 (Formal version of Definition 4.1)** *Consider a randomized algorithm $\mathcal{A}$ to solve problem* (1). *Suppose the number of communication rounds is $T$. Define information sets $\mathcal{I}_{t+1}$, $\mathcal{I}_{t+1}^0$ and $\mathcal{I}_{t+1}^1$. Here $\mathcal{I}_{t+1}$ denotes all the information $\mathcal{A}$ obtains after round $t$, while $\mathcal{I}_{t+1}^0$ and $\mathcal{I}_{t+1}^1$ denote the information before and after (possible) anchor point updating during round $t$, respectively. The algorithm updates the information set by the following procedure.*

1. *Choose a distribution $\mathcal{D}$ over $[n]$ with $q_i = \mathbb{P}_{Z \sim \mathcal{D}}(Z = i) > 0$, a positive number $p \leq c_0/n$[13] and the initial points $\boldsymbol{x}_0$. Specify a master note 1 and assume $\max_{2 \leq i \leq n} q_i \leq q_0/n$. Node 1 sends $\boldsymbol{x}_0$ to all the other nodes and other nodes send $h_{f_i}^{\mathrm{P}}(\boldsymbol{x}_0, \gamma_0)$ back to node 1. Initialize the information set as $\mathcal{I}_0 := \mathrm{span}\{\boldsymbol{x}_0, \nabla f_i(\boldsymbol{x}_0), \mathrm{prox}_{f_i}^{\gamma_0}(\boldsymbol{x}_0) \mid 1 \leq i \leq n\}$ and set $t = 0$ and $\tilde{\boldsymbol{x}}_0 = \boldsymbol{x}_0$.*

2. *Sample $i_t \sim \mathcal{D}$. Node 1 sends $\tilde{\boldsymbol{x}}_t$ to node $i_t$ and node $i_t$ sends $h_{f_{i_t}}^{\mathrm{P}}(\tilde{\boldsymbol{x}}_t, \gamma_t)$ back to node 1. Update the information set*

$$\mathcal{I}_{t+1}^0 := \mathrm{span}\{\boldsymbol{y}, \nabla f_{i_t}(\tilde{\boldsymbol{x}}_t), \mathrm{prox}_{f_{i_t}}^{\gamma_t}(\tilde{\boldsymbol{x}}_t) \mid \boldsymbol{y} \in \mathcal{I}_t\} \tag{25}$$

3. *Update the information set $\mathcal{I}_{t+1}^1$ and choose $\boldsymbol{x}_{t+1} \in \mathcal{I}_{t+1}^1$ following the linear-span protocol*

$$\boldsymbol{x}_{t+1} \in \mathcal{I}_{t+1}^1 := \mathrm{span}\{\boldsymbol{y}, \nabla f_1(\boldsymbol{z}), \mathrm{prox}_{f_1}^{\gamma_t'}(\boldsymbol{w}) \mid \boldsymbol{y}, \boldsymbol{z}, \boldsymbol{w} \in \mathcal{I}_{t+1}^0\}. \tag{26}$$

4. *Sample a Bernoulli random variable $a_t$ with expectation equal to $p$. If $a_t = 1$, go to step 5 (update the anchor point); otherwise, set $\tilde{\boldsymbol{x}}_{t+1} = \boldsymbol{x}_{t+1}$, $\mathcal{I}_{t+1} = \mathcal{I}_{t+1}^1$ and go to step 6 (do not update the anchor point).*

5. *Sample $j_t \sim \mathcal{D}$. Node 1 sends some $\boldsymbol{y}_{t+1} \in \mathcal{I}_{t+1}^1$ to node $j_t$ and node $j_t$ sends $h_{f_{j_t}}^{\mathrm{P}}(\boldsymbol{y}_{t+1}, \gamma_{t+1}'')$ back to node 1. Obtain the anchor point $\tilde{\boldsymbol{y}}_{t+1}$ by*

$$\tilde{\boldsymbol{y}}_{t+1} \in \mathrm{span}\{\boldsymbol{y}, \nabla f_1(\boldsymbol{z}), \mathrm{prox}_{f_1}^{\gamma_t'}(\boldsymbol{w}), \nabla f_{j_t}(\boldsymbol{y}_{t+1}), \mathrm{prox}_{f_{j_t}}^{\gamma_{t+1}''}(\boldsymbol{y}_{t+1}) \mid \boldsymbol{y}, \boldsymbol{z}, \boldsymbol{w} \in \mathcal{I}_{t+1}^1\}. \tag{27}$$

*Then node 1 sends the anchor point $\tilde{\boldsymbol{y}}_{t+1}$ to all the other nodes and other nodes send $h_{f_i}^{\mathrm{P}}(\tilde{\boldsymbol{y}}_{t+1}, \gamma_{t+1})$ back to node 1. Update the information set and obtain $\tilde{\boldsymbol{x}}_{t+1}$ by*

$$\tilde{\boldsymbol{x}}_{t+1} \in \mathcal{I}_{t+1} := \mathrm{span}\{\boldsymbol{y}, \tilde{\boldsymbol{y}}_{t+1}, \nabla f_i(\tilde{\boldsymbol{y}}_{t+1}), \mathrm{prox}_{f_i}^{\gamma_{t+1}}(\tilde{\boldsymbol{y}}_{t+1}) \mid \boldsymbol{y} \in \mathcal{I}_{t+1}^1, 1 \leq i \leq n\}. \tag{28}$$

6. *If $t = T - 1$, output some point in $\mathcal{I}_T$; otherwise, set $t \leftarrow t + 1$ and go back to step 2.*

*Here all the random variables $i_t$, $j_t$ and $a_t$ with $0 \leq t \leq T - 1$ are mutually independent, and the step sizes of the proximal operator $\gamma_t$, $\gamma_t'$ and $\gamma_t''$ are positive numbers.*

---

[12]See Lemma F.2

[13]To include catalyst accelerated algorithms, we also need $p \geq c_1/n$ for some $c_1 > 0$ (see footnote 17). To analyze Algorithm 2, $p \leq c_0/n$ is enough.

Now we explain this definition and show that Algorithm 2 (with Algorithm 1 as a part) satisfies our definition.

**Initialization.** In our definition, step 1 is the initialization step. Without loss of generality, we can assume $\boldsymbol{x}_0 = \boldsymbol{0}$ and node 1 is the master node. Otherwise, it suffices to consider $\{\tilde{f}_i(\boldsymbol{x}) = f_i(\boldsymbol{x} + \boldsymbol{x}_0)\}_{i=1}^n$ and exchange the indices between node 1 and the master node. In Algorithm 2, the distribution $\mathcal{D}$ is $\mathrm{Unif}([n])$[14], and $p = 1/n$. In the initialization stage, the algorithm needs to calculate the full gradient of the initial point $\boldsymbol{x}_0$, whose communication cost is $2(n-1)$.

We note that Definition F.1 enjoys a loopless structure while Algorithm 2 has two loops. In fact, when $p$ is fixed, a loopless algorithm is equivalent to a two-loop one with the inner loop size obeying $\mathrm{Geom}(p)$[15].

**Analysis of one communication round.** In each communication round, whether to calculate the full gradient depends on a coin toss with success probability $p$, as shown in step 4.

**The case $a_t = 0$.** We first focus on the case where the full gradient need not be calculated. Such a scenario corresponds to an iteration of Algorithm 1. Each communication round start with step 2. In this step, the algorithm samples a local node, with which the master node communicates. And the communication cost is 2. $\tilde{\boldsymbol{x}}_t$ in this step corresponds to $\boldsymbol{x}_t$ in Algorithm 1. In step 3, the master node calculates the next point based on the current information set $\mathcal{I}_{t+1}^0$ as well as the PIFO $h_{f_1}^{\mathrm{P}}$. This corresponds to line 7 in Algorithm 1. Indeed, the subproblem (7) can be rewritten as finding

$$\arg\min_{\boldsymbol{x} \in \mathbb{R}^d} A_\theta^t(\boldsymbol{x}) = \arg\min_{\boldsymbol{x} \in \mathbb{R}^d} \left\{ \frac{1}{2\theta} \|\boldsymbol{x} - \boldsymbol{x}_t + \theta[\nabla f_{i_t}(\boldsymbol{x}_t) - \boldsymbol{g}_t - \nabla f_1(\boldsymbol{x}_t)]\|^2 + f_1(\boldsymbol{x}) \right\}$$
$$= \mathrm{prox}_{f_1}^\theta \left( \boldsymbol{x}_t - \theta[\nabla f_{i_t}(\boldsymbol{x}_t) - \boldsymbol{g}_t - \nabla f_1(\boldsymbol{x}_t)] \right).$$

If the algorithm has access to the PIFO $h_{f_1}^{\mathrm{P}}$, then the subproblem (7) can be exactly solved by one step of (26). Otherwise, one can apply (26) recursively without the proximal information (i.e., only using the IFO $h_{f_1}^{\mathrm{I}}$), e.g., (accelerated) gradient methods, to find an approximate solution of (7)[16]

**The case $a_t = 1$.** When $a_t = 1$ in step 4, the algorithm needs to perform step 5, which corresponds to an outer iteration of Algorithm 2. Before calculating the full gradient, the algorithm first samples a local node $j_t$ again and the master node communicates the information about $\boldsymbol{y}_{t+1}$ with this node. Here $\boldsymbol{y}_{t+1}$ corresponds to $\boldsymbol{y}_{k+1}$ in Algorithm 2, and the communication cost is 2. Then the master node calculates $\tilde{\boldsymbol{y}}_{t+1}$ by (27), which corresponds to $\boldsymbol{x}_{k+1}$ (of the next iteration) in Algorithm 2. That is to say, lines 7, 8 and 4 (of the next iteration) in Algorithm 2 can be summarized as (27). Then the master node communicates with all the other nodes the information about $\tilde{\boldsymbol{x}}_{t+1}$, and the communication cost is $2(n-1)$. In (28), the algorithm picks up $\tilde{\boldsymbol{x}}_{t+1}$ as the starting point of the next round. In Algorithm 2, $\tilde{\boldsymbol{x}}_{t+1}$ is the same to $\tilde{\boldsymbol{y}}_{t+1}$.

**Communication cost.** From the above analysis, the communication cost in step 5 is $2(n-1)$. Since we assume $q \le c_0/n$, step 5 is performed infrequently and the expected communication cost is (at most) $2(n-1) \cdot p \approx 2c_0$ for a sufficiently large $n$. As a result, the total communication cost of a round is roughly $2 + 2c_0$ in expectation. After $T$ rounds, the expected communication cost is roughly $2(n-1) + 2(1 + c_0)T$. As a result, we can use the number of rounds to measure communication complexity.

**The linear-span protocol and information set.** In Definition F.1, we focus on loopless algorithms based on the linear-span protocol. One can check that many methods, e.g., KatyushaX [6], L-SVRG and L-Katyusha [34], Loopless SARAH [42] and SVRP[17] [33], satisfy our definition. And this class of algorithms is sufficiently large in that the upper and lower bounds have matched for most cases

---

[14]When analyzing computational complexity, this distribution can also depend on the smoothness of each component function [57, 6]

[15]See Proposition A.2.

[16]Such a modification makes no difference to subsequent analysis. See Remark F.3.

[17]For Catalyzed SVRP in their paper, we can slightly modify it without affecting the gradient or communication complexity. Specifically, we remove the full gradient step at the beginning of the inner loop and do not update the current point until the full gradient is calculated. The number of additional communication rounds is $1/p = \Theta(n)$ in expectation, as long as $p = \Theta(1/n)$. Since in each inner loop, the algorithm must calculate the full gradient, whose gradient or communication complexity is also $\Theta(n)$, such a modification would not affect the total complexity.

[25]. Built on the linear-span protocol, the information set $\mathcal{I}_{t+1}$, a linear subspace of the whole space, gathers all the gradient and proximal information obtained by $t$ rounds of communication and includes all the possible points generated by the algorithm after round $t$. Clearly, the sequence $\{\mathcal{I}_t\}_{t=0}^{T}$ is nondecreasing in the sense that $\mathcal{I}_t \subseteq \mathcal{I}_{t'}$ for any $t' > t$.

## F.2 Details of Section 4.2

Recall that in Section 4.2, we consider the following class of matrices

$$
\boldsymbol{B}(m, \zeta) = \begin{bmatrix} 1 & -1 & & & \\ & 1 & -1 & & \\ & & \ddots & \ddots & \\ & & & 1 & -1 \\ & & & & \zeta \end{bmatrix} \in \mathbb{R}^{m \times m}.
$$

And one can check the matrix $\boldsymbol{A}(m, \zeta)$ is a tridiagonal matrix, i.e.,

$$
\boldsymbol{A}(m, \zeta) := \boldsymbol{B}(m, \zeta)^{\top} \boldsymbol{B}(m, \zeta) := \begin{bmatrix} 1 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & \zeta^2 + 1 \end{bmatrix} \in \mathbb{R}^{m \times m}.
$$

With the hard instance constructed in (13), we have the following lemma, which is a modification of Lemma 6.1 in Han et al. [25] in that the partitions of the index sets are slightly different.

**Lemma F.2** *Suppose that $n \geq 3$, $m \geq 3$, $\gamma$ is an arbitrarily positive number and $\boldsymbol{x} \in \mathcal{F}_k$ for some $0 \leq k < m$. If $k = 0$, we have*

$$
\nabla r_i(\boldsymbol{x}), \ \mathrm{prox}_{r_i}^{\gamma}(\boldsymbol{x}) \in \begin{cases} \mathcal{F}_1, & \text{if } i = 1, \\ \mathcal{F}_0, & \text{otherwise.} \end{cases}
$$

*If $k > 0$, we have*

$$
\nabla r_i(\boldsymbol{x}), \ \mathrm{prox}_{r_i}^{\gamma}(\boldsymbol{x}) \in \begin{cases} \mathcal{F}_{k+1}, & \text{if } k \in \mathcal{L}_i, \\ \mathcal{F}_k, & \text{otherwise.} \end{cases}
$$

*Here $\{\mathcal{F}_k\}_{k=0}^{m}$ are defined as $\mathcal{F}_0 = \{\boldsymbol{0}\}$ and $\mathcal{F}_k = \mathrm{span}\{\boldsymbol{e}_1, \boldsymbol{e}_2, \dots, \boldsymbol{e}_k\}$ for $1 \leq k \leq m$, and we omit the parameters of $r_i$ to simplify the notation.*

Lemma F.2 tells us that if the current point $\boldsymbol{x}$ lies in some subspace of $\mathbb{R}^m$, only one component can provide the information of the next dimension by gradient or proximal information. In this sense, PIFO cannot provide more information than IFO. Thus, when we focus on the communication complexity of an algorithm, it makes no difference to distinguish between IFO and PIFO. And we can assume the algorithm has access to PIFO without loss of generality. Moreover, when $k > 0$, the oracle of $r_1$ can never provide any information on the next dimension. The proof of Lemma F.2 is deferred to Appendix F.4.

With Lemma F.2, Lemma 4.3 is a natural corollary and the proof is deferred to Appendix F.5.

**Remark F.3** *Recall that in steps 2 and 3, the difference between $\mathcal{I}_{t+1}^1$ and $\mathcal{I}_{t+1}^0$ only resides in $h_{f_1}^{\mathrm{P}}(\boldsymbol{x}, \gamma)$ (or $h_{r_1}^{\mathrm{P}}(\boldsymbol{x}, \gamma)$ when we consider problem (13)) for $\boldsymbol{x} \in \mathcal{I}_t^0$, while Lemma F.2 implies that $h_{r_1}^{\mathrm{P}}(\boldsymbol{x}, \gamma)$ would not expand the information set as long as $\boldsymbol{x} \neq \boldsymbol{0}$ [18]. This demonstrates that applying (26) recursively would not affect the analysis of communication complexity.*

The next result is a corollary of Lemma 4.3 and the proof is deferred to Appendix F.6.

---

[18] In the proof of Lemma 4.3 in Appendix F.5, we show that $\mathcal{I}_{t+1}^1 = \mathcal{I}_{t+1}^0$

**Corollary F.4** *Define the random variables $T_0 = -1$,*

$$T_k := \min_t\{t : t > T_{k-1}, 3k-2 \in \mathcal{L}_{i_t} \text{ or } a_t = 1\} \text{ for } 1 \le k \le (m-1)/3, \tag{29}$$

*and $Y_k := T_k - T_{k-1}$. Then we have (i) $\mathcal{I}_{t+1} \subseteq \mathcal{F}_{3k-2}$ for any $t < T_k$; (ii) the $Y_k$ are mutually independent; (iii) $Y_k \sim \mathrm{Geom}(q_{k'} + p - pq_{k'})$ with $k' \equiv 3k-1 \pmod{(n-1)}$, $2 \le k' \le n$.*

Corollary F.4 claims that $T_k$ is the smallest index of the communication round after which the information set can be expanded to $\mathcal{F}_{3k+2}$. Moreover, $T_k$ can be decomposed into the sum of independent geometric random variables. With Lemma A.6, which gives a concentration result for the sum of geometric random variables, we have the following proposition, whose proof is deferred to Appendix F.7.

**Proposition F.5** *Let $0 \le M \le (m-2)/3$ and $N = \frac{n(M+1)}{4(q_0+c_0)} + 1$ with $q_0$ and $c_0$ defined in Definition 4.1. Suppose we use an algorithm $\mathcal{A}$ satisfying Definition F.1 to solve problem (13). After $N$ round of communication, the algorithm obtains the information set $\mathcal{I}_N$. Then we have $\mathcal{I}_N \subseteq \mathcal{F}_{3M+1} \subset \mathcal{F}_{m-1}$. Moreoever, if*

$$\min_{\boldsymbol{x} \in \mathcal{F}_{3M+1}} r(\boldsymbol{x}) - \min_{\boldsymbol{x} \in \mathbb{R}^m} r(\boldsymbol{x}) \ge 9\epsilon, \tag{30}$$

*we have*

$$\mathbb{E}\min_{\boldsymbol{x} \in \mathcal{I}_N} r(\boldsymbol{x}) - \min_{\boldsymbol{x} \in \mathbb{R}^m} r(\boldsymbol{x}) \ge \epsilon.$$

Proposition F.5 specifies the number of communication rounds needed to find an $\epsilon$-suboptimal solution under the condition (30). Roughly speaking, the condition requires that the exact solution of problem (13) does not lie in some subspace of $\mathcal{F}_{m-1}$

Now we come back to the hard instance (13). Recall that $r$ is $c$-strongly convex and $r_i$'s satisfy $\sqrt{8n+4}$-aveSS. We need to properly scale the function class $\{r_i\}_{i=1}^n$ such that it satisfies Assumption 1. Note that rescaling does not influence Lemma F.2. Thus Proposition F.5 still holds for any rescaled version of problem (13). Specially, we consider the following problem

$$\min_{\boldsymbol{x} \in \mathbb{R}^m} f^{\mathrm{h}}(\boldsymbol{x}) := \frac{1}{n}\sum_{i=1}^n f_i^{\mathrm{h}}(\boldsymbol{x}) \quad \text{where} \quad f_i^{\mathrm{h}}(\boldsymbol{x}) := \lambda\, r(\boldsymbol{x}/\beta; m, \zeta, c),$$

$$\lambda = \frac{4\Delta}{\rho-1}, \; \beta = \frac{4}{\rho-1}\sqrt{\frac{\Delta}{\mu(\rho+1)}}, \; \zeta = \sqrt{\frac{2}{1+\rho}} \text{ and } c = \frac{4}{\rho^2-1} \text{ with } \rho = \sqrt{\frac{2\delta/\mu}{\sqrt{2n+1}}}+1. \tag{31}$$

Here $n, \delta, \mu$ and $\Delta$ are given parameters. As shown in the next Proposition, $\delta$ is the AveSS parameter, $\mu$ is the strong convexity parameter and $\Delta$ is the function value gap between the initial point and the solution.

**Proposition F.6** *The problem defined in (31) with $n \ge 3$ and $m \ge 3$ has the following properties.*

1. *$f^{\mathrm{h}}$ is $\mu$-strongly convex and $f_i^{\mathrm{h}}$'s satisfy $\delta$-AveSS.*

2. *Let $q = \frac{\rho-1}{\rho+1}$. The minimizer of $f^{\mathrm{h}}$ is $\boldsymbol{x}_* = \frac{\beta(\rho+1)}{2}(q, q^2, \ldots, q^m)^\top$ and $f^{\mathrm{h}}(\boldsymbol{0}) - f^{\mathrm{h}}(\boldsymbol{x}_*) = \Delta$.*

3. *For $0 \le k \le m-1$, we have*

$$\min_{\boldsymbol{x} \in \mathcal{F}_k} f^{\mathrm{h}}(\boldsymbol{x}) - f^{\mathrm{h}}(\boldsymbol{x}_*) \ge \Delta q^{2k} \text{ and } \min_{\boldsymbol{x} \in \mathcal{F}_k} \|\boldsymbol{x} - \boldsymbol{x}_*\|^2 \ge \frac{4\Delta}{\mu(\rho+1)}q^{2k}. \tag{32}$$

Property 2 shows that the minimizer of problem (31) has all elements nonzero. Thus, it does not lie in any subspace $\mathcal{F}_k$ for $k < m$. As a result, we cannot obtain an approximate solution up to an arbitrarily small accuracy, unless we get an iterate with the last element nonzero, as claimed by Property 3. This implies problem (31) satisfies the condition 30.

Combining Propositions F.5 and F.6, we can establish the lower bound of the communication complexity.

**Theorem F.7 (Formal version of Theorem 4.4)** *Suppose we use any algorithm $\mathcal{A}$ satisfying Definition F.1 to solve the minimization problem (31) and the following conditions hold*

$$n \geq 3 \;\; and \;\; \epsilon \leq \frac{\Delta}{9} \cdot q^3, \;\; with \; q = \frac{\rho - 1}{\rho + 1}, \rho = \sqrt{\frac{2\delta/\mu}{\sqrt{2n+1}}} + 1.$$

*Set $m = \left\lfloor \frac{\log(\Delta/(9\epsilon))}{2\log(1/q)} + 2 \right\rfloor$. In order to find $\hat{x}$ such that $\mathbb{E} f^{\mathrm{h}}(\hat{x}) - \min_{x \in \mathbb{R}^m} f^{\mathrm{h}}(x) < \epsilon$, the communication complexity in expectation is*

$$\begin{cases} \Omega\left(n + n^{3/4}\sqrt{\delta/\mu}\log(1/\epsilon)\right), & for \; \frac{\delta}{\mu} = \Omega(\sqrt{n}), \\ \Omega\left(n + \frac{n\log(1/\epsilon)}{1 + (\log(\mu\sqrt{n}/\delta))_+}\right), & for \; \frac{\delta}{\mu} = \mathcal{O}(\sqrt{n}). \end{cases}$$

Recall that in (32), $\min_{x \in \mathcal{F}_k} f^{\mathrm{h}}(x) - f^{\mathrm{h}}(x_*)$ and $\min_{x \in \mathcal{F}_k} \|x - x_*\|^2$ are both lower bounded by $q^{2k}$ multiplied with some constants. If we want to find $\hat{x}$ such that $\mathbb{E} \|\hat{x} - x_*\|^2 < \epsilon$, the communication complexity is the same. The proof of Theorem F.7 is deferred to Appendix F.9.

## F.3 Proof of Proposition 4.2

Proof: For convenience of notation, we omit the dependence of $r_i$, $r$, $B$ and $b_l$ on the parameters $m$, $\zeta$ and $c$. With the definition of $\{r_i\}_{i=1}^n$ and $r$, we have

$$\nabla(r_i - r)(x) = \begin{cases} -\sum_{l=1}^m b_l b_l^\top x - (n-1)e_1, & i = 1, \\ n\sum_{l \in \mathcal{L}_i} b_l b_l^\top x - \sum_{l=1}^m b_l b_l^\top x + e_1, & i \neq 1. \end{cases} \tag{33}$$

From the definition of $b_l$, we have

$$b_l^\top b_l = \begin{cases} 2, & 1 \leq l \leq m-1, \\ \zeta^2, & l = m, \end{cases} \qquad b_l^\top b_{l+1} = \begin{cases} -1, & 1 \leq l \leq m-2, \\ -\zeta, & l = m-1, \end{cases} \tag{34}$$

and $b_l^\top b_{l'} = 0$ for any $|l - l'| \geq 2$. Since $n \geq 3$, this implies $b_l^\top b_{l'} = 0$ for any $l, l' \in \mathcal{L}_i$ and $l \neq l'$. Define $b_0 = b_{m+1} = \mathbf{0}$ for ease of notation. Let $A_i = \sum_{l \in \mathcal{L}_i} b_l b_l^\top$, $A = \sum_{i=2}^n A_i = \sum_{l=1}^m b_l b_l^\top$. Then for any $x, y \in \mathbb{R}^m$ and $u = x - y$, we have

$$\sum_{i=1}^n \|\nabla(r_i - r)(x) - \nabla(r_i - r)(y)\|^2 \overset{(33)}{=} \|Au\|^2 + \sum_{i=2}^n \|nA_iu - Au\|^2$$

$$= \|Au\|^2 + \sum_{i=2}^n n^2\|A_iu\|^2 - 2n\sum_{i=2}^n (A_iu)^\top Au + (n-1)\|Au\|^2 = n^2\sum_{i=2}^n \|A_iu\|^2 - n\|Au\|^2.$$

Note that by Eq. (34), $\|A_iu\|^2 = \left\|\sum_{l \in \mathcal{L}_i} b_l b_l^\top u\right\|^2 = \sum_{l \in \mathcal{L}_i} (b_l^\top u)^2 b_l^\top b_l$ and

$$\|Au\|^2 = \left\|\sum_{l=1}^m b_l b_l^\top u\right\|^2 = \sum_{l=1}^m (b_l^\top u)^2 b_l^\top b_l + (b_l^\top u)(b_{l+1}^\top u)b_{l+1}^\top b_l + (b_l^\top u)(b_{l-1}^\top u)b_{l-1}^\top b_l$$

$$= \sum_{l=1}^m (b_l^\top u)^2 b_l^\top b_l + 2(b_l^\top u)(b_{l+1}^\top u)b_{l+1}^\top b_l,$$

where the final equality uses $b_0 = b_{m+1} = \mathbf{0}$. Hence, we get

$$\frac{1}{n}\sum_{i=1}^n \|\nabla(r_i - r)(x) - \nabla(r_i - r)(y)\|^2$$

$$= n\sum_{l=1}^m (b_l^\top u)^2 b_l^\top b_l - \left[\sum_{l=1}^m (b_l^\top u)^2 b_l^\top b_l + 2(b_l^\top u)(b_{l+1}^\top u)b_{l+1}^\top b_l\right]$$

$$= (n-1)\sum_{l=1}^m (b_l^\top u)^2 b_l^\top b_l - 2\sum_{l=1}^m (b_l^\top u)(b_{l+1}^\top u)b_{l+1}^\top b_l. \tag{35}$$

Recall that $0 < \zeta \leq \sqrt{2}$. Then (34) implies $\boldsymbol{b}_l^\top \boldsymbol{b}_l \leq 2$ and $|\boldsymbol{b}_l^\top \boldsymbol{b}_{l+1}| \leq \sqrt{2}$ for any $l$. Substituting these into (35) and using Cauchy's inequality, we have

$$\frac{1}{n} \sum_{i=1}^n \|\nabla(r_i - r)(\boldsymbol{x}) - \nabla(r_i - r)(\boldsymbol{y})\|^2 \leq 2(n-1) \sum_{l=1}^m (\boldsymbol{b}_l^\top \boldsymbol{u})^2 + \sqrt{2} \sum_{l=1}^m \left[ (\boldsymbol{b}_l^\top \boldsymbol{u})^2 + (\boldsymbol{b}_{l+1}^\top \boldsymbol{u})^2 \right]$$

$$\leq \left( 2n + 2\sqrt{2} - 2 \right) \sum_{l=1}^m (\boldsymbol{b}_l^\top \boldsymbol{u})^2 \leq (2n+1) \sum_{l=1}^m (\boldsymbol{b}_l^\top \boldsymbol{u})^2.$$

Notice that $(\boldsymbol{b}_l^\top \boldsymbol{u})^2 = (u_l - u_{l+1})^2 \leq 2(u_l^2 + u_{l+1}^2)$ for $1 \leq l \leq m-1$ and $(\boldsymbol{b}_m^\top \boldsymbol{u})^2 = (\zeta u_m)^2 \leq 2u_m^2$. This implies

$$\frac{1}{n} \sum_{i=1}^n \|\nabla(r_i - r)(\boldsymbol{x}) - \nabla(r_i - r)(\boldsymbol{y})\|^2 \leq 4(2n+1) \|\boldsymbol{u}\|^2 = (8n+4) \|\boldsymbol{x} - \boldsymbol{y}\|^2.$$

As a result, $r_i$'s satisfy $\sqrt{8n+4}$-AveSS. $\qquad\square$

### F.4 Proof of Lemma F.2

Proof: For convenience of notation, we omit the dependence of $r_i$, $r$, $\boldsymbol{B}$ and $\boldsymbol{b}_l$ on the parameters $m$, $\zeta$ and $c$.

1) First, we focus on the gradient of the $r_i$. Recall that

$$r_i(\boldsymbol{x}) = \begin{cases} \frac{c}{2} \|\boldsymbol{x}\|^2 - n \langle \boldsymbol{e}_1, \boldsymbol{x} \rangle, & \text{for } i = 1, \\ \frac{n}{2} \sum_{l \in \mathcal{L}_i} \boldsymbol{x}^\top \boldsymbol{b}_l \boldsymbol{b}_l^\top \boldsymbol{x} + \frac{c}{2} \|\boldsymbol{x}\|^2, & \text{for } i \neq 1. \end{cases} \Rightarrow \nabla r_i(\boldsymbol{x}) = \begin{cases} c\boldsymbol{x} - n\boldsymbol{e}_1, & \text{for } i = 1, \\ n \sum_{l \in \mathcal{L}_i} \boldsymbol{b}_l \boldsymbol{b}_l^\top \boldsymbol{x} + c\boldsymbol{x}, & \text{for } i \neq 1. \end{cases}$$

(i): If $\boldsymbol{x} \in \mathcal{F}_0$, i.e., $\boldsymbol{x} = \boldsymbol{0}$, we have $\nabla r_1(\boldsymbol{0}) = -n\boldsymbol{e}_1 \in \mathcal{F}_1$ and $\nabla r_i(\boldsymbol{0}) = \boldsymbol{0}$ for $i \neq 1$. (ii): If $\boldsymbol{x} = (x_1, \ldots, x_m) \in \mathcal{F}_k$ for $1 \leq k < m$, we have $\nabla r_1(\boldsymbol{x}) = c\boldsymbol{x} - n\boldsymbol{e}_1 \in \mathcal{F}_k$. As for $i \neq 1$, we need to examine $\boldsymbol{b}_l \boldsymbol{b}_l^\top \boldsymbol{x}$. One can check

$$\boldsymbol{b}_l \boldsymbol{b}_l^\top \boldsymbol{x} = \begin{cases} (x_l - x_{l+1})(\boldsymbol{e}_l - \boldsymbol{e}_{l+1}), & 1 \leq l \leq m-1, \\ \zeta^2 x_m \boldsymbol{e}_m, & l = m. \end{cases}$$

For $\boldsymbol{x} \in \mathcal{F}_k$, we have

$$\boldsymbol{b}_l \boldsymbol{b}_l^\top \boldsymbol{x} \in \begin{cases} \mathcal{F}_k, & l \neq k, \\ \mathcal{F}_{k+1}, & l = k. \end{cases} \tag{36}$$

As a result, if $k \in \mathcal{L}_i$, then $\nabla r_i(\boldsymbol{x}) \in \mathcal{F}_{k+1}$; otherwise, $\nabla r_i(\boldsymbol{x}) \in \mathcal{F}_k$.

2) Now we turn to the proximal operator. (i) For $i = 1$, it is easy to verify $\mathrm{prox}_{r_1}^\gamma(\boldsymbol{x}) = (1/\gamma + c)^{-1}(\boldsymbol{x}/\gamma + n\boldsymbol{e}_1)$. Thus, if $\boldsymbol{x} \in \mathcal{F}_0$, $\mathrm{prox}_{r_1}^\gamma(\boldsymbol{x}) \in \mathcal{F}_1$; if $\boldsymbol{x} \in \mathcal{F}_k$ for $k \geq 1$, $\mathrm{prox}_{r_1}^\gamma(\boldsymbol{x}) \in \mathcal{F}_k$. (ii) For $i \neq 1$, we define $\boldsymbol{u}_i := \mathrm{prox}_{r_i}^\gamma(\boldsymbol{x})$ for simplicity. Then $\boldsymbol{u}_i$ satisfies the following equation

$$\left[ n\gamma \boldsymbol{B}_i^\top \boldsymbol{B}_i + (c\gamma + 1) \boldsymbol{I} \right] \boldsymbol{u}_i = \boldsymbol{x}, \; \boldsymbol{B}_i := \sum_{l \in \mathcal{L}_i} \boldsymbol{e}_l \boldsymbol{b}_l^\top.$$

Note that $\boldsymbol{B}_i^\top \boldsymbol{B}_i = \sum_{l \in \mathcal{L}_i} \boldsymbol{b}_l \boldsymbol{b}_l^\top$. By the Sherman-Morrison-Woodbury formula, we get

$$\left( \boldsymbol{I} + \tilde{c} \boldsymbol{B}_i^\top \boldsymbol{B}_i \right)^{-1} = \boldsymbol{I} - \boldsymbol{B}_i^\top \left( \frac{1}{\tilde{c}} \boldsymbol{I} + \boldsymbol{B}_i \boldsymbol{B}_i^\top \right)^{-1} \boldsymbol{B}_i, \forall \tilde{c} \neq 0.$$

In the proof of Proposition 4.2, we have shown that $\boldsymbol{b}_l^\top \boldsymbol{b}_{l'} = 0$ for any $|l - l'| \geq 2$ and consequently $\boldsymbol{b}_l^\top \boldsymbol{b}_{l'} = 0$ for any $l, l' \in \mathcal{L}_i$ and $l \neq l'$. Thus, $\boldsymbol{B}_i \boldsymbol{B}_i^\top = \sum_{l \in \mathcal{L}_i} \boldsymbol{b}_l^\top \boldsymbol{b}_l \boldsymbol{e}_l \boldsymbol{e}_l^\top$ is a diagonal matrix. Then we can denote $\boldsymbol{D}_i = \left( \frac{c\gamma + 1}{n\gamma} \boldsymbol{I} + \boldsymbol{B}_i \boldsymbol{B}_i^\top \right)^{-1} = \sum_{l=1}^m d_{i,l} \boldsymbol{e}_l \boldsymbol{e}_l^\top$ and obtain

$$\boldsymbol{u}_i = \left[ n\gamma \boldsymbol{B}_i^\top \boldsymbol{B}_i + (c\gamma + 1) \boldsymbol{I} \right]^{-1} \boldsymbol{x} = \frac{1}{c\gamma + 1} \left( \boldsymbol{I} + \frac{n\gamma}{c\gamma + 1} \boldsymbol{B}_i^\top \boldsymbol{B}_i \right)^{-1} \boldsymbol{x} = \frac{\boldsymbol{x} - \boldsymbol{B}_i^\top \boldsymbol{D}_i \boldsymbol{B}_i \boldsymbol{x}}{c\gamma + 1}.$$

Then we have $\boldsymbol{B}_i^\top \boldsymbol{D}_i \boldsymbol{B}_i \boldsymbol{x} = \sum_{l \in \mathcal{L}_i} d_{i,l} \boldsymbol{b}_l \boldsymbol{b}_l^\top \boldsymbol{x}$. For $\boldsymbol{x} \in \mathcal{F}_k$, by (36), if $k \in \mathcal{L}_i$, then $\boldsymbol{u}_i \in \mathcal{F}_{k+1}$; otherwise, $\boldsymbol{u}_i \in \mathcal{F}_k$. This completes the proof. $\qquad\square$

## F.5 Proof of Lemma 4.3

Proof: Since we can assume $\boldsymbol{x}_0 = \boldsymbol{0}$, Lemma F.2 implies $\nabla r_1 \in \mathcal{F}_1$. Then from step 1, we have $\mathcal{I}_0 = \mathcal{F}_1$[19].

Then we focus on the second claim and examine how many dimensions of the information set we can increase after a round of communication.

Since we choose node 1 as the master node, we have access to $\nabla r_1$ and $\mathrm{prox}_{r_1}^\gamma$ in each communication round. Recall that we set $\mathcal{L}_1$ as the empty set. Lemma F.2 guarantees that the information provided by $r_1$ can never expand the information set unless the information set only contains $\boldsymbol{0}$. Thus, (26) does not affect the information set, i.e., $\mathcal{I}_{t+1}^1 = \mathcal{I}_{t+1}^0$ for any $t \geq 0$.

When $a_t = 0$, only (25) can expand the information set. By Lemma F.2, if $i_t$ satisfies $k \in \mathcal{L}_{i_t}$, we have $\mathcal{I}_{t+1}^1 = \mathcal{I}_{t+1}^0 \subseteq \mathcal{F}_{k+1}$. Otherwise, we still have $\mathcal{I}_{t+1}^1 = \mathcal{I}_{t+1}^0 \subseteq \mathcal{F}_k$. Then step 4 in Definition F.1 implies $\mathcal{I}_{t+1} = \mathcal{I}_{t+1}^1$.

When $a_t = 1$, from the above analysis, we always have $\mathcal{I}_{t+1} \subseteq \mathcal{F}_{k+1}$. By Lemma F.2, (27) could expand the information set by at most one dimension. It follows that $\tilde{\boldsymbol{y}}_{t+1} \in \mathcal{F}_{k+2}$. Using Lemma F.2 again yields $\mathcal{I}_{t+1} \subseteq \mathcal{F}_{k+3}$. □

## F.6 Proofs of Corollary F.4

Proof: We prove the first claim by induction on $t$. That is to say, we prove that for any integer $t \geq -1$, $\mathcal{I}_{t+1} \subseteq \mathcal{F}_{3k-2}$ for any $k$ satisfying $t < T_k$. Define $k(t)$ as the positive integer such that $T_{k(t)-1} \leq t < T_{k(t)}$. From the monotonicity of $\mathcal{F}_\cdot$, it suffices to prove $\mathcal{I}_{t+1} \subseteq \mathcal{F}_{3k(t)-2}$.

By Lemma F.2, we have $\mathcal{I}_0 = \mathcal{F}_1$ and $k(-1) = 1$. The claim holds for $t = -1$. Suppose that $\mathcal{I}_{t+1} \subseteq \mathcal{F}_{3k(t)-2}$. If $3k(t) - 2 \in \mathcal{L}_{i_{t+1}}$ or $a_{t+1} = 1$, Lemma 4.3 together with (29) implies $k(t+1) = k(t) + 1$ and $\mathcal{I}_{t+2} \subseteq \mathcal{F}_{3k(t)+1} = \mathcal{F}_{3k(t+1)-1}$. Otherwise, we still have $k(t+1) = k(t)$ and $\mathcal{I}_{t+2} \subseteq \mathcal{F}_{3k(t)-1} = \mathcal{F}_{3k(t+1)-1}$.

For the second claim, the independence of $\{Y_k\}_{k \geq 1}$ is natural consequence of the independence of $\{(i_t, a_t)\}_{t \geq 1}$.

For the last one, note that $3k - 2 \in \mathcal{L}_{i_t}$ is equivalent to $i_t \equiv 3k-1 (\mathrm{mod}\ (n-1))$ for $2 \leq i_t \leq n$. Then we have for $k' \equiv 3k-1 (\mathrm{mod}\ (n-1)), 2 \leq k' \leq n$,

$$
\begin{aligned}
\mathbb{P}(T_k - T_{k-1} = s) &= \mathbb{P}(i_{T_{k-1}+1} \neq k', \ldots, i_{T_{k-1}+s-1} \neq k', a_{T_{k-1}+1} = \cdots = a_{T_{k-1}+s-1} = 0, \\
&\qquad i_{T_{k-1}+s} = k' \text{ or } a_{T_{k-1}+s} = 1) \\
&\overset{(i)}{=} [(1 - q_{k'})(1-p)]^{s-1} [1 - (1 - q_{k'})(1-p)],
\end{aligned}
$$

where $(i)$ is due to the independence of $\{(i_t, a_t)\}_{t \geq 1}$. So $Y_k = T_k - T_{k-1}$ is a geometric random variable with success probability $1 - (1 - q_{k'})(1 - p) = q_{k'} + p - q_{k'}p$. □

## F.7 Proof of Proposition F.5

Proof: By Corollary F.4, if $N - 1 < T_{M+1}$, then $\mathcal{I}_N \subseteq \mathcal{F}_{3M+1} \subseteq \mathcal{F}_{m-1}$. Thus we have

$$
\mathbb{E} \min_{\boldsymbol{x} \in \mathcal{I}_N} r(\boldsymbol{x}) - \min_{\boldsymbol{x} \in \mathbb{R}^m} r(\boldsymbol{x}) \geq \mathbb{E} \left[ \min_{\boldsymbol{x} \in \mathcal{I}_N} r(\boldsymbol{x}) - \min_{\boldsymbol{x} \in \mathbb{R}^m} r(\boldsymbol{x}) \middle| N - 1 < T_{M+1} \right] \mathbb{P}(N - 1 < T_{M+1})
$$

$$
\geq \mathbb{E} \left[ \min_{\boldsymbol{x} \in \mathcal{F}_{3M+1}} r(\boldsymbol{x}) - \min_{\boldsymbol{x} \in \mathbb{R}^m} r(\boldsymbol{x}) \middle| N - 1 < T_{M+1} \right] \mathbb{P}(N - 1 < T_{M+1}) \geq 9\epsilon \, \mathbb{P}(N - 1 < T_{M+1}).
$$

By Corollary F.4 again, $T_{M+1}$ can be written as $T_{M+1} = \sum_{l=1}^{M+1} Y_l$, where $\{Y_l\}_{1 \leq l \leq M+1}$ are independent random variables, and $Y_l \sim \mathrm{Geom}(\tilde{q}_l)$ with $\tilde{q}_l = q_{l'} + p - pq_{l'}$, $l' \equiv 3l-1 (\mathrm{mod}\ (n-1))$ and $2 \leq l' \leq n$. Moreover, Definition F.1 guarantees $\max_{2 \leq l' \leq n} q_l' \leq q_0/n$ and $p \leq c_0/n$. Then we have $\sum_{l=1}^{M+1} \tilde{q}_l \leq (q_0 + c_0)(M+1)/n$. Therefore, by Lemma A.6, we have

$$
\mathbb{P}(T_{M+1} > N - 1) = \mathbb{P}\left( \sum_{l=1}^{M+1} Y_l > \frac{n(M+1)}{4(q_0 + c_0)} \right) \geq \mathbb{P}\left( \sum_{l=1}^{M+1} Y_l > \frac{(M+1)^2}{4 \sum_{l=1}^{M+1} \tilde{q}_l} \right) \geq \frac{1}{9},
$$

---

[19]In the definition of the information set, each $f_i$ is replaced by $r_i$ here.

which implies our desired result. □

## F.8 Proof of Proposition F.6

Proof: **Property 1.** By Proposition 4.2, we have $f^{\mathrm{h}}$ is $\lambda c/\beta^2$-strongly convex and $f_i^{\mathrm{h}}$'s satisfy $\lambda\sqrt{8n+4}/\beta^2$-AveSS. One can check $\lambda c/\beta^2 = \mu$ and $\lambda\sqrt{8n+4}/\beta^2 = \delta$.

**Property 2.** Let $\xi := \lambda/\beta^2 = \mu(\rho^2 - 1)/4$. We have

$$f^{\mathrm{h}}(\boldsymbol{x}) = \frac{\xi}{2}\boldsymbol{x}^\top \boldsymbol{A}(m,\zeta)\,\boldsymbol{x} + \frac{\mu}{2}\|\boldsymbol{x}\|^2 - \xi\beta\,\langle \boldsymbol{e}_1, \boldsymbol{x}\rangle.$$

Letting $\nabla f^{\mathrm{h}}(\boldsymbol{x}) = \boldsymbol{0}$ yields $(\xi\boldsymbol{A}(m,\zeta) + \mu\boldsymbol{I})\,\boldsymbol{x} = \xi\beta\boldsymbol{e}_1$, or equivalently.

$$\begin{bmatrix} 1 + \frac{\mu}{\xi} & -1 & & & \\ -1 & 2 + \frac{\mu}{\xi} & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 + \frac{\mu}{\xi} & -1 \\ & & & -1 & \zeta^2 + 1 + \frac{\mu}{\xi} \end{bmatrix} \boldsymbol{x} = \begin{bmatrix} \beta \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}. \tag{37}$$

Since $q = \frac{\rho-1}{\rho+1}$, we get $2 + \frac{\mu}{\xi} = \frac{2\rho^2+2}{\rho^2-1} = q + \frac{1}{q}$ and $\zeta^2 + 1 + \frac{\mu}{\xi} = \frac{\rho+1}{\rho-1} = \frac{1}{q}$. We solve (37) by

$$x_{m-1} - \frac{x_m}{q} = 0,\ x_k - \left(q + \frac{1}{q}\right)x_{k+1} + x_{k+2} = 0,\ k \in [m-2],\ \left(q + \frac{1}{q} - 1\right)x_1 - x_2 = \beta.$$

Thus, $\boldsymbol{x}_* = \frac{\beta}{1-q}(q, q^2, \ldots, q^m)^\top$ and $f^{\mathrm{h}}(\boldsymbol{x}_*) = -\frac{\xi\beta\langle \boldsymbol{e}_1, \boldsymbol{x}_*\rangle}{2} = -\frac{\xi\beta^2 q}{2(1-q)} = -\frac{\lambda(\rho-1)}{4} \overset{(31)}{=} -\Delta$.

**Property 3.** If $\boldsymbol{x} \in \mathcal{F}_k$, $1 \leq k < m$, then $x_{k+1} = x_{k+2} = \cdots = x_m = 0$. Let $\boldsymbol{y}$ denote the first $k$ coordinates of $\boldsymbol{x}$ and $\boldsymbol{A}_k$ denote the first $k$ rows and columns of $\boldsymbol{A}(m,\zeta)$. Then for any $\boldsymbol{x} \in \mathcal{F}_k$, we can rewrite $f^{\mathrm{h}}(\boldsymbol{x})$ as

$$\hat{f}^{\mathrm{h}}(\boldsymbol{y}) := f^{\mathrm{h}}(\boldsymbol{x}) = \frac{\xi}{2}\boldsymbol{y}^\top \boldsymbol{A}_k \boldsymbol{y} + \frac{\mu}{2}\|\boldsymbol{x}\|^2 - \xi\beta\,\langle \hat{\boldsymbol{e}}_1, \boldsymbol{y}\rangle,$$

where $\hat{\boldsymbol{e}}_1$ is the first $k$ coordinates of $\boldsymbol{e}_1$. Let $\nabla f_k(\boldsymbol{y}) = \boldsymbol{0}$, that is

$$\begin{bmatrix} 1 + \frac{\mu}{\xi} & -1 & & & \\ -1 & 2 + \frac{\mu}{\xi} & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 + \frac{\mu}{\xi} & -1 \\ & & & -1 & 2 + \frac{\mu}{\xi} \end{bmatrix} \boldsymbol{y} = \begin{bmatrix} \beta \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}.$$

Similarly, we need to solve

$$x_{m-1} = \left(q + \frac{1}{q}\right)x_m,\ x_k - \left(q + \frac{1}{q}\right)x_{k+1} + x_{k+2} = 0,\ k \in [m-2],\ \left(q + \frac{1}{q} - 1\right)x_1 - x_2 = \beta.$$

By some computation, one can check the solution is

$$\boldsymbol{y}_* = \frac{\beta q^{k+1}}{2(q-1)(1+q^{2k+1})}\left(q^{-k} - q^k, q^{-(k-1)} - q^{k-1}, \ldots, q^{-1} - q^1\right)^\top.$$

Thus, we have

$$f^{\mathrm{h}}(\boldsymbol{y}_*) = -\frac{\xi\beta\,\langle \hat{\boldsymbol{e}}_1, \boldsymbol{y}_*\rangle}{2} = \frac{\xi\beta^2 q}{2(1-q)}\cdot\frac{1-q^{2k}}{1+q^{2k+1}} = \frac{\lambda(\rho-1)}{4}\cdot\frac{1-q^{2k}}{1+q^{2k+1}} = \frac{(1-q^{2k})\Delta}{1+q^{2k+1}},$$

and by $q < 1$, we further have that

$$\min_{\boldsymbol{x}\in\mathcal{F}_k} f^{\mathrm{h}}(\boldsymbol{x}) - \min_{\boldsymbol{x}\in\mathbb{R}^m} f^{\mathrm{h}}(\boldsymbol{x}) = f^{\mathrm{h}}(\boldsymbol{y}_*) - f^{\mathrm{h}}(\boldsymbol{x}_*) = \Delta\left(1 - \frac{1-q^{2k}}{1+q^{2k+1}}\right) = \frac{(1+q)q^{2k}\Delta}{1+q^{2k+1}} \geq \Delta q^{2k}.$$

Moreover, recall that $\boldsymbol{x}_* = \frac{\beta(\rho+1)}{2}(q, q^2, \ldots, q^m)^\top$. Then we have

$$\min_{\boldsymbol{x}\in\mathcal{F}_k} \|\boldsymbol{x} - \boldsymbol{x}^\star\|^2 = \frac{\beta^2(\rho+1)^2}{4}\sum_{i=k+1}^m q^{2i} \geq \frac{\beta^2(\rho+1)^2 q^2}{4}q^{2k} = \frac{4\Delta}{\mu(\rho+1)}q^{2k}.$$

This completes the proof. □

## F.9 Proof of Theorem F.7

Proof: Let $q = \frac{\rho - 1}{\rho + 1}$ and $M = \left\lfloor \frac{\log(\Delta/9\epsilon)}{6 \log 1/q} - \frac{1}{3} \right\rfloor$. From the condition on $\epsilon$ and the definition of $m$, one can check $0 \le M \le (m-2)/3$ and $m \ge 3$. Moreover, we have $3M + 1 \le \frac{\log(9\epsilon/\Delta)}{2 \log q}$. Then by Proposition F.5, after $N = \frac{n(M+1)}{4(q_0 + c_0)} + 1$ rounds of communication, the information set satisfies $\mathcal{I}_N \subseteq \mathcal{F}_{3M+1} \subseteq \mathcal{F}_{m-1}$. The third property of Propostion F.6 implies

$$\min_{\boldsymbol{x} \in \mathcal{F}_{3M+1}} f^{\mathrm{h}}(\boldsymbol{x}) - \min_{\boldsymbol{x} \in \mathbb{R}^m} f(\boldsymbol{x}) \ge \Delta q^{6M+2} \ge 9\epsilon.$$

Then by Proposition F.5 again, in order to find $\hat{\boldsymbol{x}}$ such that $\mathbb{E} f^{\mathrm{h}}(\hat{\boldsymbol{x}}) - \min_{\boldsymbol{x} \in \mathbb{R}^m} f^{\mathrm{h}}(\boldsymbol{x}) < \epsilon$, the algorithm $\mathcal{A}$ needs at least $N$ communication rounds.

Now we give a lower bound $N$. According to whether $2\delta/\mu$ is larger than $\sqrt{2n+1}$, we divide the analysis into two cases.

**Case 1: $2\delta/\mu \ge \sqrt{2n+1}$.** Then $\rho \ge \sqrt{2}$. By inequality $0 < \log(1+x) \le x, \forall x > 0$, we get

$$\frac{1}{\log \frac{1}{q}} = \frac{1}{\log\left(1 + \frac{2}{\rho - 1}\right)} \ge \frac{1}{\frac{2}{\rho - 1}} = \frac{\rho - 1}{2} = \frac{1}{2}\left(\sqrt{\frac{2\delta/\mu}{\sqrt{2n+1}} + 1} - 1\right) \ge \frac{\sqrt{2\delta/\mu}}{6\sqrt[4]{2n+1}},$$

where the final inequality uses $\sqrt{t+1} - 1 \ge \sqrt{t}/3, \forall t \ge 1$. Moreover, the condition on $\epsilon$ implies $\log \frac{\Delta}{9\epsilon} \ge 3 \log \frac{1}{q}$. Then we have

$$M + 1 \ge \frac{\log \frac{\Delta}{9\epsilon}}{6 \log \frac{1}{q}} - \frac{1}{3} \ge \frac{\log \frac{\Delta}{9\epsilon}}{18 \log \frac{1}{q}} \ge \frac{\sqrt{2\delta/\mu}}{108\sqrt[4]{2n+1}} \log \frac{\Delta}{9\epsilon} = \Omega\left(\frac{\sqrt{\delta/\mu}}{n^{1/4}} \log \frac{1}{\epsilon}\right).$$

It follows that $N = \frac{n(M+1)}{4(q_0 + c_0)} + 1 = \Omega\left(n^{3/4}\sqrt{\delta/\mu}\log(1/\epsilon)\right)$. The total communication cost in expectation is of the order $\Theta(n + N) = \Omega\left(n + n^{3/4}\sqrt{\delta/\mu}\log(1/\epsilon)\right)$.

**Case 2: $2\delta/\mu < \sqrt{2n+1}$.** In this case, we have $1 \le \rho < \sqrt{2}$ and consequently

$$\log \frac{1}{q} = \log\left(1 + \frac{2}{\rho - 1}\right) \overset{(i)}{\le} \log\left(1 + \frac{3\sqrt{2n+1}}{\delta/\mu}\right) \le \log\left(\frac{7\mu\sqrt{2n+1}}{2\delta}\right) \le 2 + \log\left(\frac{\mu\sqrt{2n+1}}{2\delta}\right),$$

where $(i)$ uses $\sqrt{t+1} - 1 \ge t/3, \forall 0 < t \le 1$ and $\rho = \sqrt{\frac{2\delta/\mu}{\sqrt{2n+1}} + 1}$, Moreover, the condition on $\epsilon$ implies $\log \frac{\Delta}{9\epsilon} \ge 3 \log \frac{1}{q}$. Then we have

$$M + 1 \ge \frac{\log \frac{\Delta}{9\epsilon}}{6 \log \frac{1}{q}} - \frac{1}{3} \ge \frac{\log \frac{\Delta}{9\epsilon}}{18 \log \frac{1}{q}} = \Omega\left(\frac{\log(1/\epsilon)}{1 + (\log(\mu\sqrt{n}/\delta))_+}\right).$$

where $(a)_+$ denote $\max\{a, 0\}$. It follows that $N = \frac{n(M+1)}{4(q_0 + c_0)} + 1 = \Omega\left(\frac{n \log(1/\epsilon)}{1 + (\log(\mu\sqrt{n}/\delta))_+}\right)$. The total communication cost in expectation is of the order $\Theta(n + N) = \Omega\left(n + \frac{n \log(1/\epsilon)}{1 + (\log(\mu\sqrt{n}/\delta))_+}\right)$. $\square$

## G  Experiment Details

We show some detail of our numerical experiments in this section. The computation of problem-dependent parameters is defined as follows. Since the objective is

$$f(\boldsymbol{x}) = \frac{1}{n} \sum_{i=1}^{n} \left[ f_i(\boldsymbol{x}) := \frac{1}{m} \sum_{j=1}^{m} \left(\boldsymbol{z}_{i,j}^{\top} \boldsymbol{x} - y_{i,j}\right)^2 + \frac{\mu}{2} \|\boldsymbol{x}\|^2 \right],$$

Let $\boldsymbol{Z}_i = (\boldsymbol{z}_{i,1}, \cdots, \boldsymbol{z}_{i,m})/\sqrt{m/2}, \boldsymbol{y}_i = (y_{i,1}, \ldots, y_{i,m})^{\top}/\sqrt{m/2}$. We reformulate $f_i, i \in [n]$ into

$$f_i(\boldsymbol{x}) = \frac{1}{2} \left\|\boldsymbol{Z}_i^{\top} \boldsymbol{x} - \boldsymbol{y}_i\right\|^2 + \frac{\mu}{2} \|\boldsymbol{x}\|^2, \nabla^2 f_i(\boldsymbol{x}) = \boldsymbol{Z}_i \boldsymbol{Z}_i^{\top} + \mu \boldsymbol{I}_d.$$

Table 2: The choices of interpolation $\tau = s\tau_0$ in experiments.

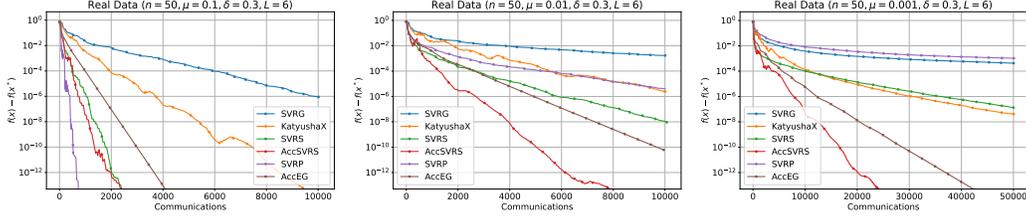|  |  | Katyusha X | | | AccSVRS | | |
|---|---|---|---|---|---|---|---|
| Synthetic data | $\mu$ | 1 | 0.1 | 0.01 | 1 | 0.1 | 0.01 |
|  | s | 1 | 2 | 5 | 2 | 5 | 10 |
| Real data | $\mu$ | 0.1 | 0.01 | 0.001 | 0.1 | 0.01 | 0.001 |
|  | s | 1 | 1 | 2 | 2 | 0.5 | 0.5 |



Figure 2: Numerical experiments on real data. The corresponding coefficients are shown in the title of each graph. We plot the function gap on a log scale versus the number of communication steps, where one exchange of vectors counts as a communication step.

Thus, we obtain the smoothness of each $f_i$ is $L_i = \|\boldsymbol{Z}_i\|^2 + \mu$, and $L := \max_{i \in [n]} L_i$. Obviously, $f$ is $\mu$-strongly convex.

For the synthetic data, we first generate a random symmetric matrix $\boldsymbol{Z}_0 \in \mathbb{R}^{d \times d}$ with $d = 100$ and $\|\boldsymbol{Z}_0\| = 3000$, then we add a perturbed symmetric matrix $\boldsymbol{N}_i, \forall i \in [n]$ with $n = 400, \|\boldsymbol{N}_i\| \approx 30$ to obtain $\boldsymbol{Z}_i = \boldsymbol{Z}_0 + \boldsymbol{N}_i$. We also add a correction $\lambda_{\min}(\boldsymbol{Z}_i)\boldsymbol{I}_d$ to $\boldsymbol{Z}_i$ to further make $\boldsymbol{Z}_i \succeq 0$. Finally, we recompute the center matrix $\boldsymbol{Z} = \sum_{i=1}^{n} \boldsymbol{Z}_i/n$ and $\delta$-average similarity coefficient following AveHS in Eq. (6) as

$$\delta = \sqrt{\frac{1}{n} \sum_{i=1}^{n} \|\boldsymbol{Z}_i - \boldsymbol{Z}\|^2}.$$

We use the analytic solution obtained by the proximal step since

$$\operatorname{prox}_{f_1}^{\theta}(\boldsymbol{x}_0) := \underset{\boldsymbol{x} \in \mathbb{R}^d}{\arg\min} f_1(\boldsymbol{x}) + \frac{1}{2\theta}\|\boldsymbol{x} - \boldsymbol{x}_0\|^2 = \left[\boldsymbol{Z}_1\boldsymbol{Z}_1^\top + \left(\mu + \frac{1}{\theta}\right)\boldsymbol{I}_d\right]^{-1}\left(\boldsymbol{Z}_1\boldsymbol{y} + \frac{\boldsymbol{x}_0}{\theta}\right).$$

For Katyusha X [6, Fact 4.2], and AccSVRS (Thm 3.6), we scale the interpolation coefficient $\tau = s\tau_0$ with $s \in \{0.5, 1, 2, 5, 10\}$ and $\tau_0$ is the theoretical value. The finally used scaling $s$ is shown in Table 2. The initial points of all methods are the same, which are sampled from $\operatorname{Unif}(\mathcal{S}^{d-1})$.

We also run the real data 'a9a' from LIBSVM library [16], where we split it into $n = 50$ datasets with the data size $m = 600$. The results are shown in Figure 2, and we can observe similar behavior of our methods.