

# Appendix

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Mining for Smoothness Information</b>	<b>2</b>
<b>3</b>	<b>Motivation and Contributions</b>	<b>3</b>
<b>4</b>	<b>New Communication-Efficient Methods Exploiting Matrix Smoothness</b>	<b>5</b>
4.1	DCGD+ . . . . .	6
4.2	Variance reduction: DIANA+ . . . . .	6
4.3	Acceleration with variance reduction: ADIANA+ . . . . .	7
<b>5</b>	<b>Improvements Over the Original Methods</b>	<b>8</b>
<b>6</b>	<b>Experiments</b>	<b>9</b>
<b>A</b>	<b>Conclusions, Extensions and Future Work</b>	<b>18</b>
<b>B</b>	<b>Limitations</b>	<b>19</b>
<b>C</b>	<b>Extra Experiments and Experimental Setup</b>	<b>20</b>
C.1	Proposed and usual sparsification techniques for the 3 distributed methods. . . . .	20
C.2	Variance reduction with new sparsification and importance sampling . . . . .	21
C.3	The effect of sparsification level $\tau$ on the convergence rate . . . . .	21
<b>D</b>	<b>Table of Frequently Used Notation</b>	<b>23</b>
<b>E</b>	<b>Theory in the Single Node Case: RCD as Sketched Gradient Descent (SkGD)</b>	<b>24</b>
E.1	‘NSync . . . . .	24
E.2	Sketched Gradient Descent (SkGD) . . . . .	25
E.3	CGD+ . . . . .	26
<b>F</b>	<b>Lower Bounds for Sketches as Linear Compression Operators</b>	<b>28</b>
F.1	Fixed sketches . . . . .	28
F.2	Random sketches . . . . .	29
F.3	Optimal sketches . . . . .	30
F.4	Random sketches with linear constraints . . . . .	31
F.5	Variance against communication trade-off . . . . .	31
<b>G</b>	<b>Proofs</b>	<b>33</b>
G.1	Proof of Theorem 8 . . . . .	33

G.2	Proof of Theorem 12 . . . . .	33
G.3	Proof of Theorem 2 . . . . .	34
G.4	Proof of Theorem 3 . . . . .	36
G.5	Proof of Theorem 4 . . . . .	38
<b>H</b>	<b>Improvements Over The Original Methods</b>	<b>43</b>
H.1	Importance sampling for DCGD+ . . . . .	43
H.2	Importance sampling for DIANA+ . . . . .	44
H.3	Independent sampling for ADIANA+ . . . . .	45
<b>I</b>	<b>Variance Reduction: ISEGA+</b>	<b>47</b>
<b>J</b>	<b>Variance Reduction with Bi-directional Compression: DIANA++</b>	<b>50</b>

## A Conclusions, Extensions and Future Work

In this paper we have proposed a novel gradient sparsification technique for distributed optimization and demonstrated that it allows one to properly exploit the smoothness structure of the local objective. We have shown that the proposed matrix-smoothness-aware sparsification can be coupled with both the variance reduction and acceleration, providing further speedup in terms of the convergence rate and the total bits transmitted from workers to server. Next, we list possible extensions of our work that we believe can or should be done in the future:

- **Subsampling the local objective.** While DCGD+, DIANA+ and ADIANA+ all require an access to the full local gradient from each machine at every iteration, we believe this requirement can be easily dropped. In particular, the local objective can be further subsampled and extra variance reduction can be employed on top of these methods, similarly to as done for ISAEGA [Hanzely and Richtárik, 2019b].
- **Greedy sparsification.** Notice that the sparsified local gradient can be seen as a randomized coordinate descent estimator of a given machine. However, greedy coordinate descent was shown to outperform randomized coordinate descent in certain scenarios [Nutini et al., 2017]. Therefore, one might pose a question whether a greedy sparsification might work for distributed optimization.
- **Bi-directional sparsification.** As we also mention in Section B, one drawback of our approach<sup>7</sup> is that only worker→server communication is sparse. It would be very interesting to develop a bi-directional sparsification capable of properly exploiting the smoothness matrices. For this matter, in Section J we develop and analyze DIANA++ method employing bi-directional matrix-smoothness-aware sparsification and twofold variance reduction.
- **Weakly convex and non-convex cases.** While we state our theory for the strongly convex case (i.e., Assumption 2), it can be rather easily extended to weakly convex case (i.e.,  $\mu = 0$ ). However, obtaining an efficient smoothness matrix aware sparsification for non-convex optimization remains an open problem.

---

<sup>7</sup>In fact, this is a drawback of the vast majority of compression methods from the literature. A notable exception is DoubleSqueeze [Tang et al., 2019] which compresses the server→worker communication too.

## B Limitations

Next, we discuss main limitations of our approach.

- The server is required to store matrices  $\mathbf{L}_i^{1/2}$  for all machines  $i \in [n]$  and multiply them by sparse updates  $\mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^k)$  in each iteration. Therefore, our method is not expected to be practical when  $d$  is large and matrices  $\mathbf{L}_i$  are not of a special structure so that they are cheap to store and so that  $\mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^k)$  can be evaluated cheaply.<sup>8</sup> On the other hand, our strategy is still practical when i)  $d$  is small or ii)  $\mathbf{L}_i$  is of a special structure such as low rank or diagonal. In particular, diagonal  $\mathbf{L}_i$  requires only  $\mathcal{O}(\tau)$  extra computation per each node (which is negligible), while attaining a rate which is never worse compared to the naive sparsification.
- Except DIANA++ method presented in Section J, we sparsify only the communication from the workers to server. Sparsifying workers→server communication only is very common in the area of distributed optimization as the workers→server communication is significantly more expensive compared to the server→workers communication. Such a phenomenon can be assigned to the fact that the server is broadcasting the same vector to all workers, and thus the server→workers communication can be implemented more efficiently.

**Remark 6.** *The overhead that comes from the computation of  $\mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^k)$  is not an issue in general. Given that  $\mathbf{L}_i$  is of rank  $r$ , one requires  $\mathcal{O}(d^2 r)$  flops to precompute SVD of  $\mathbf{L}_i$ . Given that SVD of  $\mathbf{L}_i$  is known, the evaluation of  $\mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^k)$  takes only  $\mathcal{O}(d^2 r)$  flops. While the cost of computing  $\nabla f_i(x^k)$  varies depending on the application, we can expect it to takes at least  $\Omega(d^2 r)$  flops for the application of generalized linear models (i.e., logistic regression). Next, we shall mention that evaluating  $\mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^k)$  comes at  $\mathcal{O}(d)$  cost when  $\mathbf{L}_i$  is diagonal.*

---

<sup>8</sup>For example, if  $\mathbf{L}_i$  is of rank  $r$ , for all  $i$ , we require extra  $\mathcal{O}(ndr)$  storage and  $\mathcal{O}(ndr)$  flops at the server at each iteration.

## C Extra Experiments and Experimental Setup

As mentioned, all experiments of this work are performed on logistic regression objective with LibSVM data. In particular, the objective is given as

$$f_i(x) := \frac{1}{m_i} \sum_{j=1}^{m_i} \log(1 + \exp((\mathbf{A}_{im})_{j,:}x \cdot (b_{im})_j)) + \frac{\mu}{2} \|x\|^2,$$

where  $\mathbf{A}_{im} \in \mathbb{R}^{d_{im} \times d}$  is the data matrix with corresponding labels  $b_{im} \in \mathbb{R}^{d_{im}}$ . In our case, we did split the randomly reshuffled datasets into equal chunks among workers in each case so that  $m_i = m_j$  for all  $i, j \leq n$ . The data matrix  $\mathbf{A}$  was normalized so that each datapoint has a norm equal to  $\frac{1}{2}$ . Lastly, we have chosen  $\mu = 10^{-3}$  for all experiments.

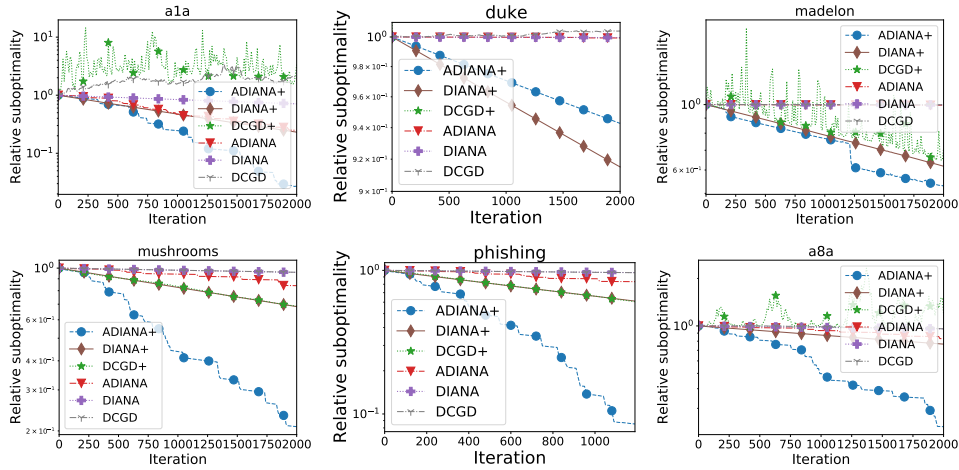
For each of the datasets, we have selected a specific number of workers given by Table 3. Each of the method was run with theory supported parameters with an exception of the ADIANA+, where we have omitted several constant factors for the sake of practicality.

**Table 3: Datasets.**

Dataset	# datapoints	$d$	$n$	$m_i$ (Fig 3, 4)	$m_i$ (remaining figures)
a1a	1 605	123	107	1	15
mushrooms	8 124	112	12	1	677
phishing	11 055	68	11	1	1 005
madelon	2 000	500	4	1	500
duke	44	7 129	4	1	11
a8a	22 696	123	8	1	2837

### C.1 Proposed and usual sparsification techniques for the 3 distributed methods.

In the experiment from the main body (Figure 1) we used the largest possible number of workers for each dataset. For the sake of robustness, we present an equivalent experiment to the one from Figure 1, but we use a moderate number of the workers given by Table 3. Figure 3 presents the result. We set the  $x$ -axis as the iteration this time. We do so to properly see the results as in some cases, the initial communication for our methods was larger than the communication during the actual algorithm run in the reported time frame.



**Figure 3:** Same as Figure 1, but smaller number of workers (see Table 3).

## C.2 Variance reduction with new sparsification and importance sampling

Here present an equivalent experiment to the one from Figure 2, but we use a moderate number of the workers given by Table 3.

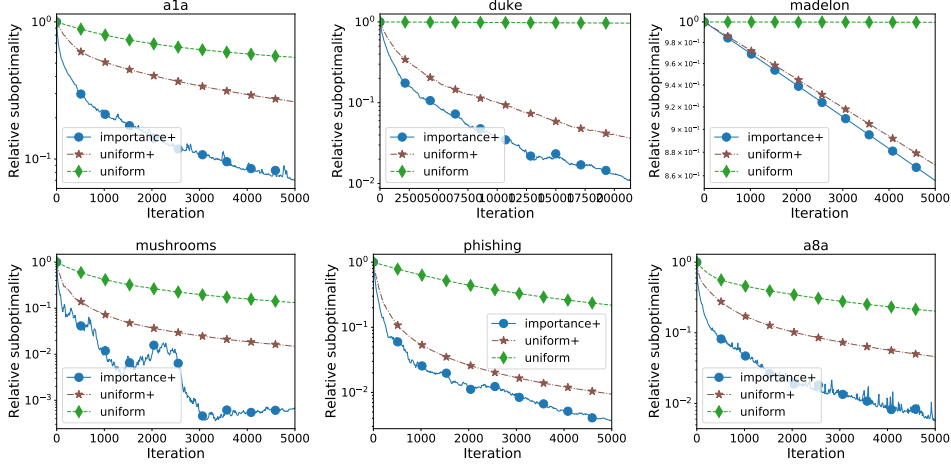


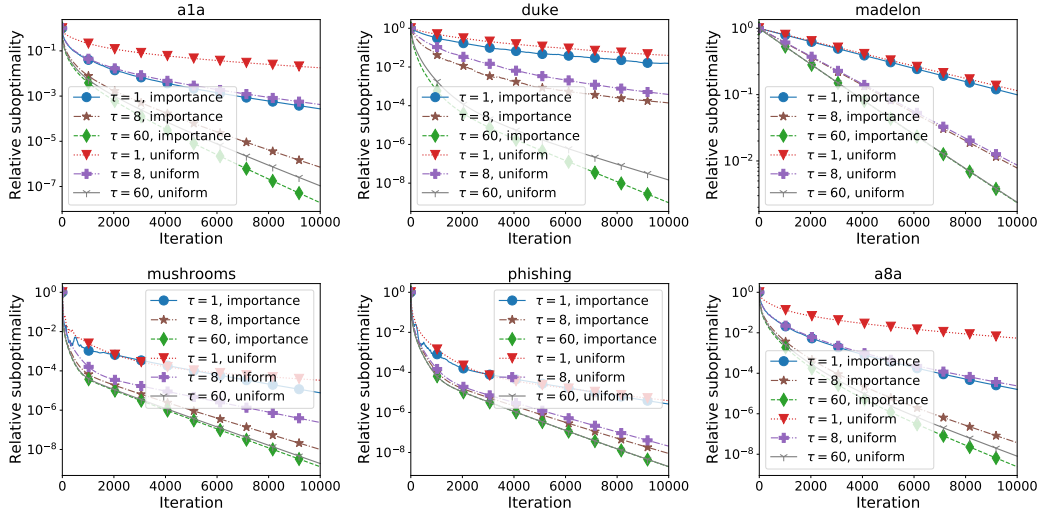
Figure 4: Same as Figure 2, but smaller number of workers (see Table 3).

## C.3 The effect of sparsification level $\tau$ on the convergence rate

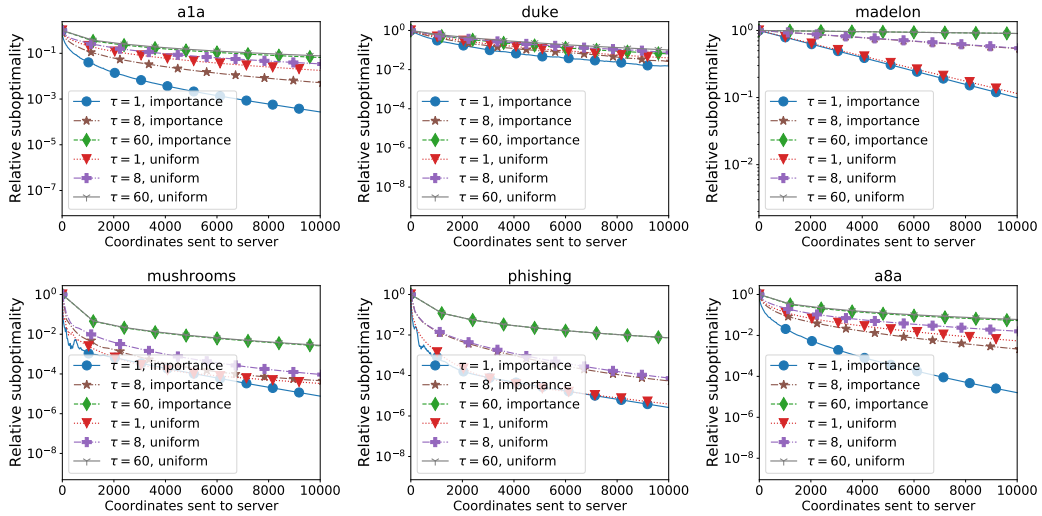
In this experiment, we study the effect of sparsification level  $\tau$  on the convergence rate. Informally speaking, our theory suggests that the sparsification does not hurt the convergence rate unless  $\tau$  is smaller than some constant. The value of such constant depends on various factors such as the type of sampling and the specific smoothness structure of the objective.

To contrast this with known results, Mishchenko et al. [2020] show that the sparsification does not hurt ISEGA significantly (a method with sparsification unaware of smoothness matrix) as soon as  $\tau n \geq d$ . Addmittedly, Mishchenko et al. [2020] assume identical smoothness constants for both  $f$  and  $f_i$ , so such a conclusion is slightly imprecise. In our case, ignoring the  $\tilde{\omega}_{\max}$  factor, the rate is dominated by the sparsification factors only if  $L = \mathcal{O}\left(\frac{\tilde{L}_{\max}}{n}\right)$ .

The results are presented in Figure 5 (Iteration vs Residual) and Figure 6 (Communication vs Residual). As expected, we see that the sparsification only hurts the iteration complexity when  $\tau$  is below certain treshold which is smaller for the uniform sampling compared to the importance sampling. Consequently, DIANA+ is capable of significantly reducing the worker->server communication at no cost in terms of the total iteration complexity.



**Figure 5:** Effect of  $\tau$  on the convergence speed of DIANA+ (Algorithm 2).



**Figure 6:** Same as Figure 5, but  $x$ -axis corresponds to the coordinates sent to the server instead of the iteration.

## D Table of Frequently Used Notation

**Table 4:** Notation used throughout the paper

Symbol	Description	Reference
$d$	dimension of the model $x \in \mathbb{R}^d$	(20)
$\mu$	strong convexity parameter of $f$	Asm. 2
$\mathbf{L}$	smoothness matrix of $f$	Asm. 1
$\mathbf{L}_{ij}$	the element at $i$ th row and $j$ th column of $\mathbf{L}$	-
$\mathbf{L}_i$	smoothness matrix of $f_i$	Asm. 1
$L_i$	smoothness constant of $f_i(x)$ , i.e., $L_i = \lambda_{\max}(\mathbf{L}_i)$	-
$L$	smoothness constant of $f$ , i.e., $L = \lambda_{\max}(\mathbf{L})$	-
$S$	random sampling (subset) of coordinates $[d] := \{1, 2, \dots, d\}$	-
$p_{jl}, p_j$	$p_{jl} := \text{Prob}(\{j, l\} \subseteq S)$ , $p_j := p_{jj}$	-
$\mathbf{P}$	the probability matrix $(p_{jl})_{j,l=1}^d$ associated with random sampling $S$	(8)
$v_i$	ESO parameters associated with $f$ and $S$ jointly	-
$\mathbf{C}$	diagonal sketch matrix with $i$ th random variable $c_i = 1/p_i$ if $i \in S$ and 0 otherwise	(6)
$\omega$	variance of general compression operator $\mathcal{C}$ , i.e. $\mathbb{E}[\ \mathcal{C}(x) - x\ ^2] \leq \omega\ x\ ^2$ , $\forall x \in \mathbb{R}^d$	-
$\bar{\mathbf{C}}, \bar{\mathbf{C}}_i^k$	$\bar{\mathbf{C}} := \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2}$ , $\bar{\mathbf{C}}_i^k = \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2}$	-
$\mathbf{I}, \mathbf{E}$	the identity matrix and the matrix with all entries equal to 1	-
$\bar{\mathbf{P}}, \tilde{\mathbf{P}}$	$\bar{\mathbf{P}} = \text{Diag}(1/p) \mathbf{P} \text{Diag}(1/p)$ with entries $\bar{p}_{ij} = \frac{p_{ij}}{p_i p_j}$ , and $\tilde{\mathbf{P}} = \bar{\mathbf{P}} - \mathbf{E}$	(8)
$\bar{\mathcal{L}}, \tilde{\mathcal{L}}$	expected smoothness constants $\bar{\mathcal{L}} = \lambda_{\max}(\bar{\mathbf{P}} \circ \mathbf{L})$ , $\tilde{\mathcal{L}} = \lambda_{\max}(\tilde{\mathbf{P}} \circ \mathbf{L})$	-
$n$	number of parallel machines in distributed setting	(1)
$\mathbf{C}_i, \mathbf{P}_i, \bar{\mathbf{P}}_i, \tilde{\mathbf{P}}_i$	diagonal sketch matrix and probability matrices for $i$ th worker	(6), (8)
$p_{i;j}, \bar{p}_{i;j}, \tilde{p}_{i;j}$	$j$ -th diagonal element of $\mathbf{P}_i, \bar{\mathbf{P}}_i, \tilde{\mathbf{P}}_i$	-
$\omega_i$	variance of compression operator induced by $\mathbf{C}_i$ , i.e. $\omega_i = \max_{1 \leq j \leq d} \frac{1}{p_{i;j}} - 1$	-
$\omega_{\max}$	$\max_{1 \leq i \leq n} \omega_i = \max_{1 \leq i \leq n} \max_{1 \leq j \leq d} \frac{1}{p_{i;j}} - 1$	(11)
$\bar{\mathcal{L}}_i, \tilde{\mathcal{L}}_i$	expected smoothness constants, $\bar{\mathcal{L}}_i = \lambda_{\max}(\bar{\mathbf{P}}_i \circ \mathbf{L}_i)$ , $\tilde{\mathcal{L}}_i = \lambda_{\max}(\tilde{\mathbf{P}}_i \circ \mathbf{L}_i)$	-
$\bar{\mathcal{L}}_{\max}, \tilde{\mathcal{L}}_{\max}$	$\bar{\mathcal{L}}_{\max} = \max_{1 \leq i \leq n} \lambda_{\max}(\bar{\mathbf{P}}_i \circ \mathbf{L}_i)$ , $\tilde{\mathcal{L}}_{\max} = \max_{1 \leq i \leq n} \lambda_{\max}(\tilde{\mathbf{P}}_i \circ \mathbf{L}_i)$	(9)
$\nu, \nu_s$	Parameters describing matrices $\mathbf{L}_i$ , $\nu := \frac{\sum_{i=1}^n L_i}{\max_{i \in [n]} L_i}$ , $\nu_s := \max_{i \in [n]} \frac{\sum_{j=1}^d L_{i;j}^{1/s}}{\max_{j \in [d]} L_{i;j}^{1/s}}$	(13)



## E Theory in the Single Node Case: RCD as Sketched Gradient Descent (SkGD)

In single node setup, matrix smoothness assumption and arbitrary samplings have been considered mainly in the context of coordinate descent methods. For example, randomized sampling  $S = \{j\}, j \in [d]$  with arbitrary probabilities  $p_j \in (0, 1]$  reduces to standard *Randomized Coordinate Descent (RCD)* algorithms [Nesterov, 2012, Richtárik and Takáč, 2014]. Parallel and mini-batch variants arise when the sampling  $S$  contains more than one coordinate [Bradley et al., 2011, Richtárik and Takáč, 2016b]. The first coordinate descent method analyzed with arbitrary sampling and under  $\mathbf{L}$ -smoothness assumption is the 'NSync algorithm [Richtárik and Takáč, 2016a, Qu and Richtárik, 2016a,b] considered for strongly convex losses. In the same general setup, Hanzely and Richtárik [2019a] developed and analyzed *Accelerated Coordinate Descent*. Recently, Hanzely et al. [2018] developed a variance reduced coordinate descent algorithm, *SEGA (SkEtched GrAdient)*, which uses general sketch matrices and handles non-separable proximal terms in contrast to traditional coordinate descent methods. This idea of gradient sketching then extended to *Generalized Jacobian Sketching (GJS)* algorithm providing a unified theory for first-order methods with variance reduced [Hanzely and Richtárik, 2019b].

Consider the unconstrained optimization problem

$$\min_{x \in \mathbb{R}^d} f(x), \quad (20)$$

with very large dimension  $d$  and assume that function  $f$  is  $\mathbf{L}$ -smooth. In this setting, the state-of-art methods are *Randomized Coordinate Descent (RCD)* type methods where in each iteration only a few coordinates get updated. Here we present new theories for RCD with arbitrary sampling paradigm, which are new and follow the idea of sketches. We will view RCD as a special case of *Compressed Gradient Descent (CGD)* with sketches (6).

### E.1 'NSync

First, we recall the first coordinate descent type algorithm, 'NSync [Richtárik and Takáč, 2016a], using arbitrary sampling. Let  $S \subseteq [d]$  be an arbitrary (proper) sampling<sup>9</sup> of coordinates such that  $p_j := \text{Prob}(j \in S) > 0, j = 1, 2, \dots, d$ . For a vector  $h \in \mathbb{R}^d$ , let  $h_S \in \mathbb{R}^d$  be the vector coinciding with  $h$  at coordinates  $j \in S$  and zeros everywhere else. Denote by  $\circ$  the Hadamard (i.e. element-wise) product. Given an arbitrary sampling  $S$  and smoothness matrix  $\mathbf{L}$ , let  $v = (v_1, v_2, \dots, v_d)$  be positive constants satisfying the *Expected Separable Overapproximation (ESO)* inequality

$$\mathbf{P} \circ \mathbf{L} \preceq \text{Diag}(p \circ v), \quad (21)$$

where  $\mathbf{P}$  is the probability matrix associated with sampling  $S$  having entries  $p_{jl} := \text{Prob}(\{j, l\} \subseteq S)$ ,  $p_j = p_{jj}$ . Analogous to (8), let  $\tilde{\mathbf{P}} = \bar{\mathbf{P}} - \mathbf{E}$ .

---

#### Algorithm 4 'NSync [Richtárik and Takáč, 2016a]

---

- 1: **Input:** Initial point  $x^0 \in \mathbb{R}^d$ , random sampling  $S$ , step size parameters  $v$ , current point  $x^k$
  - 2: Sample random set of coordinates  $S_k \sim S$
  - 3: Update selected coordinates  $x^{k+1} = x^k - \frac{1}{v} \circ \nabla f(x)_{S_k}$
- 

**Theorem 7** ('NSync, [Richtárik and Takáč, 2016a]). *Let Assumptions 1, 2 hold and  $v \sim \text{ESO}(f, S)$  be the vector of ESO parameters associated with function  $f$  and sampling  $S$ . Then the iterates  $\{x^k\}$  of 'NSync converge as follows*

$$\mathbb{E}[f(x^k)] - f(x^*) \leq \left(1 - \min_{1 \leq j \leq d} \frac{p_j \mu}{v_j}\right)^k \Delta_f,$$

where  $\Delta_f = f(x^0) - f(x^*)$ .

---

<sup>9</sup>only proper samplings are considered in this work

Thus, ‘Nsync gives an iteration complexity

$$\max_{1 \leq j \leq d} \frac{v_j}{p_j \mu} \log \frac{\Delta_f}{\varepsilon}. \quad (22)$$

In case of serial sampling, namely  $|S| = 1$  a.s., we have  $\mathbf{P} = \mathbf{Diag}(p_1, p_2, \dots, p_d)$ . Hence ESO holds with  $v_j = \mathbf{L}_{jj}$  and iteration complexity becomes  $\max_j \frac{\mathbf{L}_{jj}}{p_j \mu} \log \frac{\Delta_f}{\varepsilon}$ . This leads to the optimal probabilities  $p_j = \frac{\mathbf{L}_{jj}}{\sum_l \mathbf{L}_{ll}}$  yielding iteration complexity  $\frac{\sum_l \mathbf{L}_{ll}}{\mu} \log \frac{\Delta_f}{\varepsilon}$ .

## E.2 Sketched Gradient Descent (SkGD)

Let us view RCD methods as a special case of *Compressed Gradient Descent (CGD)* with linear and diagonal sketch  $\mathbf{C}$  defined in (6) and consider random sparsification operator  $\mathcal{C}$  induced by random diagonal sketch  $\mathbf{C}$ , namely  $\mathcal{C}(x) = \mathbf{C}x$ ,  $x \in \mathbb{R}^d$ . Clearly,  $\mathcal{C}$  is an unbiased compression (i.e.  $\mathbb{E}[\mathcal{C}(x)] = x$ ) with variance  $\omega = \max_{1 \leq j \leq d} \frac{1}{p_j} - 1$ :

$$\mathbb{E}[\|\mathbf{C}x - x\|_2^2] = x^\top \mathbb{E}[\mathbf{C}^2 - \mathbf{I}]x \leq \omega \|x\|_2^2. \quad (23)$$

---

### Algorithm 5 SKGD

---

- 1: **Input:** Initial point  $x^0 \in \mathbb{R}^d$ , diagonal sketch  $\mathbf{C}$ , step size  $\gamma$ , current point  $x^k$
  - 2:  $x^{k+1} = x^k - \gamma \mathbf{C} \nabla f(x^k)$
- 

**Theorem 8** (see G.1). *Let Assumptions 1, 2 hold and  $S$  be any proper sampling with probability matrix  $\mathbf{P}$ . Then, for the step-size  $0 < \gamma \leq \lambda_{\max}^{-1}(\bar{\mathbf{P}} \circ \mathbf{L})$ , the iterates  $\{x^k\}$  of Algorithm 5 converge as follows*

$$\mathbb{E}[f(x^k)] - f(x^*) \leq (1 - \gamma\mu)^k \Delta_f.$$

The following lemma shows that, both ‘Nsync and SkGD provide the same theoretical guarantees.

**Lemma 9.**

$$\min_{v: \mathbf{P} \circ \mathbf{L} \leq \mathbf{Diag}(v \circ p)} \max_{1 \leq j \leq d} \frac{v_j}{p_j} = \lambda_{\max}(\bar{\mathbf{P}} \circ \mathbf{L}).$$

*Proof.* If parameters  $v$  satisfy ESO inequality (21), then parameters defined by

$$v'_i := p_i \max_j \frac{v_j}{p_j} \geq v_i, \quad 1 \leq i \leq d$$

also satisfy ESO inequality and give the same iteration complexity as

$$\lambda := \max_i \frac{v_i}{p_i} = \max_i \frac{v'_i}{p_i}.$$

In particular, this implies that instead of searching for  $d$  parameters  $v_1, \dots, v_d$  satisfying ESO inequality  $\mathbf{P} \circ \mathbf{L} \leq \mathbf{Diag}(v \circ p)$  it suffices to find one scalar  $\lambda > 0$  such that  $\mathbf{P} \circ \mathbf{L} \leq \mathbf{Diag}(\lambda p \circ p)$  and set  $v_i = \lambda p_i$  for all  $i \in [d]$ . The optimal (smallest) value of the scaling factor is

$$\lambda = \lambda_{\max}(\mathbf{Diag}(1/p)(\mathbf{P} \circ \mathbf{L})\mathbf{Diag}(1/p)) = \lambda_{\max}((\mathbf{Diag}(1/p)\mathbf{P}\mathbf{Diag}(1/p)) \circ \mathbf{L}) = \lambda_{\max}(\bar{\mathbf{P}} \circ \mathbf{L}).$$

Notice that with the choice of  $v = \lambda p$ , iteration complexities as well as the update rules of both methods coincide.  $\square$

One difference between these two methods is that, the update direction  $\frac{1}{v} \circ \nabla f(x)_S$  of ‘Nsync is biased in general as opposed to unbiased direction  $\frac{1}{p} \circ \nabla f(x)_S$  of SkGD.

Note that the rate and the analysis of Theorem 8 is with respect to functional values (i.e.  $f(x^k) - f^*$ ). Natural question is to develop an analysis based on iterates of the algorithm (i.e.  $\|x^k - x^*\|^2$ ). Below, we provide such analysis under slightly different conditions on  $f$  and with weighted distances. Formally, let, instead of  $\mathbf{L}$ -smoothness and  $\mu$ -convexity, assume

$$\mu \|x - x^*\|_{\mathbf{L}}^2 + \|\nabla f(x)\|^2 \leq 2 \langle \nabla f(x), (x - x^*) \rangle_{\mathbf{L}}. \quad (24)$$

Notice that the following is true just by combining  $\mathbf{L}$ -smoothness and  $\mu$ -convexity:

$$\mu \|x - x^*\|^2 + \|\nabla f(x)\|_{\mathbf{L}^\dagger}^2 \leq 2\langle \nabla f(x), (x - x^*) \rangle. \quad (25)$$

However, in general, inequalities (24) and (25) are not equivalent.

**Theorem 10.** *Let instead of  $\mathbf{L}$ -smoothness and  $\mu$ -convexity assume (24) holds. Then, for the step-size  $0 < \gamma \leq \lambda_{\max}^{-1}(\bar{\mathbf{P}} \circ \mathbf{L})$ , the iterates  $\{x^k\}$  of Algorithm 5 converge as follows*

$$\mathbb{E} [\|x^k - x^*\|_{\mathbf{L}}^2] \leq (1 - \gamma\mu)^k \|x^0 - x^*\|_{\mathbf{L}}^2.$$

*Proof.* Consider the improvement of the algorithm in a single iteration  $x^+ = x - \gamma \mathbf{C} \nabla f(x)$ .

$$\begin{aligned} \mathbb{E} [\|x^+ - x^*\|_{\mathbf{L}}^2] &= \mathbb{E} [\|x - x^* - \gamma \mathbf{C} \nabla f(x)\|_{\mathbf{L}}^2] \\ &= \|x - x^*\|_{\mathbf{L}}^2 - 2\gamma \langle x - x^*, \nabla f(x) \rangle_{\mathbf{L}} + \gamma^2 \mathbb{E} [\|\mathbf{C} \nabla f(x)\|_{\mathbf{L}}^2] \\ &= \|x - x^*\|_{\mathbf{L}}^2 - 2\gamma \langle x - x^*, \nabla f(x) \rangle_{\mathbf{L}} + \gamma^2 \|\nabla f(x)\|_{\mathbb{E}[\mathbf{C} \mathbf{L} \mathbf{C}]}^2 \\ &\stackrel{(36)}{=} \|x - x^*\|_{\mathbf{L}}^2 - 2\gamma \langle x - x^*, \nabla f(x) \rangle_{\mathbf{L}} + \gamma^2 \|\nabla f(x)\|_{\bar{\mathbf{P}} \circ \mathbf{L}}^2 \\ &\leq \|x - x^*\|_{\mathbf{L}}^2 - 2\gamma \langle x - x^*, \nabla f(x) \rangle_{\mathbf{L}} + \gamma^2 \lambda_{\max}(\bar{\mathbf{P}} \circ \mathbf{L}) \|\nabla f(x)\|^2 \\ &\leq \|x - x^*\|_{\mathbf{L}}^2 - 2\gamma \langle x - x^*, \nabla f(x) \rangle_{\mathbf{L}} + \gamma \|\nabla f(x)\|^2 \\ &\stackrel{(24)}{\leq} (1 - \gamma\mu) \|x - x^*\|_{\mathbf{L}}^2. \end{aligned}$$

□

### E.3 CGD+

Here we introduce a new variant of CGD with non-diagonal matrix  $\bar{\mathbf{C}} := \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2}$ , which works with any proximable regularizer  $R(x)$ . In this case the method converges to the neighborhood of the solution. Recall that the proximal operator is defined as follows:

$$\text{prox}_R(x) = \arg \min_{u \in \mathbb{R}^d} \left( R(u) + \frac{1}{2} \|u - x\|^2 \right). \quad (26)$$

Define expected smoothness constants

$$\bar{\mathcal{L}} = \lambda_{\max}(\bar{\mathbf{P}} \circ \mathbf{L}), \quad \tilde{\mathcal{L}} = \lambda_{\max}(\tilde{\mathbf{P}} \circ \mathbf{L}).$$

The following lemma reveals the relationship between these constants.

**Lemma 11.** *Let  $L = \lambda_{\max}(\mathbf{L})$ . Then  $L \leq \bar{\mathcal{L}} \leq L + \tilde{\mathcal{L}}$ .*

*Proof.* First, positive semi-definiteness of  $\mathbf{P}$  was proved in Theorem 3.1 [Qu and Richtárik, 2016]. As  $\text{Diag}(1/p)$  is positive definite, then  $\bar{\mathbf{P}}$  is positive semi-definite too. Since Hadamard product  $\circ$  preserves positive semi-definiteness, we have that  $\bar{\mathbf{P}} \circ \mathbf{L} \succeq 0$ . It follows from Lemma 17 that

$$\mathbb{E} [\mathbf{L}^{1/2} (\bar{\mathbf{C}} - \mathbf{I})^\top (\bar{\mathbf{C}} - \mathbf{I}) \mathbf{L}^{1/2}] = \mathbf{L}^{1/2} \mathbf{L}^{\dagger 1/2} (\bar{\mathbf{P}} \circ \mathbf{L}) \mathbf{L}^{\dagger 1/2} \mathbf{L}^{1/2}.$$

Hence the left hand side as well as  $\bar{\mathbf{P}} \circ \mathbf{L}$  are symmetric and positive semidefinite. In particular,  $\bar{\mathbf{P}} \circ \mathbf{L} \succeq \mathbf{L}$ . Hence  $L = \lambda_{\max}(\mathbf{L}) \leq \lambda_{\max}(\bar{\mathbf{P}} \circ \mathbf{L}) = \bar{\mathcal{L}}$ . The upper bound follows from the convexity of  $\lambda_{\max}$  as  $\tilde{\mathcal{L}} = \lambda_{\max}(\tilde{\mathbf{P}} \circ \mathbf{L}) = \lambda_{\max}(\mathbf{L} + \tilde{\mathbf{P}} \circ \mathbf{L}) \leq \lambda_{\max}(\mathbf{L}) + \lambda_{\max}(\tilde{\mathbf{P}} \circ \mathbf{L}) = L + \tilde{\mathcal{L}}$ . □

---

#### Algorithm 6 CGD+

---

- 1: **Input:** Initial point  $x^0 \in \mathbb{R}^d$ , sketch matrix  $\bar{\mathbf{C}} = \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2}$ , step size  $\gamma$ , current point  $x^k$
  - 2:  $x^{k+1} = \text{prox}_{\gamma R}(x^k - \gamma \bar{\mathbf{C}} \nabla f(x^k))$
- 

With the new sketch  $\bar{\mathbf{C}}$  in Algorithm 6 we are able to perform the analysis with respect to iterates in standard norm, under strong convexity and  $\mathbf{L}$ -smoothness, allowing any proximable regularizer.

**Table 5:** Original and proposed new methods for both single node and distributed setups.

ORIGINAL	*NSYNC	CGD	DCGD	DIANA	ADIANA
NEW	SkGD (ALG.5)	CGD+ (ALG.6)	DCGD+ (ALG.1)	DIANA+ (ALG.2)	ADIANA+ (ALG.3)
PROXIMAL	✗	✓	✓	✓	✓
DISTRIBUTED	✗	✗	✓	✓	✓
VARIANCE REDUCED	✗	✗	✗	✓	✓
ACCELERATED	✗	✗	✗	✗	✓

**Table 6:** Complexity of new methods with hidden log factors and constants.

Method	Iteration Complexity
SkGD (Algorithm 5)	$\frac{\bar{\mathcal{L}}}{\mu}$
CGD+ (Algorithm 6)	$\frac{\bar{\mathcal{L}}}{\mu} + \frac{\tilde{\mathcal{L}}}{\mu^2 \varepsilon}$
DCGD+ (Algorithm 1)	$\frac{L}{\mu} + \frac{\tilde{\mathcal{L}}_{\max}}{\mu n} + \frac{\tilde{\mathcal{L}}_{\max}}{\mu^2 n \varepsilon}$
DIANA+ (Algorithm 2)	$\omega_{\max} + \frac{L}{\mu} + \frac{\tilde{\mathcal{L}}_{\max}}{\mu n}$
ADIANA+ (Algorithm 3)	$\begin{cases} \omega_{\max} + \sqrt{\omega_{\max} \frac{\tilde{\mathcal{L}}_{\max}}{\mu n}} & \text{if } nL \leq \tilde{\mathcal{L}}_{\max} \\ \omega_{\max} + \sqrt{\frac{L}{\mu}} + \sqrt{\omega_{\max} \sqrt{\frac{\tilde{\mathcal{L}}_{\max}}{\mu n}} \sqrt{\frac{L}{\mu}}} & \text{if } nL > \tilde{\mathcal{L}}_{\max}. \end{cases}$

**Theorem 12** (see G.2). *Let Assumptions 1, 2 hold and  $S$  be a sampling with probability matrix  $\mathbf{P}$ . Then, for the step-size  $0 < \gamma \leq 1/2\bar{\mathcal{L}}$ , the iterates  $\{x^k\}$  of Algorithm 6 converge as follows*

$$\mathbb{E} [\|x^k - x^*\|^2] \leq (1 - \gamma\mu)^k \|x^0 - x^*\|^2 + \frac{2\gamma\tilde{\mathcal{L}}}{\mu} \|\nabla f(x^*)\|_{\mathbf{L}^\dagger}^2.$$

## F Lower Bounds for Sketches as Linear Compression Operators

Here we investigate general sketch matrices  $\mathbf{S}$  as a linear compression operators. The motivation of this is to understand the trade-off between communication and variance of linear compressors. The notation, used in this section only, slightly deviates from the paper but otherwise is consistent throughout the section.

Consider compression of vectors  $x \in \mathbb{R}^d$  allowing approximation error in exchange for less bits of communication. Let compression operator  $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$  be composed of some linear encoder  $E(x) = \mathbf{S}x$  with  $s \times d$  sketch matrix  $\mathbf{S}$  and an arbitrary decoder  $D: \mathbb{R}^s \rightarrow \mathbb{R}^d$ , so that  $\mathcal{C}(x) = D(\mathbf{S}x)$ . Throughout we consider the space  $\mathbb{R}^d$  equipped with an inner product together with its induced norm given by some symmetric and positive definite matrix  $\mathbf{B}$  of size  $d \times d$  as follows

$$\langle x, y \rangle_{\mathbf{B}} = x^{\top} \mathbf{B} y, \quad \|x\|_{\mathbf{B}} = \sqrt{\langle x, x \rangle_{\mathbf{B}}}, \quad x, y \in \mathbb{R}^d.$$

In general, we let matrix  $\mathbf{S}$ , number of rows  $s$  and decoder  $D$  to be random, while the matrix  $\mathbf{B}$  will be fixed throughout the analysis. Since we consider only linear encoders, we may assume  $\|x\|_{\mathbf{B}} = 1$ .

### F.1 Fixed sketches

We first analyze the case where the sketch matrix  $\mathbf{S}$  is fixed and hence the compression operator  $\mathcal{C}$  is deterministic. The analysis then we will lead us on a more usefull result for random sketches. The decoder  $D$  receiving vector  $y = \mathbf{S}x$  should be able to reconstruct  $\hat{x} = D(y)$  so to minimize the squared error

$$\alpha(\mathbf{S}) := \sup_{\|x\|_{\mathbf{B}}=1} \|\mathcal{C}(x) - x\|_{\mathbf{B}}^2 = \sup_{\|x\|_{\mathbf{B}}=1} \|D(\mathbf{S}x) - x\|_{\mathbf{B}}^2 \leq 1.$$

The following lemma shows the optimal strategy for the decoder and possible values for  $\alpha(\mathbf{S})$ .

**Lemma 13.** *For a fixed sketch  $\mathbf{S}$  the optimal reconstruction from  $y = \mathbf{S}x$  is*

$$D^*(y) = \mathbf{S}^{\dagger \mathbf{B}} y \equiv \mathbf{B}^{-1} \mathbf{S}^{\top} (\mathbf{S} \mathbf{B}^{-1} \mathbf{S}^{\top})^{\dagger} y, \quad (27)$$

where  $\cdot^{\dagger}$  indicates the Moore–Penrose inverse of a matrix. Furthermore, if  $\ker(\mathbf{S}) = \{0\}$  then  $\alpha(\mathbf{S}) = 0$  as in this case  $D^*(\mathbf{S}x) = x$  for any  $x \in \mathbb{R}^d$ . Otherwise, if  $\ker(\mathbf{S}) \neq \{0\}$ , then  $\alpha(\mathbf{S}) = 1$ .

*Proof.* Let  $\ker(\mathbf{S}) = \{z: \mathbf{S}z = 0\}$  be the kernel of  $\mathbf{S}$  and  $x^{\dagger \mathbf{B}} = \mathbf{S}^{\dagger \mathbf{B}} y$  be the minimal  $\mathbf{B}$ -norm solution to the system  $\mathbf{S}z = y$  so that the set of all solutions is  $x^{\dagger \mathbf{B}} + \ker(\mathbf{S})$ :

$$x^{\dagger \mathbf{B}} = \arg \min_{x: \mathbf{S}x=y} \|x\|_{\mathbf{B}}^2 = \mathbf{S}^{\dagger \mathbf{B}} y = \mathbf{B}^{-1/2} \left( \mathbf{S} \mathbf{B}^{-1/2} \right)^{\dagger} y,$$

Denote by

$$\hat{S}(x) := (x^{\dagger \mathbf{B}} + \ker(\mathbf{S})) \cap \{z \in \mathbb{R}^d: \|z\|_{\mathbf{B}} = 1\}$$

the intersection of the affine set of solutions and the unit sphere. Notice that initial vector  $x \in \hat{S}(x)$  as it has unit  $\mathbf{B}$ -norm and satisfies  $\mathbf{S}x = y$ . Now the cost of sending  $\mathbf{S}x$  instead of original  $x$ , is the uncertainty that the decoder has to deal with by estimating the original vector within the set  $\hat{S}$  so to minimize  $\alpha$ . We first show that  $x^{\mathbf{S}} := 2x^{\dagger \mathbf{B}} - x \in \hat{S}(x)$ , which is equivalent to

$$x^{\dagger \mathbf{B}} - x \in \ker(\mathbf{S}) \quad \text{and} \quad \|2x^{\dagger \mathbf{B}} - x\|_{\mathbf{B}}^2 = 1.$$

The first claim follows from the fact that both  $x$  and  $x^{\dagger \mathbf{B}}$  are solutions to  $\mathbf{S}z = y$ , namely  $\mathbf{S}x^{\dagger \mathbf{B}} = y = \mathbf{S}x$ . Expanding the square in the second claim we get  $\langle x^{\dagger \mathbf{B}}, x^{\dagger \mathbf{B}} - x \rangle_{\mathbf{B}} = 0$  which holds as  $x^{\dagger \mathbf{B}}$  is the minimal  $\mathbf{B}$ -norm solution. Therefore the vector  $y$  the decoder receives does not differentiate between  $x$  and  $x^{\mathbf{S}}$ . This implies that for any choice of  $\hat{x}$  of the decoder

$$\max \left( \|\hat{x} - x\|_{\mathbf{B}}^2, \|\hat{x} - x^{\mathbf{S}}\|_{\mathbf{B}}^2 \right) \geq \frac{1}{4} \left( \|\hat{x} - x\|_{\mathbf{B}} + \|\hat{x} - x^{\mathbf{S}}\|_{\mathbf{B}} \right)^2 \geq \frac{1}{4} \|x^{\mathbf{S}} - x\|_{\mathbf{B}}^2 = \|x^{\dagger \mathbf{B}} - x\|_{\mathbf{B}}^2$$

squared-error is unavoidable for the couple  $x, x^{\mathbf{S}}$  and the optimal choice is  $\hat{x} = x^{\dagger \mathbf{B}}$ . Thus, the optimal decoding strategy to  $y = \mathbf{S}x$  is  $D^*(y) = x^{\dagger \mathbf{B}}$  given in (27). Now, if  $\ker(\mathbf{S}) \neq \{0\}$  then we could pick the initial vector  $x$  from the kernel space, i.e.  $x \in \ker(\mathbf{S})$  and  $\|x\|_{\mathbf{B}} = 1$ . Then we would have  $x^{\dagger \mathbf{B}} = 0$  and hence the minimal squared-error  $\alpha(\mathbf{S}) = 1$ . On the other hand, if  $\ker(\mathbf{S}) = \{0\}$ , then  $x^{\dagger \mathbf{B}} = x$  as the system  $\mathbf{S}z = y$  has unique solution.  $\square$

To conclude for fixed sketches, notice that,  $x$  and  $x^S$  are in symmetry in this analysis. Indeed, if the initial vector was  $x^S$  as opposed to  $x$ , then  $\mathbf{S}x = \mathbf{S}x^S$ , hence  $x^{S^\dagger \mathbf{B}} = x^{\dagger \mathbf{B}}$  and  $x^{SS} = x$ . Therefore, the analysis of Lemma 13 leads to the following lower bound for any decoder  $D$  and initial vector  $x \in \mathbb{R}^d$

$$\max_{z=x, x^S} \|\mathcal{C}(z) - z\|_{\mathbf{B}}^2 \geq \|x^{\dagger \mathbf{B}} - x\|_{\mathbf{B}}^2 = 1 - \|x^{\dagger \mathbf{B}}\|_{\mathbf{B}}^2 = 1 - \|\mathbf{Z}x\|_{\mathbf{B}}^2, \quad (28)$$

where we used orthogonality  $\langle x^{\dagger \mathbf{B}}, x^{\dagger \mathbf{B}} - x \rangle_{\mathbf{B}} = 0$  and defined the random matrix  $\mathbf{Z} = \mathbf{Z}(\mathbf{S})$  via

$$\mathbf{Z} := \mathbf{S}^{\dagger \mathbf{B}} \mathbf{S} = \mathbf{B}^{-1/2} \left( \mathbf{S} \mathbf{B}^{-1/2} \right)^{\dagger} \mathbf{S} = \mathbf{B}^{-1} \mathbf{S}^{\top} \left( \mathbf{S} \mathbf{B}^{-1} \mathbf{S}^{\top} \right)^{\dagger} \mathbf{S}.$$

## F.2 Random sketches

Now we turn to the general case when sketch matrix  $\mathbf{S}$  is random and drawn from some distribution  $\mathcal{D}$ , to which both encoder and decoder have access. The number of rows  $s$  of  $\mathbf{S}$  can also be random. In this case, the decoder  $D$  upon receiving random vector  $y = \mathbf{S}x$  should estimate possibly randomized  $\hat{x} = D(y)$  so to minimize the expected square error

$$\alpha(\mathcal{D}) := \sup_{\|x\|_{\mathbf{B}}=1} \mathbb{E} [\|\mathcal{C}(x) - x\|_{\mathbf{B}}^2] \leq 1, \quad (29)$$

where  $\mathcal{C}(x) = D(\mathbf{S}x)$  is a random mapping with a source of randomness coming from the distribution  $\mathcal{D}$  and decoder  $D$ . Below we prove a lower bound for  $\alpha(\mathcal{D})$ .

**Theorem 14.** *Let  $\mathcal{D}$  be some distribution over  $s \times d$  matrices  $\mathbf{S}$  allowing variable number of rows  $s \in [d]$ . Then for any (possibly randomized) compression operator  $\mathcal{C}(x) = D(\mathbf{S}x)$  with i.i.d. samples  $\mathbf{S} \sim \mathcal{D}$  and  $x \in \mathbb{R}^d$  the following lower bound holds*

$$\alpha(\mathcal{D}) + \mathbb{E}_{\mathcal{D}} [r/d] \geq 1, \quad (30)$$

where  $r = \text{rank}(\mathbf{S})$  is the number of independent rows in  $\mathbf{S}$ .

*Proof.* Based on the lower bound (28) obtained from the deterministic case, decoder cannot avoid the error  $1 - \|\mathbf{Z}x\|_{\mathbf{B}}^2$  even in the case of knowing what sketch the encoder used. Therefore minimal expected error  $1 - \mathbb{E}_{\mathbf{S} \sim \mathcal{D}} \|\mathbf{Z}x\|_{\mathbf{B}}^2$  is unavoidable for any initial  $x$ . This leads to the following bound

$$\begin{aligned} 1 - \alpha(\mathcal{D}) &\leq \inf_{\|x\|_{\mathbf{B}}=1} \mathbb{E}_{\mathcal{D}} [\|\mathbf{Z}x\|_{\mathbf{B}}^2] \\ &= \inf_{\|x\|_{\mathbf{B}}=1} \mathbb{E}_{\mathcal{D}} [x^{\top} \mathbf{Z}^{\top} \mathbf{B} \mathbf{Z} x] \\ &\stackrel{z=\mathbf{B}^{1/2}x}{=} \inf_{\|z\|=1} \mathbb{E}_{\mathcal{D}} [z^{\top} \mathbf{B}^{-1/2} \mathbf{Z}^{\top} \mathbf{B} \mathbf{Z} \mathbf{B}^{-1/2} z] \\ &= \inf_{\|z\|=1} z^{\top} \mathbb{E}_{\mathcal{D}} [\mathbf{B}^{-1/2} \mathbf{Z}^{\top} \mathbf{B} \mathbf{Z} \mathbf{B}^{-1/2}] z \\ &= \lambda_{\min} \left( \mathbb{E}_{\mathcal{D}} [\mathbf{B}^{-1/2} \mathbf{Z}^{\top} \mathbf{B} \mathbf{Z} \mathbf{B}^{-1/2}] \right) \\ &= \lambda_{\min} \left( \mathbb{E}_{\mathcal{D}} [\mathbf{B}^{-1} \mathbf{Z}^{\top} \mathbf{B} \mathbf{Z}] \right) \\ &= \lambda_{\min} \left( \mathbb{E}_{\mathcal{D}} [\mathbf{B}^{-1} \mathbf{S}^{\top} (\mathbf{S} \mathbf{B}^{-1} \mathbf{S}^{\top})^{\dagger} \mathbf{S}] \right) \\ &= \lambda_{\min} (\mathbb{E}_{\mathcal{D}} [\mathbf{Z}]), \end{aligned}$$

where the expectation is with respect to  $\mathbf{S} \sim \mathcal{D}$ . Thus, we obtained the following lower bound:

$$\alpha(\mathcal{D}) + \lambda_{\min} (\mathbb{E}_{\mathcal{D}} [\mathbf{S}^{\dagger \mathbf{B}} \mathbf{S}]) \geq 1. \quad (31)$$

To prove the inequality (30), it is enough to establish the following upper bound for the minimal eigenvalue

$$\lambda_{\min} (\mathbb{E}_{\mathcal{D}} [\mathbf{Z}]) \leq \mathbb{E}_{\mathcal{D}} [r/d].$$

We follow the proof of Lemma 4.2 of Gower and Richtárik [2015] to prove this inequality. It can be easily checked that, using the properties of pseudo-inverse,  $\mathbf{Z} = \mathbf{S}^{\dagger \mathbf{B}} \mathbf{S}$  is an idempotent matrix for any  $\mathbf{S}$ , namely  $\mathbf{Z}^2 = \mathbf{Z}$ . This implies that all eigenvalues of  $\mathbf{Z}$  are either 0 or 1 as they must satisfy

the same relation  $\lambda^2 = \lambda$ . Trace  $\text{tr}(\mathbf{Z})$  of such matrices coincides with the number of non-zero eigenvalues, which also shows the rank:

$$\text{tr}(\mathbf{Z}) = \sum_{i=1}^d \lambda_i(\mathbf{Z}) = \#\{i \in [d]: \lambda_i(\mathbf{Z}) \neq 0\} = \text{rank}(\mathbf{Z}). \quad (32)$$

From the properties of pseudo-inverse it follows that  $\text{rank}(\mathbf{A}^\dagger \mathbf{A}) = \text{rank}(\mathbf{A}^\dagger) = \text{rank}(\mathbf{A})$  for any matrix  $\mathbf{A}$ . Hence

$$\begin{aligned} \text{rank}(\mathbf{Z}) &= \text{rank}(\mathbf{S}^\dagger \mathbf{B} \mathbf{S}) = \text{rank}\left(\mathbf{B}^{-1/2} \left(\mathbf{S} \mathbf{B}^{-1/2}\right)^\dagger \mathbf{S}\right) \\ &= \text{rank}\left(\left(\mathbf{S} \mathbf{B}^{-1/2}\right)^\dagger \mathbf{S} \mathbf{B}^{-1/2}\right) = \text{rank}\left(\mathbf{S} \mathbf{B}^{-1/2}\right) = \text{rank}(\mathbf{S}) = r. \end{aligned}$$

Combining with (32) we get  $\text{tr}(\mathbf{Z}) = r$ . The purpose of expressing the rank as a trace is that in contrast to rank, trace and expectation operators are commutative, which basically follows from the linearity of the expectation:

$$\text{tr}(\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]) = \mathbb{E}_{\mathcal{D}}[\text{tr}(\mathbf{Z})]. \quad (33)$$

Using (32), (33) and  $\text{tr}(\mathbf{Z}) = r$ , we conclude

$$\lambda_{\min}(\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]) \leq \frac{1}{d} \sum_{i=1}^d \lambda_i(\mathbb{E}_{\mathcal{D}}[\mathbf{Z}]) = \frac{\text{tr}(\mathbb{E}_{\mathcal{D}}[\mathbf{Z}])}{d} = \frac{\mathbb{E}_{\mathcal{D}}[\text{tr}(\mathbf{Z})]}{d} = \frac{\mathbb{E}_{\mathcal{D}}[r]}{d},$$

which completes the proof.  $\square$

### F.3 Optimal sketches

With the knowledge of this new lower bound, here we construct a distribution  $\mathcal{D}$  of sketches that will achieve equality in (30). Let  $\mathbf{B} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$  be the eigendecomposition of the symmetric matrix  $\mathbf{B}$ , where  $\mathbf{\Lambda}$  is diagonal with eigenvalues and  $\mathbf{Q}$  is orthogonal with eigenvectors as columns. Let  $\mathbf{C}$  be the diagonal sketch of size  $d \times d$  corresponding to random sparsification with probabilities  $p = (p_i)_{i=1}^d$ , namely

$$\mathbf{C} = \text{Diag}(c), \quad c_i = \begin{cases} 1 & \text{with prob. } p_i, \\ 0 & \text{with prob. } 1 - p_i. \end{cases}$$

Define a distribution  $\mathcal{D} = \mathcal{D}_p$  of sketches as  $\mathbf{S} = \mathbf{C} \mathbf{Q}^\top$  and notice that

$$\mathbb{E}_{\mathcal{D}}[\text{rank}(\mathbf{S})] = \mathbb{E}_{\mathcal{D}}[\text{rank}(\mathbf{C})] = \mathbb{E}_{\mathcal{D}}[\#\{i \in [d]: c_i = 1\}] = \mathbb{E}_{\mathcal{D}}\left[\sum_{i=1}^d c_i\right] = \sum_{i=1}^d \mathbb{E}_{\mathcal{D}}[c_i] = \sum_{i=1}^d p_i.$$

Therefore,  $\mathbb{E}_{\mathcal{D}}[r/d] = \frac{1}{d} \sum p_i$ . With decoder  $D(x) = \mathbf{Q}x$  we get a compression operator  $\mathcal{C}(x) = \mathbf{Q} \mathbf{S} x$ . Next, we compute  $\alpha(\mathcal{D})$  as follows

$$\begin{aligned} \alpha(\mathcal{D}) &= \sup_{\|x\|_{\mathbf{B}}=1} \mathbb{E}[\|\mathcal{C}(x) - x\|_{\mathbf{B}}^2] \\ &= \sup_{\|x\|_{\mathbf{B}}=1} \mathbb{E}[\|\mathbf{Q} \mathbf{S} x - x\|_{\mathbf{B}}^2] \\ &= \sup_{x^\top \mathbf{B} x = 1} \mathbb{E}[x^\top (\mathbf{I} - \mathbf{Q} \mathbf{S})^\top \mathbf{B} (\mathbf{I} - \mathbf{Q} \mathbf{S}) x] \\ &= \sup_{x^\top \mathbf{Q} \mathbf{C} \mathbf{Q}^\top x = 1} x^\top \mathbb{E}[(\mathbf{I} - \mathbf{Q} \mathbf{C} \mathbf{Q}^\top) \mathbf{B} (\mathbf{I} - \mathbf{Q} \mathbf{C} \mathbf{Q}^\top)] x \\ &= \sup_{(\mathbf{Q}^\top x)^\top \mathbf{\Lambda} (\mathbf{Q}^\top x)} (\mathbf{Q}^\top x)^\top \mathbb{E}[(\mathbf{I} - \mathbf{C}) \mathbf{Q}^\top \mathbf{B} \mathbf{Q} (\mathbf{I} - \mathbf{C})] (\mathbf{Q}^\top x) \\ &\stackrel{y=\mathbf{Q}^\top x}{=} \sup_{y^\top \mathbf{\Lambda} y = 1} y^\top \mathbb{E}[(\mathbf{I} - \mathbf{C}) \mathbf{\Lambda} (\mathbf{I} - \mathbf{C})] y \\ &= \sup_{y^\top \mathbf{\Lambda} y = 1} (\mathbf{\Lambda}^{1/2} y)^\top \mathbb{E}[(\mathbf{I} - \mathbf{C})^2] (\mathbf{\Lambda}^{1/2} y) \\ &\stackrel{z=\mathbf{\Lambda}^{1/2} y}{=} \sup_{\|z\|=1} z^\top \cdot \text{Diag}(1 - p) \cdot z \\ &= \max_{1 \leq i \leq d} (1 - p_i) = 1 - \min_{1 \leq i \leq d} p_i. \end{aligned}$$

Hence

$$1 \leq \alpha(\mathcal{D}) + \mathbb{E}_{\mathcal{D}}[r/d] = 1 - \min_{1 \leq i \leq d} p_i + \frac{1}{d} \sum_{i=1}^d p_i,$$

and equality occurs if and only if all probabilities  $p_i$  are equal to some  $q \in [0, 1]$ . Thus, the optimal sketches are obtained by rotating the coordinate basis to the basis of eigenvectors of  $\mathbf{Q}$  (i.e.  $x \rightarrow \mathbf{Q}^\top x$ ), and then randomly sparsify coordinates with diagonal sketch matrix  $\mathbf{C}$  (i.e.  $\mathbf{Q}^\top x \rightarrow \mathbf{CQ}^\top x = \mathbf{S}x$ ). We summarize this result in the following theorem.

**Theorem 15.** *Let  $\mathbf{B} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$  be the eigendecomposition of  $\mathbf{B}$  of induced norm,  $q \in [0, 1]$  and  $\mathbf{C}$  be random diagonal sketch corresponding to the random  $q$ -sparsifier. Then sketches  $\mathbf{S} = \mathbf{CQ}^\top$  are optimal with respect to variance against rank trade-off (30) with squared error  $\alpha = 1 - q$  and expected rank  $\mathbb{E}[r] = qd$ .*

#### F.4 Random sketches with linear constraints

In this part we extend the theory of compressing vectors  $x \in \mathbb{R}^d$  with an additional linear constraint  $x \in \text{Range}(\mathbf{A})$  for some  $d \times d'$  matrix  $\mathbf{A}$ . Such scenarios occur when to-be-compressed vectors are the gradients of  $f(w) = \phi(\mathbf{A}^\top w)$ , for which  $\nabla f(w) = \mathbf{A} \nabla \phi(\mathbf{A}^\top w) \in \text{Range}(\mathbf{A})$ . Without loss of generality, we may assume that  $\mathbf{A}$  is of full column rank and consequently  $d' = \dim \text{Range}(\mathbf{A}) = \text{rank}(\mathbf{A})$ . The constraint  $x \in \text{Range}(\mathbf{A})$  then can be equivalently written as  $x = \mathbf{A}x'$  for some  $x' \in \mathbb{R}^{d'}$ . The induced inner product and norm on  $\text{Range}(\mathbf{A})$  is then given by the matrix  $\mathbf{A}^\top \mathbf{B} \mathbf{A}$  as

$$\langle x, y \rangle_{\mathbf{B}} = \langle \mathbf{A}x', \mathbf{A}y' \rangle_{\mathbf{B}} = \langle x', y' \rangle_{\mathbf{A}^\top \mathbf{B} \mathbf{A}}, \quad x = \mathbf{A}x', y = \mathbf{A}y'.$$

Notice that, since  $\mathbf{S}x = \mathbf{S}\mathbf{A}x'$ , communication of  $x \in \mathbb{R}^d$  with sketches  $\mathbf{S}$  reduces to communication of  $x' \in \mathbb{R}^{d'}$  with sketches  $\mathbf{S}\mathbf{A}$ . Thus, the additional constraint  $x \in \text{Range}(\mathbf{A}) \subset \mathbb{R}^d$  reduces the problem to lower  $d'$ -dimension with sketches  $\mathbf{S}\mathbf{A}$ ,  $\mathbf{S} \sim \mathcal{D}$  and norm induced by  $\mathbf{A}^\top \mathbf{B} \mathbf{A}$ .

#### F.5 Variance against communication trade-off

The obtained lower bound (30) can be easily translated in terms of the number of bits. Assuming each float takes 32 bits to encode and there is no redundant row in  $\mathbf{S}$  (i.e.  $s = r$ ), then  $\mathbf{S}x \in \mathbb{R}^r$  can be communicated with up to  $b = 32r$  bits. Therefore, the lower bound (30) can be written as

$$\alpha + \frac{\mathbb{E}[b]}{32d} \geq 1, \quad (34)$$

which (ignoring the expectation) is exponentially stronger than the lower bound  $\alpha \cdot 4^{b/d} \geq 1$  obtained for general compressors in [Safaryan et al., 2020]. We visualize the comparison of these two lower bounds in Figure 7. Furthermore, denote by  $\beta := \mathbb{E}[b]/32d$  the expected communication reduction factor and recall that  $\alpha$  is the portion of the expected lost of information. With this notation the above lower bound (34) turns to the following simple inequality

$$\alpha + \beta \geq 1,$$

showing the trade-off between information lost and communication reduction for linear compressors; namely more reduction in communication leads to bigger information loss and vice versa. In one extreme, when all  $32d$  bits are sent, no reduction in communication is made ( $\beta = 1$ ) and no information is lost ( $\alpha = 0$ ). In other extreme, when no bits gets transferred ( $\beta = 0$ ) we loose all information ( $\alpha = 1$ ).

To conclude this section, let us investigate the optimality of random  $q$ -sparsifier with respect to the lower bound (34). Recall that random  $q$ -sparsifier is optimal with respect to (30). Let  $q \in (0, 1)$ , and  $k$  be the (random) number of non-zero entries of sparsified vector. Clearly,  $\mathbb{E}[k] = qd$  and to encode any  $k$ -sparse vector one needs  $b = 32k + \log_2 \binom{d}{k}$  bits. As we know from Theorem 15, the squared error  $\alpha = 1 - q$ . Therefore

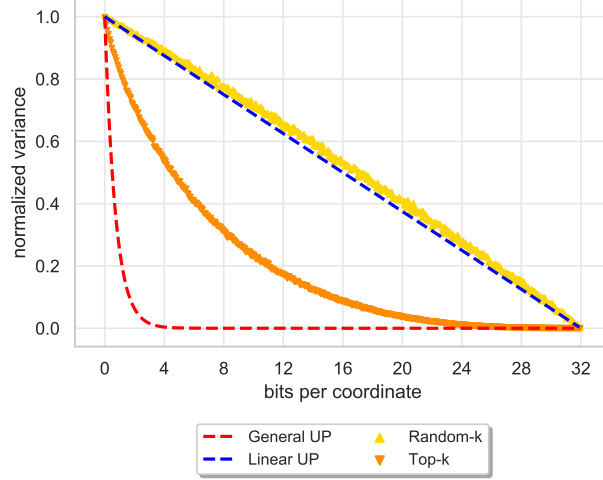
$$\alpha + \beta = 1 - q + \frac{1}{32d} \mathbb{E} \left[ 32k + \log_2 \binom{d}{k} \right] = 1 + \frac{1}{32d} \mathbb{E} \left[ \log_2 \binom{d}{k} \right] \leq 1 + \frac{1}{32} \mathbb{E} \left[ H_2 \left( \frac{k}{d} \right) \right] \leq 1 + \frac{H_2(q)}{32}.$$

The first inequality follows from the following estimate (only upper bound) for binomial coefficients

$$\frac{2^{dH_2(\tau)}}{\sqrt{8d\tau(1-\tau)}} \leq \binom{d}{\tau d} \leq \frac{2^{dH_2(\tau)}}{\sqrt{2\pi d\tau(1-\tau)}}, \quad 0 < \tau < 1,$$



where  $H_2(\tau) = -\tau \log_2 \tau - (1 - \tau) \log_2 (1 - \tau)$  is the binary entropy function in bits. The second inequality follows from concavity  $H_2$  function and the Jensen's inequality. Because of the symmetry around  $\tau = 1/2$  (namely  $H_2(1 - \tau) = H_2(\tau)$ ) and concavity of the function  $H_2$ , one can show that the maximum is achieved at  $\tau = 1/2$  and  $H_2(1/2) = 1$ . Thus, in the worst case we have  $\alpha + \beta \leq 33/32$  upper bound, when roughly half of the entries are chosen uniformly at random. For other values of  $q$ , it is even closer to the optimum; numerically  $H_2(\tau) \approx (4\tau(1 - \tau))^{3/4}$ ,  $0 \leq \tau \leq 1$ .



**Figure 7:** Comparison of general uncertainty principle  $\alpha \cdot 4^{b/d} \geq 1$  (dashed red line) of Safaryan et al. [2020] against the new linear version (34) (dashed blue line). Each color represents one compression method: yellow for usual random sparsification with uniform probabilities and orange for greedy sparsification (a.k.a Top- $k$  sparsification). Each triangle marker indicates one particular  $d = 10^3$  dimensional vector randomly generated from Gaussian distribution, which subsequently gets compressed by the compression operator mentioned in the legend.

## G Proofs

### G.1 Proof of Theorem 8

Using smoothness of  $f$ , we have

$$\begin{aligned}
\mathbb{E}f(x^{k+1}) &= \mathbb{E}f(x^k - \gamma \mathbf{C} \nabla f(x^k)) \\
&\leq f(x^k) - \gamma \langle \nabla f(x^k), \mathbb{E}[\mathbf{C} \nabla f(x^k)] \rangle + \frac{\gamma^2}{2} \mathbb{E}[\|\mathbf{C} \nabla f(x^k)\|_{\mathbf{L}}^2] \\
&= f(x^k) - \gamma \|\nabla f(x^k)\|^2 + \frac{\gamma^2}{2} \|\nabla f(x^k)\|_{\mathbb{E}[\mathbf{CLC}]}^2 \\
&\leq f(x^k) - \gamma (2 - \gamma \lambda_{\max}(\mathbb{E}[\mathbf{CLC}])) \cdot \frac{1}{2} \|\nabla f(x^k)\|^2.
\end{aligned} \tag{35}$$

Computing the expectation inside, we get

$$\mathbb{E}[\mathbf{CLC}] = \mathbb{E}\left[(c_i c_j \mathbf{L}_{ij})_{i,j=1}^d\right] = \left(\frac{p_{ij} \mathbf{L}_{ij}}{p_i p_j}\right)_{i,j=1}^d = (\text{Diag}(1/p) \mathbf{P} \text{Diag}(1/p)) \circ \mathbf{L} = \bar{\mathbf{P}} \circ \mathbf{L}. \tag{36}$$

Therefore, using the bound for the step size  $\gamma$  and strong convexity of  $f$ , we get

$$\begin{aligned}
\mathbb{E}[f(x^{k+1}) - f(x^*)] &\leq (f(x^k) - f(x^*)) - \gamma (2 - \gamma \lambda_{\max}(\bar{\mathbf{P}} \circ \mathbf{L})) \cdot \frac{1}{2} \|\nabla f(x^k)\|^2 \\
&\leq (f(x^k) - f(x^*)) - \frac{\gamma}{2} \|\nabla f(x^k)\|^2 \\
&\leq (1 - \gamma \mu) (f(x^k) - f(x^*)),
\end{aligned} \tag{37}$$

repeated application of which completes the proof.

### G.2 Proof of Theorem 12

The following lemmas will be useful to handle the computation with pseudo-inverses.

**Lemma 16** (Lemma E.2 and E.3 [Hanzely and Richtárik, 2019b]). *If  $f$  is convex and  $\mathbf{L}$ -smooth, then for any  $x, y \in \mathbb{R}^d$*

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^\dagger}^2. \tag{38}$$

*If, in addition,  $f$  is bounded below, then  $\nabla f(x) \in \text{Range}(\mathbf{L}^\dagger) = \text{Range}(\mathbf{L})$  for all  $x \in \mathbb{R}^d$ .*

**Lemma 17.** *With  $\bar{\mathbf{C}} = \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger/2}$ , the following holds*

$$\mathbb{E}[\mathbf{L}^{1/2} (\bar{\mathbf{C}} - \mathbf{I})^\top (\bar{\mathbf{C}} - \mathbf{I}) \mathbf{L}^{1/2}] = \mathbf{L}^{1/2} \mathbf{L}^{\dagger/2} (\tilde{\mathbf{P}} \circ \mathbf{L}) \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2}. \tag{39}$$

*Proof.* Using the property  $\mathbf{L}^{1/2} \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2} = \mathbf{L}^{1/2}$  of pseudoinverse, we have

$$\begin{aligned}
\mathbb{E}[\mathbf{L}^{1/2} (\bar{\mathbf{C}} - \mathbf{I})^\top (\bar{\mathbf{C}} - \mathbf{I}) \mathbf{L}^{1/2}] &= \mathbb{E}[\mathbf{L}^{1/2} (\mathbf{L}^{\dagger/2} \mathbf{C} \mathbf{L}^{1/2} - \mathbf{I}) (\mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger/2} - \mathbf{I}) \mathbf{L}^{1/2}] \\
&= \mathbb{E}[\mathbf{L}^{1/2} (\mathbf{L}^{\dagger/2} \mathbf{C} \mathbf{L} \mathbf{C} \mathbf{L}^{\dagger/2} - \mathbf{L}^{\dagger/2} \mathbf{C} \mathbf{L}^{1/2} - \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger/2} + \mathbf{I}) \mathbf{L}^{1/2}] \\
&\stackrel{(36)}{=} \mathbf{L}^{1/2} (\mathbf{L}^{\dagger/2} (\bar{\mathbf{P}} \circ \mathbf{L}) \mathbf{L}^{\dagger/2} - \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2} - \mathbf{L}^{1/2} \mathbf{L}^{\dagger/2} + \mathbf{I}) \mathbf{L}^{1/2} \\
&= \mathbf{L}^{1/2} (\mathbf{L}^{\dagger/2} (\bar{\mathbf{P}} \circ \mathbf{L}) \mathbf{L}^{\dagger/2} - \mathbf{L}^{\dagger/2} \mathbf{L} \mathbf{L}^{\dagger/2}) \mathbf{L}^{1/2} \\
&\quad + \mathbf{L}^{1/2} (\mathbf{L}^{\dagger/2} \mathbf{L} \mathbf{L}^{\dagger/2} - \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2} - \mathbf{L}^{1/2} \mathbf{L}^{\dagger/2} + \mathbf{I}) \mathbf{L}^{1/2} \\
&= \mathbf{L}^{1/2} \mathbf{L}^{\dagger/2} (\bar{\mathbf{P}} \circ \mathbf{L} - \mathbf{L}) \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2} + \mathbf{L}^{1/2} (\mathbf{I} - \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2}) (\mathbf{I} - \mathbf{L}^{1/2} \mathbf{L}^{\dagger/2}) \mathbf{L}^{1/2} \\
&= \mathbf{L}^{1/2} \mathbf{L}^{\dagger/2} (\tilde{\mathbf{P}} \circ \mathbf{L}) \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2} + (\mathbf{L}^{1/2} - \mathbf{L}^{1/2} \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2}) (\mathbf{L}^{1/2} - \mathbf{L}^{1/2} \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2}) \\
&= \mathbf{L}^{1/2} \mathbf{L}^{\dagger/2} (\tilde{\mathbf{P}} \circ \mathbf{L}) \mathbf{L}^{\dagger/2} \mathbf{L}^{1/2}.
\end{aligned}$$

□

For convenience we skip iteration count  $k$ , and write  $x, x^+$  instead of  $x^k, x^{k+1}$ . Using non-expansiveness of the prox operator we get

$$\begin{aligned}
\mathbb{E} [\|x^+ - x^*\|^2] &\leq \mathbb{E} \left[ \|x - x^* - \gamma \left( \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2} \right) \nabla f(x) + \gamma \nabla f(x^*) \|^2 \right] \\
&= \|x - x^*\|^2 - 2\gamma \langle x - x^*, \nabla f(x) - \nabla f(x^*) \rangle + \gamma^2 \mathbb{E} \left[ \left\| \left( \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2} \right) \nabla f(x) - \nabla f(x^*) \right\|^2 \right] \\
&\leq \|x - x^*\|^2 - 2\gamma \langle x - x^*, \nabla f(x) - \nabla f(x^*) \rangle \\
&\quad + 2\gamma^2 \mathbb{E} \left[ \left\| \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2} (\nabla f(x) - \nabla f(x^*)) \right\|^2 \right] + 2\gamma^2 \mathbb{E} \left[ \left\| \left( \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2} - \mathbf{I} \right) \nabla f(x^*) \right\|^2 \right] \\
&\leq \|x - x^*\|^2 - 2\gamma \langle x - x^*, \nabla f(x) - \nabla f(x^*) \rangle \\
&\quad + 2\gamma^2 \lambda_{\max}(\mathbb{E} [\mathbf{C} \mathbf{L} \mathbf{C}]) \left\| \mathbf{L}^{\dagger 1/2} (\nabla f(x) - \nabla f(x^*)) \right\|^2 + 2\gamma^2 \mathbb{E} \left[ \left\| \left( \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2} - \mathbf{I} \right) \nabla f(x^*) \right\|^2 \right] \\
&\stackrel{(36), (40)}{\leq} \|x - x^*\|^2 - 2\gamma \langle x - x^*, \nabla f(x) - \nabla f(x^*) \rangle \\
&\quad + 2\gamma^2 \lambda_{\max}(\tilde{\mathbf{P}} \circ \mathbf{L}) \left\| \nabla f(x) - \nabla f(x^*) \right\|_{\mathbf{L}^\dagger}^2 + 2\gamma^2 \lambda_{\max}(\tilde{\mathbf{P}} \circ \mathbf{L}) \left\| \nabla f(x^*) \right\|_{\mathbf{L}^\dagger}^2 \\
&= \|x - x^*\|^2 - 2\gamma \langle x - x^*, \nabla f(x) - \nabla f(x^*) \rangle + 2\gamma^2 \tilde{\mathcal{L}} \left\| \nabla f(x) - \nabla f(x^*) \right\|_{\mathbf{L}^\dagger}^2 + 2\gamma^2 \tilde{\mathcal{L}} \left\| \nabla f(x^*) \right\|_{\mathbf{L}^\dagger}^2,
\end{aligned}$$

where we used  $\mathbb{E} [\mathbf{C} \mathbf{L} \mathbf{C}] = \tilde{\mathbf{P}} \circ \mathbf{L}$  based on (36) and for the last term we used Lemma 16 to represent  $\nabla f(x^*) = \mathbf{L}^{1/2} g_*$  and then applied Lemma 17

$$\begin{aligned}
\mathbb{E} \left[ \left\| \left( \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2} - \mathbf{I} \right) \nabla f(x^*) \right\|^2 \right] &= \mathbb{E} \left[ g_*^\top \mathbf{L}^{1/2} \left( \mathbf{L}^{\dagger 1/2} \mathbf{C} \mathbf{L}^{1/2} - \mathbf{I} \right) \left( \mathbf{L}^{1/2} \mathbf{C} \mathbf{L}^{\dagger 1/2} - \mathbf{I} \right) \mathbf{L}^{1/2} g_* \right] \\
&= \nabla f(x^*)^\top \left( \mathbf{L}^{\dagger 1/2} \left( \tilde{\mathbf{P}} \circ \mathbf{L} \right) \mathbf{L}^{\dagger 1/2} \right) \nabla f(x^*) \\
&\leq \lambda_{\max}(\tilde{\mathbf{P}} \circ \mathbf{L}) \left\| \nabla f(x^*) \right\|_{\mathbf{L}^\dagger}^2.
\end{aligned} \tag{40}$$

Using the bound on step size  $\gamma \leq 1/2\tilde{\mathcal{L}}$ , strong convexity of  $f$  and (38), we continue as follows

$$\begin{aligned}
\mathbb{E} [\|x^+ - x^*\|^2] &\leq \|x - x^*\|^2 - \gamma \langle x - x^*, \nabla f(x) - \nabla f(x^*) \rangle \\
&\quad - \gamma \left( \langle x - x^*, \nabla f(x) - \nabla f(x^*) \rangle - \left\| \nabla f(x) - \nabla f(x^*) \right\|_{\mathbf{L}^\dagger}^2 \right) \\
&\quad + 2\gamma^2 \tilde{\mathcal{L}} \left\| \nabla f(x^*) \right\|_{\mathbf{L}^\dagger}^2 \\
&\stackrel{(38)}{\leq} (1 - \gamma\mu) \|x - x^*\|^2 + 2\gamma^2 \tilde{\mathcal{L}} \left\| \nabla f(x^*) \right\|_{\mathbf{L}^\dagger}^2.
\end{aligned}$$

Telescoping the above inequality, we complete the proof.

### G.3 Proof of Theorem 2

**Proof technique.** First we show the unbiasedness of  $g^k$ . As smoothness matrices  $\mathbf{L}_i$  are not necessarily invertible, terms like  $\mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2}$  show up in the analysis and block chains of cancellations. This part is handled by the fact that gradients  $\nabla f_i(x)$  of an  $\mathbf{L}_i$ -smooth function are constraint to remain in  $\text{Range } \mathbf{L}_i$  and the mapping associated with the matrix  $\mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2}$  is identity on the subspace  $\text{Range } (\mathbf{L}_i)$ . Second part is the tight estimation of  $\mathbb{E}_k \|g^k - \nabla f(x^*)\|^2$ , which describes the progress of the method in the presence of stochasticity. Key part is getting the decomposition

$$\mathbb{E}_k [\|g^k - \nabla f(x^*)\|^2] = \left\| \nabla f(x^k) - \nabla f(x^*) \right\|^2 + \frac{1}{n^2} \sum_{i=1}^n \left\| \nabla f_i(x^k) \right\|_{\mathbf{L}_i^{\dagger 1/2} (\tilde{\mathbf{P}}_i \circ \mathbf{L}_i) \mathbf{L}_i^{\dagger 1/2}}^2, \tag{41}$$

which shows the exact interaction between random sketches and local smoothness. We complete the proof using the unified convergence theory of Gorbunov et al. [2020a].

**Proof.**

In this proof we skip the iteration count  $k$  to simplify the notation. Define

$$\begin{aligned}
\mathbf{M}_i &:= \mathbf{L}_i^{1/2} \mathbb{E} [(\bar{\mathbf{C}}_i - \mathbf{I})^\top (\bar{\mathbf{C}}_i - \mathbf{I})] \mathbf{L}_i^{1/2} \\
&\stackrel{(39)}{=} \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\tilde{\mathbf{P}}_i \circ \mathbf{L}_i) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \\
&= \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\bar{\mathbf{P}}_i \circ \mathbf{L}_i - \mathbf{L}_i) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \\
&= \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\bar{\mathbf{P}}_i \circ \mathbf{L}_i) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} - \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \\
&= \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\bar{\mathbf{P}}_i \circ \mathbf{L}_i) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} - \mathbf{L}_i \\
&= \mathbf{L}_i^{1/2} (\mathbb{E} [\bar{\mathbf{C}}_i^\top \bar{\mathbf{C}}_i] - \mathbf{I}) \mathbf{L}_i^{1/2}.
\end{aligned} \tag{42}$$

We are going to estimate the moment  $\mathbb{E} [\|g(x) - \nabla f(x^*)\|^2]$  and show the following bound for the gradient estimator  $g(x) = \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{C}}_i \nabla f_i(x)$  (see line 5 of Algorithm 1):

$$\mathbb{E} [\|g(x) - \nabla f(x^*)\|^2] \leq 2 \left( L + \frac{2\tilde{\mathcal{L}}}{n} \right) D_f(x, x^*) + \frac{2\sigma^*}{n}.$$

Due to Lemma 16, we have  $\nabla f_i(x) = \mathbf{L}_i^{1/2} r_i$  for some  $r_i$ . Therefore

$$\mathbb{E} [\bar{\mathbf{C}}_i \nabla f_i(x)] = \mathbb{E} [\mathbf{L}_i^{1/2} \mathbf{C}_i \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} r_i] = \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} r_i = \mathbf{L}_i^{1/2} r_i = \nabla f_i(x), \tag{43}$$

which implies unbiasedness of the estimator  $g(x)$ , namely  $\mathbb{E} [g(x)] = \nabla f(x)$ . Next, note that

$$\begin{aligned}
\mathbb{E} [\|g(x) - \nabla f(x^*)\|^2] &= \mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{C}}_i \nabla f_i(x) - \nabla f(x^*) \right\|^2 \right] \\
&= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|\bar{\mathbf{C}}_i \nabla f_i(x) - \nabla f(x^*)\|^2] + \frac{1}{n^2} \sum_{i \neq j} \mathbb{E} \langle \bar{\mathbf{C}}_i \nabla f_i(x) - \nabla f(x^*), \bar{\mathbf{C}}_j \nabla f_j(x) - \nabla f(x^*) \rangle \\
&= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|\bar{\mathbf{C}}_i \nabla f_i(x)\|^2] + \|\nabla f(x^*)\|^2 - 2\mathbb{E} \langle \bar{\mathbf{C}}_i \nabla f_i(x), \nabla f(x^*) \rangle + \frac{1}{n^2} \sum_{i \neq j} \langle \nabla f_i(x) - \nabla f(x^*), \nabla f_j(x) - \nabla f(x^*) \rangle \\
&= \frac{1}{n^2} \sum_{i=1}^n \|\nabla f_i(x)\|_{\mathbb{E}[\bar{\mathbf{C}}_i^\top \bar{\mathbf{C}}_i]}^2 + \|\nabla f(x^*)\|^2 - 2\langle \nabla f_i(x), \nabla f(x^*) \rangle + \|\nabla f(x) - \nabla f(x^*)\|^2 - \frac{1}{n^2} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x^*)\|^2 \\
&= \frac{1}{n^2} \sum_{i=1}^n \left\| \mathbf{L}_i^{1/2} r_i \right\|_{\mathbb{E}[\bar{\mathbf{C}}_i^\top \bar{\mathbf{C}}_i] - \mathbf{I}}^2 + \frac{1}{n^2} \sum_{i=1}^n \|\nabla f_i(x)\|^2 + \|\nabla f(x^*)\|^2 - 2\langle \nabla f_i(x), \nabla f(x^*) \rangle \\
&\quad + \|\nabla f(x) - \nabla f(x^*)\|^2 - \frac{1}{n^2} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x^*)\|^2 \\
&= \|\nabla f(x) - \nabla f(x^*)\|^2 + \frac{1}{n^2} \sum_{i=1}^n \|r_i\|_{\mathbf{M}_i}^2 \\
&= \|\nabla f(x) - \nabla f(x^*)\|^2 + \frac{1}{n^2} \sum_{i=1}^n \|r_i\|_{\mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\tilde{\mathbf{P}}_i \circ \mathbf{L}_i) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2}}^2 \\
&= \|\nabla f(x) - \nabla f(x^*)\|^2 + \frac{1}{n^2} \sum_{i=1}^n \left\| \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x) \right\|_{\tilde{\mathbf{P}}_i \circ \mathbf{L}_i}^2,
\end{aligned}$$

which gives as the following decomposition

$$\mathbb{E} [\|g(x) - \nabla f(x^*)\|^2] = \|\nabla f(x) - \nabla f(x^*)\|^2 + \frac{1}{n^2} \sum_{i=1}^n \left\| \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x) \right\|_{\tilde{\mathbf{P}}_i \circ \mathbf{L}_i}^2. \tag{44}$$

For the first term it can be bounded using convexity and smoothness of  $f$ , namely  $\|\nabla f(x) - \nabla f(x^*)\|^2 \leq 2LD_f(x, x^*)$ . For the second term we proceed as follows

$$\begin{aligned}
\frac{1}{n^2} \sum_{i=1}^n \left\| \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x) \right\|_{\tilde{\mathbf{P}}_i \circ \mathbf{L}_i}^2 &\leq \frac{1}{n^2} \sum_{i=1}^n \lambda_{\max}(\tilde{\mathbf{P}}_i \circ \mathbf{L}_i) \|\mathbf{L}_i^{\dagger 1/2} \nabla f_i(x)\|^2 \\
&= \frac{1}{n^2} \sum_{i=1}^n \tilde{\mathcal{L}}_i \|\nabla f_i(x)\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq \frac{2}{n^2} \sum_{i=1}^n \tilde{\mathcal{L}}_i \|\nabla f_i(x) - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + \frac{2}{n^2} \sum_{i=1}^n \tilde{\mathcal{L}}_i \|\nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq \frac{2}{n^2} \sum_{i=1}^n 2\tilde{\mathcal{L}}_i D_{f_i}(x, x^*) + \frac{2\sigma^*}{n} \\
&= \frac{4\tilde{\mathcal{L}}_{\max}}{n} D_f(x, x^*) + \frac{2\sigma^*}{n}.
\end{aligned} \tag{45}$$

Combining these two estimates, we get

$$\mathbb{E} [\|g(x) - \nabla f(x^*)\|^2] \leq 2 \left( L + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) D_f(x, x^*) + \frac{2\sigma^*}{n}.$$

It remains to apply the result of Gorbunov et al. [2020a].

#### G.4 Proof of Theorem 3

**Proof technique.** The structure of the proof resembles the one for DCGD+. With the introduced shift vectors, the unbiasedness of  $g^k$  additionally requires  $h_i^k \in \text{Range}(\mathbf{L}_i)$ . This is resolved by the initialization  $h_i^0 \in \text{Range}(\mathbf{L}_i)$  and linear update rule for  $h_i^{k+1}$  in line 5. The proof develops a decomposition similar to (41) with modified second term  $\sigma^k := \frac{1}{n} \sum_{i=1}^n \|h_i^k - \nabla f(x^*)\|_{\mathbf{L}_i^\dagger}^2$  involving shifts  $h_i^k$ . To avoid the neighborhood term in (10) and guarantee a linear convergence for  $x^k$ , we make  $\sigma^k$  converge linearly too. Key technical part of the proof is to establish contracting recurrence relation for  $\sigma^k$  which boils down to  $\mathbb{E}[\bar{\mathbf{C}}_i^\top \mathbf{L}_i^\dagger \bar{\mathbf{C}}_i] \preceq (\omega_i + 1) \mathbf{L}_i^\dagger$ . The latter bound justifies the structure of  $\bar{\mathbf{C}}_i$  as it filters the interaction between compression and smoothness mixed in the expectation and separates variance  $\omega_i$  of compression from smoothness matrix  $\mathbf{L}_i$ .

**Proof.** First, we show the unbiasedness of the estimator  $g(x^k)$ . In (43), we showed unbiasedness of  $\bar{\mathbf{C}}_i^k \nabla f_i(x^k)$  using inclusion  $\nabla f_i(x^k) \in \text{Range}(\mathbf{L}_i)$ . Assume for a moment that we also have  $h_i^k \in \text{Range}(\mathbf{L}_i)$ . Hence, in the same way we can show  $\mathbb{E}_k [\bar{\mathbf{C}}_i^k h_i^k] = h_i^k$ , which implies the unbiasedness of  $g^k$  as

$$\mathbb{E}_k [g^k] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_k [\bar{\mathbf{C}}_i^k \nabla f_i(x^k)] - \mathbb{E}_k [\bar{\mathbf{C}}_i^k h_i^k] + h_i^k = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) = \nabla f(x^k).$$

The inclusion  $h_i^k \in \text{Range}(\mathbf{L}_i)$  follows from the initialization  $h_i^0 \in \text{Range}(\mathbf{L}_i)$  (see line 1 of Algorithm 2) and linear update rule of  $h_i^{k+1} = h_i^k + \alpha \mathbf{L}_i^{1/2} \Delta_i^k$  (see line 5 of Algorithm 2). As both  $\nabla f_i(x^k)$  and  $h_i^k$  belong to  $\text{Range}(\mathbf{L}_i)$ , denote  $\nabla f_i(x^k) - h_i^k = \mathbf{L}_i^{1/2} r_i^k$ . Next we bound

$$\begin{aligned}
\mathbb{E} [\|g^k - \nabla f(x^*)\|^2] &= \|\nabla f(x^k) - \nabla f(x^*)\|^2 + \mathbb{E} [\|g^k - \nabla f(x^k)\|^2] \\
&\leq 2LD_f(x^k, x^*) + \mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{C}}_i^k (\nabla f_i(x^k) - h_i^k) + h_i^k - \nabla f_i(x^k) \right\|^2 \right] \\
&= 2LD_f(x^k, x^*) + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \left\| (\bar{\mathbf{C}}_i^k - \mathbf{I}) \mathbf{L}_i^{1/2} r_i^k \right\|^2 \right] \\
&= 2LD_f(x^k, x^*) + \frac{1}{n^2} \sum_{i=1}^n \|r_i^k\|^2_{\mathbb{E}[\mathbf{L}_i^{1/2} (\bar{\mathbf{C}}_i^k - \mathbf{I})^\top (\bar{\mathbf{C}}_i^k - \mathbf{I}) \mathbf{L}_i^{1/2}]} \\
&\stackrel{(42)}{=} 2LD_f(x^k, x^*) + \frac{1}{n^2} \sum_{i=1}^n \|r_i^k\|_{\mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\tilde{\mathbf{P}}_i \circ \mathbf{L}_i) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2}}^2 \\
&= 2LD_f(x^k, x^*) + \frac{1}{n^2} \sum_{i=1}^n \left\| \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k) \right\|_{\tilde{\mathbf{P}}_i \circ \mathbf{L}_i}^2 \\
&\leq 2LD_f(x^k, x^*) + \frac{\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|\nabla f_i(x^k) - h_i^k\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq 2LD_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|\nabla f_i(x^k) - f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq 2LD_f(x^k, x^*) + \frac{4\tilde{\mathcal{L}}_{\max}}{n} D_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \\
&= 2 \left( L + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) D_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2
\end{aligned} \tag{46}$$

Then we deduce a recurrence relation for the last term  $\sigma^k := \frac{1}{n} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2$ . For that we will need the following bounds

$$0 \preceq \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} \preceq \mathbf{I}, \tag{47}$$

which can be proved via SVD and eigenvalue decompositions. Since  $\mathbf{L}_i$  is square, symmetric and positive semidefinite, we know that singular value decomposition and eigenvalue decompositions are the same. Let  $\mathbf{L}_i^{1/2} = \mathbf{U}_i \mathbf{D}_i \mathbf{U}_i^\top$ , where  $\mathbf{D}_i$  is diagonal and  $\mathbf{U}_i$  is orthogonal so that  $\mathbf{U}_i^\top = \mathbf{U}_i^{-1}$ . Then

$$\mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} = \mathbf{U}_i \mathbf{D}_i \mathbf{U}_i^\top \mathbf{U}_i \mathbf{D}_i^{\dagger 2} \mathbf{U}_i^\top \mathbf{U}_i \mathbf{D}_i \mathbf{U}_i^\top = \mathbf{U}_i \left( \mathbf{D}_i \mathbf{D}_i^{\dagger 2} \mathbf{D}_i \right) \mathbf{U}_i^\top = \mathbf{U}_i \left( \mathbf{D}_i \mathbf{D}_i^\dagger \right) \mathbf{U}_i^\top,$$

which can admit eigenvalues only in  $[0, 1]$  since the matrix  $\mathbf{D}_i \mathbf{D}_i^\dagger$  is diagonal with entries either 0 or 1. Denote

$$\omega_i = \lambda_{\max} \left( \mathbb{E} [(\mathbf{C}_i^k)^2] \right) - 1 = \max_{1 \leq j \leq d} \frac{1}{p_{i;j}} - 1. \tag{48}$$

and bound each summand of  $\sigma^{k+1}$  as follows

$$\begin{aligned}
\mathbb{E}_k \left[ \|h_i^{k+1} - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \right] &= \mathbb{E}_k \left[ \|h_i^k - \nabla f_i(x^*) + \alpha \bar{\Delta}_i^k\|_{\mathbf{L}_i^\dagger}^2 \right] \\
&= \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + 2\alpha \langle h_i^k - \nabla f_i(x^*), \nabla f_i(x^k) - h_i^k \rangle_{\mathbf{L}_i^\dagger} + \alpha^2 \mathbb{E} \left[ \|\bar{\mathbf{C}}_i^k (\nabla f_i(x^k) - h_i^k)\|_{\mathbf{L}_i^\dagger}^2 \right] \\
&= \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + 2\alpha \langle h_i^k - \nabla f_i(x^*), \nabla f_i(x^k) - h_i^k \rangle_{\mathbf{L}_i^\dagger} + \alpha^2 \|\nabla f_i(x^k) - h_i^k\|_{\mathbb{E}[(\bar{\mathbf{C}}_i^k)^\top \mathbf{L}_i^\dagger \bar{\mathbf{C}}_i^k]}^2 \\
&\leq \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + 2\alpha \langle h_i^k - \nabla f_i(x^*), \nabla f_i(x^k) - h_i^k \rangle_{\mathbf{L}_i^\dagger} + \alpha^2 \|\nabla f_i(x^k) - h_i^k\|_{\mathbf{L}_i^{\dagger 1/2} \mathbb{E}[(\mathbf{C}_i^k)^2] \mathbf{L}_i^{\dagger 1/2}}^2 \\
&\leq \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + 2\alpha \langle h_i^k - \nabla f_i(x^*), \nabla f_i(x^k) - h_i^k \rangle_{\mathbf{L}_i^\dagger} + \alpha^2 (1 + \omega_i) \|\nabla f_i(x^k) - h_i^k\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + 2\alpha \langle h_i^k - \nabla f_i(x^*), \nabla f_i(x^k) - h_i^k \rangle_{\mathbf{L}_i^\dagger} + \alpha \|\nabla f_i(x^k) - h_i^k\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq (1 - \alpha) \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + \alpha \|\nabla f_i(x^k) - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2,
\end{aligned}$$

where we used bounds  $\alpha \leq \frac{1}{1+\omega_i}$  and

$$\mathbb{E} \left[ (\bar{\mathbf{C}}_i^k)^\top \mathbf{L}_i^\dagger \bar{\mathbf{C}}_i^k \right] = \mathbf{L}_i^{\dagger 1/2} \mathbb{E} \left[ \mathbf{C}_i^k \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} \mathbf{C}_i^k \right] \mathbf{L}_i^{\dagger 1/2} \preceq \mathbf{L}_i^{\dagger 1/2} \mathbb{E} \left[ (\mathbf{C}_i^k)^2 \right] \mathbf{L}_i^{\dagger 1/2}.$$

Therefore

$$\begin{aligned}
\mathbb{E}_k [\sigma^{k+1}] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_k \left[ \|h_i^{k+1} - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \right] \\
&\leq \frac{1-\alpha}{n} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + \frac{\alpha}{n} \sum_{i=1}^n \|\nabla f_i(x^k) - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq (1-\alpha)\sigma^k + \frac{2\alpha}{n} \sum_{i=1}^n D_{f_i}(x^k, x^*) \\
&= (1-\alpha)\sigma^k + 2\alpha D_f(x^k, x^*).
\end{aligned}$$

Thus, with  $\alpha \leq \frac{1}{1+\omega_{\max}}$ , the estimator  $g^k$  of Algorithm 2 satisfies

$$\begin{aligned}
\mathbb{E}_k [g^k] &= \nabla f(x^k) \\
\mathbb{E}_k [\|g^k - \nabla f(x^*)\|^2] &\leq 2 \left( L + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) D_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \sigma^k \\
\mathbb{E}_k [\sigma^{k+1}] &\leq (1-\alpha)\sigma^k + 2\alpha D_f(x^k, x^*).
\end{aligned}$$

It remains to apply Theorem 4.1 [Gorbunov et al., 2020a] with parameters  $A = L + \frac{2}{n}\tilde{\mathcal{L}}_{\max}$ ,  $B = \frac{2}{n}\tilde{\mathcal{L}}_{\max}$ ,  $\rho = \alpha$ ,  $C = \alpha$  and  $M = \frac{4}{\alpha n}\tilde{\mathcal{L}}_{\max}$ ,  $A + CM = L + \frac{6}{n}\tilde{\mathcal{L}}_{\max}$ ,  $1 + \frac{B}{M} - \rho = 1 - \frac{\alpha}{2}$ .

## G.5 Proof of Theorem 4

**Proof technique.** The additional difficulty that acceleration brings on top of variance reduction is the modified term  $H^k := \frac{1}{n} \sum_{i=1}^n \|h_i^k - \nabla f_i(w^k)\|_{\mathbf{L}_i^\dagger}^2$  controlling variance reduction process. The subtlety of  $H^k$  in contrast to  $\sigma^k$  is gradients  $\nabla f_i(w^k)$  which are not fixed. Key technical part is to reduce contracting property of  $H^k$  into upper bounding  $\mathbb{E}[(\mathbf{I} - \alpha \bar{\mathbf{C}}_i)^\top \mathbf{L}_i^\dagger (\mathbf{I} - \alpha \bar{\mathbf{C}}_i)]$  by  $(1-\alpha)\mathbf{L}_i^\dagger$  as quadratic forms in the subspace  $\text{Range}(\mathbf{L}_i)$ .

**Proof.** Following the analysis of Li et al. [2020], define

$$Z^k := \|z^k - x^*\|^2, \quad Y^k := F(y^k) - F(x^*), \quad W^k := F(w^k) - F(x^*),$$

$$H^k := \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(w^k) - h_i^k\|_{\mathbf{L}_i^\dagger}^2.$$

**Lemma 18** (Lemma 2, [Li et al., 2020]). *Let  $\eta \leq \frac{1}{2L}$ ,  $\theta_1 \leq \frac{1}{4}$ ,  $\theta_2 = \frac{1}{2}$ ,  $\gamma = \frac{\eta}{2(\theta_1 + \eta\mu)}$  and  $\beta = 1 - \gamma\mu$ . Then*

$$\begin{aligned} \mathbb{E}[Z^{k+1}] + \frac{2\gamma\beta}{\theta_1} \mathbb{E}[Y^{k+1}] &\leq \beta Z^k + (1 - \theta_1 - \theta_2) \frac{2\gamma\beta}{\theta_1} Y^k + 2\gamma\beta \frac{\theta_2}{\theta_1} W^k + \frac{\gamma\eta}{\theta_1} \mathbb{E}[\|g^k - \nabla f(x^k)\|^2] \\ &\quad - \frac{\gamma}{4n\theta_1} \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 - \frac{\gamma}{8n\theta_1} \sum_{i=1}^n \|\nabla f_i(y^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2. \end{aligned}$$

*Proof.* Proof is the same as for the original lemma except we use  $\mathbf{L}_i$ -smoothness of  $f_i$  via (38).

$$f_i(u) \geq f_i(x^k) + \langle \nabla f_i(x^k), u - x^k \rangle + \frac{1}{2} \|\nabla f_i(u) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2.$$

□

**Lemma 19** (Lemma 3, [Li et al., 2020]).

$$\mathbb{E}[W^{k+1}] = (1 - q)W^k + qY^k.$$

**Lemma 20** (Lemma 4, [Li et al., 2020]).

$$\mathbb{E}[\|g^k - \nabla f(x^k)\|^2] \leq \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 + \frac{2\tilde{\mathcal{L}}_{\max}}{n} H^k.$$

*Proof.* Let  $\nabla f_i(x^k) - h_i^k = \mathbf{L}_i^{1/2} r_i^k$ . Then

$$\begin{aligned} \mathbb{E}[\|g^k - \nabla f(x^k)\|^2] &= \mathbb{E}\left[\left\|\frac{1}{n} \sum_{i=1}^n \bar{\mathbf{C}}_i^k (\nabla f_i(x^k) - h_i^k) - (\nabla f_i(x^k) - h_i^k)\right\|^2\right] \\ &= \frac{1}{n^2} \mathbb{E}\left[\left\|\sum_{i=1}^n (\bar{\mathbf{C}}_i^k - \mathbf{I})(\nabla f_i(x^k) - h_i^k)\right\|^2\right] = \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}\left[\|(\bar{\mathbf{C}}_i^k - \mathbf{I})\mathbf{L}_i^{1/2} r_i^k\|^2\right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \|r_i^k\|_{\mathbf{L}_i^{1/2}}^2 \mathbb{E}[(\bar{\mathbf{C}}_i^k - \mathbf{I})^\top (\bar{\mathbf{C}}_i^k - \mathbf{I}) \mathbf{L}_i^{1/2}] \stackrel{(42)}{=} \frac{1}{n^2} \sum_{i=1}^n \|r_i^k\|_{\mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\tilde{\mathbf{P}}_i \circ \mathbf{L}_i) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2}}^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \|\mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k)\|_{\tilde{\mathbf{P}}_i \circ \mathbf{L}_i}^2 \leq \frac{\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|\nabla f_i(x^k) - h_i^k\|_{\mathbf{L}_i^\dagger}^2 \\ &\leq \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|\nabla f_i(x^k) - \nabla f_i(w^k)\|_{\mathbf{L}_i^\dagger}^2 + \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|\nabla f_i(w^k) - h_i^k\|_{\mathbf{L}_i^\dagger}^2. \end{aligned}$$

□

**Lemma 21** (Lemma 5, [Li et al., 2020]). *If  $\alpha \leq \frac{1}{1 + \omega_{\max}}$ , where  $\omega_{\max} = \max_{1 \leq i \leq n} \omega_i$  and  $\omega_i = \max_{1 \leq j \leq d} \frac{1}{p_{i,j}} - 1$ , then*

$$\mathbb{E}[H^{k+1}] \leq \left(1 - \frac{\alpha}{2}\right) H^k + \left(1 + \frac{2q}{\alpha}\right) \frac{2q}{n} \left( \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 + \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \right).$$



*Proof.* We start bounding the summands of  $H^{k+1}$ . Let  $\nabla f_i(w^k) - h_i^k = \mathbf{L}_i^{1/2} r_i^k$ .

$$\begin{aligned}
\mathbb{E}_k \left[ \|\nabla f_i(w^{k+1}) - h_i^{k+1}\|_{\mathbf{L}_i^\dagger}^2 \right] &= q \mathbb{E}_k \left[ \|\nabla f_i(y^k) - h_i^{k+1}\|_{\mathbf{L}_i^\dagger}^2 \right] + (1-q) \mathbb{E}_k \left[ \|\nabla f_i(w^k) - h_i^{k+1}\|_{\mathbf{L}_i^\dagger}^2 \right] \\
&\leq q \left( 1 + \frac{2q}{\alpha} \right) \|\nabla f_i(w^k) - \nabla f_i(y^k)\|_{\mathbf{L}_i^\dagger}^2 + \left( 1 - q + \left( 1 + \frac{\alpha}{2q} \right) q \right) \mathbb{E} \left[ \|\nabla f_i(w^k) - h_i^{k+1}\|_{\mathbf{L}_i^\dagger}^2 \right] \\
&= q \left( 1 + \frac{2q}{\alpha} \right) \|\nabla f_i(w^k) - \nabla f_i(y^k)\|_{\mathbf{L}_i^\dagger}^2 + \left( 1 + \frac{\alpha}{2} \right) \mathbb{E} \left[ \|\nabla f_i(w^k) - h_i^{k+1}\|_{\mathbf{L}_i^\dagger}^2 \right] \\
&= q \left( 1 + \frac{2q}{\alpha} \right) \|\nabla f_i(w^k) - \nabla f_i(y^k)\|_{\mathbf{L}_i^\dagger}^2 + \left( 1 + \frac{\alpha}{2} \right) \mathbb{E} \left[ \|(I - \alpha \bar{\mathbf{C}}_i^k)(\nabla f_i(w^k) - h_i^k)\|_{\mathbf{L}_i^\dagger}^2 \right] \\
&= q \left( 1 + \frac{2q}{\alpha} \right) \|\nabla f_i(w^k) - \nabla f_i(y^k)\|_{\mathbf{L}_i^\dagger}^2 + \left( 1 + \frac{\alpha}{2} \right) \|r_i^k\|_{\mathbf{L}_i^{1/2} \mathbb{E}[(I - \alpha \bar{\mathbf{C}}_i^k)^\top \mathbf{L}_i^\dagger (I - \alpha \bar{\mathbf{C}}_i^k)] \mathbf{L}_i^{1/2}}^2.
\end{aligned}$$

Next, we simplify the matrix of the second term.

$$\begin{aligned}
&\mathbf{L}_i^{1/2} \mathbb{E} \left[ (I - \alpha \bar{\mathbf{C}}_i^k)^\top \mathbf{L}_i^\dagger (I - \alpha \bar{\mathbf{C}}_i^k) \right] \mathbf{L}_i^{1/2} \\
&= \mathbb{E} \left[ \mathbf{L}_i^{1/2} (I - \alpha \mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i^k \mathbf{L}_i^{1/2}) \mathbf{L}_i^\dagger (I - \alpha \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2}) \mathbf{L}_i^{1/2} \right] \\
&= \mathbb{E} \left[ (\mathbf{L}_i^{1/2} - \alpha \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i^k \mathbf{L}_i^{1/2}) \mathbf{L}_i^\dagger (\mathbf{L}_i^{1/2} - \alpha \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2}) \right] \\
&= \mathbb{E} \left[ \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} - \alpha \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \right. \\
&\quad \left. - \alpha \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i^k \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} + \alpha^2 \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i^k \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \right] \\
&\stackrel{(47)}{\preceq} \mathbb{E} \left[ \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} - \alpha \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \right. \\
&\quad \left. - \alpha \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i^k \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} + \alpha^2 \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\mathbf{C}_i^k)^2 \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \right] \\
&= \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} - \alpha \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} - \alpha \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} + \alpha^2 \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbb{E}[(\mathbf{C}_i^k)^2] \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \\
&\stackrel{(48)}{\preceq} \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} - 2\alpha \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} + \alpha^2 (\omega_i + 1) \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \\
&= (1 - 2\alpha + \alpha^2 (\omega_i + 1)) \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2} \\
&\preceq (1 - \alpha) \mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2},
\end{aligned}$$

where in the last step we make use of the bound  $\alpha \leq \frac{1}{1+\omega_{\max}} = \min_{1 \leq i \leq n} \frac{1}{1+\omega_i}$ . Then we finish the recurrence as follows

$$\begin{aligned}
&\mathbb{E}_k \left[ \|\nabla f_i(w^{k+1}) - h_i^{k+1}\|_{\mathbf{L}_i^\dagger}^2 \right] \\
&\leq q \left( 1 + \frac{2q}{\alpha} \right) \|\nabla f_i(w^k) - \nabla f_i(y^k)\|_{\mathbf{L}_i^\dagger}^2 + \left( 1 + \frac{\alpha}{2} \right) \|r_i^k\|_{\mathbf{L}_i^{1/2} \mathbb{E}[(I - \alpha \bar{\mathbf{C}}_i^k)^\top \mathbf{L}_i^\dagger (I - \alpha \bar{\mathbf{C}}_i^k)] \mathbf{L}_i^{1/2}}^2 \\
&\leq q \left( 1 + \frac{2q}{\alpha} \right) \|\nabla f_i(w^k) - \nabla f_i(y^k)\|_{\mathbf{L}_i^\dagger}^2 + \left( 1 + \frac{\alpha}{2} \right) (1 - \alpha) \|r_i^k\|_{\mathbf{L}_i^{1/2} \mathbf{L}_i^\dagger \mathbf{L}_i^{1/2}}^2 \\
&= q \left( 1 + \frac{2q}{\alpha} \right) \|\nabla f_i(w^k) - \nabla f_i(y^k)\|_{\mathbf{L}_i^\dagger}^2 + \left( 1 + \frac{\alpha}{2} \right) (1 - \alpha) \|\nabla f_i(w^k) - h_i^k\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq 2q \left( 1 + \frac{2q}{\alpha} \right) \left( \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 + \|\nabla f_i(y^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \right) + \left( 1 - \frac{\alpha}{2} \right) \|\nabla f_i(w^k) - h_i^k\|_{\mathbf{L}_i^\dagger}^2.
\end{aligned}$$

Averaging over  $i \in [n]$  completes the proof.  $\square$

*Proof of Theorem 4.* Using the 4 lemmas above and  $\theta_1 \leq \frac{1}{4}$ ,  $\theta_2 = \frac{1}{2}$ , the Lyapunov function  $\Psi^{k+1}$  admits the following recurrence

$$\begin{aligned}
\mathbb{E} [\Psi^{k+1}] &:= \mathbb{E} \left[ Z^{k+1} + \frac{2\gamma\beta}{\theta_1} Y^{k+1} + 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1 q} W^{k+1} + \frac{8\gamma\eta\tilde{\mathcal{L}}_{\max}}{\alpha\theta_1 n} H^{k+1} \right] \\
&\stackrel{\text{Lemma 18}}{\leq} \beta Z^k + (1 - \theta_1 - \theta_2) \frac{2\gamma\beta}{\theta_1} Y^k + 2\gamma\beta \frac{\theta_2}{\theta_1} W^k + \frac{\gamma\eta}{\theta_1} \mathbb{E} [\|g^k - \nabla f(x^k)\|^2] \\
&\quad - \frac{\gamma}{4n\theta_1} \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 - \frac{\gamma}{8n\theta_1} \sum_{i=1}^n \|\nabla f_i(y^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \\
&\quad + \mathbb{E} \left[ 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1 q} W^{k+1} + \frac{8\gamma\eta\tilde{\mathcal{L}}_{\max}}{\alpha\theta_1 n} H^{k+1} \right] \\
&\stackrel{\text{Lemma 19}}{=} \beta Z^k + (1 - \theta_1 - \theta_2) \frac{2\gamma\beta}{\theta_1} Y^k + 2\gamma\beta \frac{\theta_2}{\theta_1} W^k + \frac{\gamma\eta}{\theta_1} \mathbb{E} [\|g^k - \nabla f(x^k)\|^2] \\
&\quad - \frac{\gamma}{4n\theta_1} \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 - \frac{\gamma}{8n\theta_1} \sum_{i=1}^n \|\nabla f_i(y^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \\
&\quad + 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1 q} (1-q) W^k + 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1} Y^k + \mathbb{E} \left[ \frac{8\gamma\eta\tilde{\mathcal{L}}_{\max}}{\alpha\theta_1 n} H^{k+1} \right] \\
&\leq \beta Z^k + \left(1 - \frac{\theta_1}{2}\right) \frac{2\gamma\beta}{\theta_1} Y^k + \left(1 - \frac{\theta_1 q}{2}\right) 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1 q} W^k \\
&\quad - \frac{\gamma}{4n\theta_1} \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 - \frac{\gamma}{8n\theta_1} \sum_{i=1}^n \|\nabla f_i(y^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \\
&\quad + \frac{\gamma\eta}{\theta_1} \mathbb{E} [\|g^k - \nabla f(x^k)\|^2] + \mathbb{E} \left[ \frac{8\gamma\eta\tilde{\mathcal{L}}_{\max}}{\alpha\theta_1 n} H^{k+1} \right] \\
&\stackrel{\text{Lemma 20}}{\leq} \beta Z^k + \left(1 - \frac{\theta_1}{2}\right) \frac{2\gamma\beta}{\theta_1} Y^k + \left(1 - \frac{\theta_1 q}{2}\right) 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1 q} W^k \\
&\quad - \frac{\gamma}{4n\theta_1} \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 - \frac{\gamma}{8n\theta_1} \sum_{i=1}^n \|\nabla f_i(y^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \\
&\quad + \frac{2\gamma\eta\tilde{\mathcal{L}}_{\max}}{\theta_1 n^2} \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 + \frac{2\gamma\eta\tilde{\mathcal{L}}_{\max}}{\theta_1 n} H^k + \mathbb{E} \left[ \frac{8\gamma\eta\tilde{\mathcal{L}}_{\max}}{\alpha\theta_1 n} H^{k+1} \right] \\
&\stackrel{\text{Lemma 21}}{\leq} \beta Z^k + \left(1 - \frac{\theta_1}{2}\right) \frac{2\gamma\beta}{\theta_1} Y^k + \left(1 - \frac{\theta_1 q}{2}\right) 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1 q} W^k \\
&\quad - \frac{\gamma}{4n\theta_1} \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 - \frac{\gamma}{8n\theta_1} \sum_{i=1}^n \|\nabla f_i(y^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \\
&\quad + \frac{2\gamma\eta\tilde{\mathcal{L}}_{\max}}{\theta_1 n^2} \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 + \frac{2\gamma\eta\tilde{\mathcal{L}}_{\max}}{\theta_1 n} H^k + \frac{8\gamma\eta\tilde{\mathcal{L}}_{\max}}{\alpha\theta_1 n} \left(1 - \frac{\alpha}{2}\right) H^k \\
&\quad + \left(1 + \frac{2q}{\alpha}\right) \frac{16\gamma\eta\tilde{\mathcal{L}}_{\max}q}{\alpha\theta_1 n^2} \left( \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 + \sum_{i=1}^n \|\nabla f_i(y^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \right) \\
&= \beta Z^k + \left(1 - \frac{\theta_1}{2}\right) \frac{2\gamma\beta}{\theta_1} Y^k + \left(1 - \frac{\theta_1 q}{2}\right) 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1 q} W^k + \left(1 - \frac{\alpha}{4}\right) \frac{8\gamma\eta\tilde{\mathcal{L}}_{\max}}{\alpha\theta_1 n} H^k \\
&\quad - \frac{\gamma}{n\theta_1} \left( \frac{1}{8} - \frac{2\eta\tilde{\mathcal{L}}_{\max}}{n} \right) \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \\
&\quad - \frac{\gamma}{n\theta_1} \left( \frac{1}{8} - \left(1 + \frac{2q}{\alpha}\right) \frac{16\eta\tilde{\mathcal{L}}_{\max}q}{\alpha n} \right) \left( \sum_{i=1}^n \|\nabla f_i(w^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 + \sum_{i=1}^n \|\nabla f_i(y^k) - \nabla f_i(x^k)\|_{\mathbf{L}_i^\dagger}^2 \right).
\end{aligned}$$

To make the last two lines disappear from the recurrence, we need to make sure

$$\frac{1}{8} - \frac{2\eta\tilde{\mathcal{L}}_{\max}}{n} \geq 0 \quad \text{and} \quad \frac{1}{8} - \left(1 + \frac{2q}{\alpha}\right) \frac{16\eta\tilde{\mathcal{L}}_{\max}q}{\alpha n} \geq 0,$$

or equivalently

$$\eta \leq \frac{n}{16\tilde{\mathcal{L}}_{\max}} \quad \text{and} \quad \eta \leq \frac{n}{64\tilde{\mathcal{L}}_{\max}} \cdot \frac{1}{\frac{2q}{\alpha} \left(\frac{2q}{\alpha} + 1\right)}.$$

Since  $\alpha \leq \frac{1}{\omega_{\max}+1}$  (see Lemma 21) and we also need to have  $\eta \leq \frac{1}{2L}$  (see Lemma 18), we can set

$$\eta = \min \left( \frac{1}{2L}, \frac{n}{64\tilde{\mathcal{L}}_{\max} (2q(\omega_{\max} + 1) + 1)^2} \right).$$

Therefore

$$\begin{aligned} \mathbb{E} [\Psi^{k+1}] &\leq \beta Z^k + \left(1 - \frac{\theta_1}{2}\right) \frac{2\gamma\beta}{\theta_1} Y^k + \left(1 - \frac{\theta_1 q}{2}\right) 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1 q} W^k + \left(1 - \frac{\alpha}{4}\right) \frac{8\gamma\eta\tilde{\mathcal{L}}_{\max}}{\alpha\theta_1 n} H^k \\ &\leq \left(1 - \frac{\eta\mu}{4\theta_1}\right) Z^k + \left(1 - \frac{\theta_1}{2}\right) \frac{2\gamma\beta}{\theta_1} Y^k + \left(1 - \frac{\theta_1 q}{2}\right) 2\gamma\beta \frac{\theta_2(1+\theta_1)}{\theta_1 q} W^k + \left(1 - \frac{\alpha}{4}\right) \frac{8\gamma\eta\tilde{\mathcal{L}}_{\max}}{\alpha\theta_1 n} H^k \\ &\leq \left(1 - \min \left\{ \frac{\alpha}{4}, \frac{q}{8}, \frac{\sqrt{\eta\mu q}}{4} \right\}\right) \Psi^k, \end{aligned}$$

where we set  $\gamma = \frac{\eta}{2(\theta_1 + \eta\mu)}$ ,  $\beta = 1 - \gamma\mu \leq 1 - \frac{\eta\mu}{4\theta_1}$  due to  $\eta\mu \leq \theta_1$ , and  $\theta_1 = \min \left\{ \frac{1}{4}, \sqrt{\frac{\eta\mu}{q}} \right\}$ .

After telescoping we get an  $\varepsilon$ -solution  $\mathbb{E} [\|z^k - x^*\|^2] \leq \varepsilon$  after

$$\max \left( 4(1 + \omega_{\max}), \frac{8}{q}, 4\sqrt{\frac{2}{\mu q} \max \left( L, \frac{32\tilde{\mathcal{L}}_{\max} (2q(\omega_{\max} + 1) + 1)^2}{n} \right)} \right) \log \frac{\Psi^0}{\varepsilon}$$

iterations. Choosing  $q = \min \left\{ 1, \frac{\max \left( 1, \sqrt{\frac{nL}{32\tilde{\mathcal{L}}_{\max}}} - 1 \right)}{2(1 + \omega_{\max})} \right\}$  we can simplify the above iteration complexity into

$$k = \begin{cases} \tilde{\mathcal{O}} \left( \omega_{\max} + \sqrt{\frac{\tilde{\mathcal{L}}_{\max}(1 + \omega_{\max})}{\mu n}} \right) & \text{if } nL \leq 128\tilde{\mathcal{L}}_{\max} \\ \tilde{\mathcal{O}} \left( 1 + \omega_{\max} + \sqrt{\frac{1 + \omega_{\max}}{\sqrt{n}} \frac{\sqrt{\tilde{\mathcal{L}}_{\max} L}}{\mu}} \right) & \text{if } 128\tilde{\mathcal{L}}_{\max} < nL \leq 32\tilde{\mathcal{L}}_{\max}(2\omega_{\max} + 3)^2 \\ \tilde{\mathcal{O}} \left( \omega_{\max} + \sqrt{\frac{L}{\mu}} \right) & \text{if } 32\tilde{\mathcal{L}}_{\max}(2\omega_{\max} + 3)^2 < nL. \end{cases}$$

Combining last two cases concludes the proof.  $\square$

## H Improvements Over The Original Methods

In this part we provide detailed derivations skipped in Section 5. Recall parameters  $\nu, \nu_s$  describing matrices  $\mathbf{L}_i$ :

$$\nu := \frac{\sum_{i=1}^n L_i}{\max_{1 \leq i \leq n} L_i}, \quad \nu_s := \max_{1 \leq i \leq n} \frac{\sum_{j=1}^d \mathbf{L}_{i,j}^{1/s}}{\max_{1 \leq j \leq d} \mathbf{L}_{i,j}^{1/s}}, \quad (49)$$

where  $L_i = \lambda_{\max}(\mathbf{L}_i)$  and we will choose  $s = 1$  or  $s = 2$ . Let  $L_{\max} = \max_{1 \leq i \leq n} L_i$ .

### H.1 Importance sampling for DCGD+

Let  $\tau = \mathbb{E}[|S_i|] = \sum_{j=1}^d p_{i,j}$  be the expected mini-batch size for the samplings  $S_i$ . Notice that convergence rate of DCGD+ depends on  $\tilde{\mathcal{L}}_{\max} = \max_{1 \leq i \leq n} \tilde{\mathcal{L}}_i$ . Since each node  $i \in [n]$  generates its own diagonal sketch  $\mathbf{C}_i$  independently from others, each node can optimize  $\tilde{\mathcal{L}}_i = \lambda_{\max}(\tilde{\mathbf{P}}_i \circ \mathbf{L}_i)$  independently based on local smoothness matrix  $\mathbf{L}_i$ . In general, minimizing  $\lambda_{\max}(\tilde{\mathbf{P}}_i \circ \mathbf{L}_i)$  with respect to probability matrix  $\tilde{\mathbf{P}}_i$  is hard. However, we can find the optimal probabilities when each node generates via an independent sampling, namely  $p_{i,jl} = p_{i,j}p_{i,l}$  if  $j \neq l$ . Then

$$\lambda_{\max}(\tilde{\mathbf{P}}_i \circ \mathbf{L}_i) = \max_{1 \leq j \leq d} \left( \frac{1}{p_{i,j}} - 1 \right) \mathbf{L}_{i,j}, \quad (50)$$

for which we can find the optimal probabilities  $p_{i,j}$ . To minimize the maximum term in (50), we should have  $(1/p_{i,j} - 1) \mathbf{L}_{i,j} = \rho_i$  for some  $\rho_i \geq 0$ . Then the solution is

$$p_{i,j} = \frac{\mathbf{L}_{i,j}}{\mathbf{L}_{i,j} + \rho_i}, \quad (51)$$

where  $\rho_i \geq 0$  is the unique solution to  $\sum_{j=1}^d \frac{\mathbf{L}_{i,j}}{\mathbf{L}_{i,j} + \rho_i} = \tau$ . The latter does not allow closed form solution for  $\rho_i$ , but it can be computed numerically using one dimensional solvers. Hence, we can efficiently compute the optimal probabilities (51). Moreover, we can deduce a simple upper bound for  $\rho_i$

$$\tau = \sum_{j=1}^d \frac{\mathbf{L}_{i,j}}{\mathbf{L}_{i,j} + \rho_i} \leq \sum_{j=1}^d \frac{\mathbf{L}_{i,j}}{\rho_i} = \frac{1}{\rho_i} \sum_{j=1}^d \mathbf{L}_{i,j}, \quad (52)$$

which gives us an upper bound for  $\tilde{\mathcal{L}}_i$  as follows

$$\tilde{\mathcal{L}}_i = \lambda_{\max}(\tilde{\mathbf{P}}_i \circ \mathbf{L}_i) = \rho_i \leq \frac{1}{\tau} \sum_{j=1}^d \mathbf{L}_{i,j} \stackrel{(49)}{\leq} \frac{\nu_1}{\tau} \mathbf{L}_{\max}. \quad (53)$$

*Proof of Remark 3.* Using the following inequalities with respect to matrix order

$$\mathbf{L} \preceq \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i, \quad \mathbf{L}_i \preceq n\mathbf{L}, \quad (54)$$

we bound  $L$  as follows

$$L = \lambda_{\max}(\mathbf{L}) \stackrel{(54)}{\leq} \lambda_{\max}\left(\frac{1}{n} \sum_{i=1}^n \mathbf{L}_i\right) \leq \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) = \frac{1}{n} \sum_{i=1}^n L_i \stackrel{(49)}{\leq} \frac{\nu}{n} L_{\max}. \quad (55)$$

Fix  $\tau = \sum_{j=1}^d p_{i,j} \in [0, d]$  expected mini-batch of coordinates for all nodes  $i \in [n]$ . Then, with probabilities (51) we have

$$\frac{\tilde{\mathcal{L}}_{\max}}{n} = \frac{1}{n} \max_{1 \leq i \leq n} \tilde{\mathcal{L}}_i = \frac{1}{n} \max_{1 \leq i \leq n} \rho_i \stackrel{(53)}{\leq} \frac{\nu_1}{\tau n} \mathbf{L}_{\max} \leq \frac{\nu_1}{\tau n} L_{\max},$$

To get it upper bounded by  $L_{\max}$ , notice that  $\max_{1 \leq j \leq d} \mathbf{L}_{i;j} \leq \lambda_{\max}(\mathbf{L}_i) = L_i$ , which implies

$$\mathbf{L}_{\max} = \max_{1 \leq i \leq n} \max_{1 \leq j \leq d} \mathbf{L}_{i;j} \leq \max_{1 \leq i \leq n} L_i = L_{\max}. \quad (56)$$

Therefore

$$L + \frac{\tilde{\mathcal{L}}_{\max}}{n} \leq \left( \frac{\nu}{n} + \frac{\nu_1}{\tau n} \right) L_{\max}.$$

□

**Remark 7** (Speedup for uniform sampling). *For standard sparsification with uniform probabilities, the term affected by the compression in the complexity (consider the linear rate of DCGD for simplicity) is  $\omega_{\max} L_{\max} = \left( \frac{d}{\tau} - 1 \right) L_{\max}$ , where  $L_{\max} = \max_{i \in [n]} L_i$  is the largest smoothness constant over devices. On the other hand, in the proposed sparsification strategy we have probabilities  $p_{i;jl} = \frac{\tau^2}{d^2}$  if  $j \neq l$ , and  $p_{i;jl} = \frac{\tau}{d}$  if  $j = l$ , which implies that  $\tilde{\mathbf{P}}_i = \left( \frac{d}{\tau} - 1 \right) \mathbf{I}$ . In this case, the term affected by the compression in the complexity is*

$$\tilde{\mathcal{L}}_{\max} = \max_{i \in [n]} \lambda_{\max}(\tilde{\mathbf{P}}_i \circ \mathbf{L}_i) = \left( \frac{d}{\tau} - 1 \right) \lambda_{\max}(\text{Diag}(\mathbf{L}_i)) = \left( \frac{d}{\tau} - 1 \right) \mathbf{L}_{\max},$$

where  $\mathbf{L}_{\max} = \max_{i \in [n]} \max_{j \in [d]} \mathbf{L}_{i;j}$  is the largest diagonal element over all smoothness matrices. Now notice that  $\mathbf{L}_{\max} \leq L_{\max} \leq d \mathbf{L}_{\max}$  hold and bounds are tight. Hence, the upper bound obtained for our sparsification is always better and can be up to  $d$  times better depending on the ratio  $\frac{L_{\max}}{\mathbf{L}_{\max}} \in [1, d]$ . Thus, we can make an analogous observation between classical uniform sampling and our uniform sampling albeit with a different condition on smoothness matrices, i.e.  $\frac{L_{\max}}{\mathbf{L}_{\max}} = \Omega(d)$  instead of  $\nu_s = \mathcal{O}(1)$ .

## H.2 Importance sampling for DIANA+

To find optimal probabilities for DIANA+, we minimize  $\omega_{\max} + \frac{\tilde{\mathcal{L}}_{\max}}{\mu n}$  part of the complexity (11) when each node uses an independent sampling as for DCGD+. Definitions of  $\tilde{\mathcal{L}}_{\max}$  and  $\omega_{\max}$  imply

$$\omega_{\max} + \frac{\tilde{\mathcal{L}}_{\max}}{\mu n} = \max_{ij} \left( \frac{1}{p_{i;j}} - 1 \right) + \max_{ij} \left( \frac{1}{p_{i;j}} - 1 \right) \frac{\mathbf{L}_{i;j}}{\mu n} = \Theta \left( \max_{ij} \left( \frac{1}{p_{i;j}} - 1 \right) \left( \frac{\mathbf{L}_{i;j}}{\mu n} + 1 \right) \right). \quad (57)$$

Therefore it is equivalent to minimize the following for each node  $i \in [n]$  independently:

$$\max_{1 \leq j \leq d} \left( \frac{1}{p_{i;j}} - 1 \right) \mathbf{L}'_{i;j}, \quad \mathbf{L}'_{i;j} := \frac{\mathbf{L}_{i;j}}{\mu n} + 1 \geq 1, \quad (58)$$

This can be solved in the same way as (50). The optimal probabilities are

$$p_{i;j} = \frac{\mathbf{L}'_{i;j}}{\mathbf{L}'_{i;j} + \rho'_i} = \frac{\frac{\mathbf{L}_{i;j}}{\mu n} + 1}{\frac{\mathbf{L}_{i;j}}{\mu n} + 1 + \rho'_i} \quad (59)$$

and an upper bound for  $\rho'_i$  is analogous to (53)

$$\rho'_i \leq \frac{1}{\tau} \sum_{j=1}^d \mathbf{L}'_{i;j} = \frac{1}{\tau} \sum_{j=1}^d \left( \frac{\mathbf{L}_{i;j}}{\mu n} + 1 \right) = \frac{d}{\tau} + \frac{1}{n\tau} \sum_{j=1}^d \frac{\mathbf{L}_{i;j}}{\mu} \stackrel{(49)}{\leq} \frac{d}{\tau} + \frac{\nu_1}{n\tau} \frac{\mathbf{L}_{\max}}{\mu} \stackrel{(56)}{\leq} \frac{d}{\tau} + \frac{\nu_1}{n\tau} \frac{L_{\max}}{\mu}. \quad (60)$$

*Proof of Remark 4.* With probabilities (59) we can upper bound the complexity (11) as follows

$$\begin{aligned} \omega_{\max} + \frac{\tilde{\mathcal{L}}_{\max}}{\mu n} &\stackrel{(57)}{\leq} 2 \max_{1 \leq i \leq n} \max_{1 \leq j \leq d} \left( \frac{1}{p_{i;j}} - 1 \right) \mathbf{L}'_{i;j} \\ &\stackrel{(59)}{=} \frac{2}{\tau} \max_{1 \leq i \leq n} \rho'_i \\ &\stackrel{(60)}{\leq} \frac{2d}{\tau} + \frac{2\nu_1}{\tau n} \frac{L_{\max}}{\mu}. \end{aligned} \quad (61)$$

Combined with (55), we have

$$\omega_{\max} + \frac{L}{\mu} + \frac{\tilde{\mathcal{L}}_{\max}}{\mu n} \leq \frac{2d}{\tau} + \left( \frac{\nu}{n} + \frac{2\nu_1}{\tau n} \right) \frac{L_{\max}}{\mu}.$$

□

**Remark 8** (Improvement over standard DGD). *Let us estimate how much improvement do we get with respect to standard Distributed Gradient Descent (DGD), where each node computes full gradients  $\nabla f_i(x^k)$  and sends dense updates to the server in each iteration. The iteration complexity of DGD is  $\tilde{\mathcal{O}}(\frac{L}{\mu})$ . To compare it against the complexity (11) of DIANA+ we use the same setup as in previous remarks (namely, independent samplings with probabilities (18) and  $\tau = d/n$ ). Since  $\mathbf{L}_i \preceq n\mathbf{L}$ , we have  $L_{\max} = \max_{i \in [n]} \lambda_{\max}(\mathbf{L}_i) \leq nL$ . Hence, (19) implies*

$$\omega_{\max} + \frac{L}{\mu} + \frac{\tilde{\mathcal{L}}_{\max}}{\mu n} \leq 2n + \frac{3nL}{\mu},$$

which is  $\mathcal{O}(n)$  times bigger than the iteration complexity of DGD. However, in case of DGD, each node sends  $n$  times more bits to the server. In total, DIANA+ and DGD have the same communication complexity in the worst case. To illustrate the best complexity DIANA+ can provide, consider the special case when  $\mathbf{L}_i = \mathbf{L}$  for all  $i \in [n]$  and  $\nu_1 = \mathcal{O}(1)$ . Then, clearly  $L_{\max} = L$  and we get  $\tilde{\mathcal{O}}(n + \frac{L}{\mu})$  complexity for DIANA+, yielding up to  $n$  times speedup against DGD. Moreover, in case of diagonal matrices  $\mathbf{L}_i$ , DIANA+ spends  $n$  times less local computation on partial derivatives and guarantees additional  $n$  times speedup.

### H.3 Independent sampling for ADIANA+

For the accelerated method ADIANA+, we construct probabilities  $p_{i,j}$  similar to (51) and (59) as follows

$$p_{i,j} := \left( \frac{\mathbf{L}'_{i,j}}{\mathbf{L}'_{i,j} + \rho''_i} \right)^{1/2} = \left( \frac{\frac{\mathbf{L}_{i,j}}{\mu n} + 1}{\frac{\mathbf{L}_{i,j}}{\mu n} + 1 + \rho''_i} \right)^{1/2}, \quad \mathbf{L}'_{i,j} = \frac{\mathbf{L}_{i,j}}{\mu n} + 1 \geq 1, \quad (62)$$

where  $\rho''_i$  is determined uniquely from  $\sum_{j=1}^d \left( \frac{\mathbf{L}'_{i,j}}{\mathbf{L}'_{i,j} + \rho''_i} \right)^{1/2} = \tau$ . Notice that

$$\tau = \sum_{j=1}^d \left( \frac{\mathbf{L}'_{i,j}}{\mathbf{L}'_{i,j} + \rho''_i} \right)^{1/2} \leq \sum_{j=1}^d \left( \frac{\mathbf{L}'_{i,j}}{\rho''_i} \right)^{1/2} = \frac{1}{\sqrt{\rho''_i}} \sum_{j=1}^d \sqrt{\mathbf{L}'_{i,j}}.$$

Therefore

$$\begin{aligned} \sqrt{\rho''_i} &\leq \frac{1}{\tau} \sum_{j=1}^d \sqrt{\frac{\mathbf{L}_{i,j}}{\mu n} + 1} \leq \frac{1}{\tau} \sum_{j=1}^d \left( \sqrt{\frac{\mathbf{L}_{i,j}}{\mu n}} + 1 \right) \leq \frac{d}{\tau} + \frac{1}{\tau} \sum_{j=1}^d \sqrt{\frac{\mathbf{L}_{i,j}}{\mu n}} \\ &\stackrel{(49)}{\leq} \frac{d}{\tau} + \frac{\nu_2}{\tau} \sqrt{\frac{\mathbf{L}_{\max}}{\mu n}} \stackrel{(56)}{\leq} \frac{d}{\tau} + \frac{\nu_2}{\tau} \sqrt{\frac{L_{\max}}{\mu n}} \end{aligned} \quad (63)$$

*Proof of Remark 5.* We bound terms  $\omega_{\max}$  and  $\frac{\mathcal{L}_{\max}}{\mu n}$  using probabilities (62) as follows:

$$\omega_{\max} = \max_{i,j} \left( \frac{1}{p_{i,j}} - 1 \right) = \max_{i,j} \left( \sqrt{\frac{\rho''_i}{\mathbf{L}'_{i,j}}} + 1 - 1 \right) \leq \max_{i,j} \sqrt{\frac{\rho''_i}{\mathbf{L}'_{i,j}}} \stackrel{(62)}{\leq} \max_{i,j} \sqrt{\rho''_i} \stackrel{(63)}{\leq} \frac{d}{\tau} + \frac{\nu_2}{\tau} \sqrt{\frac{L_{\max}}{\mu n}}. \quad (64)$$

$$\frac{\mathcal{L}_{\max}}{\mu n} \stackrel{(50)}{=} \max_{i,j} \left( \frac{1}{p_{i,j}} - 1 \right) \frac{\mathbf{L}_{i,j}}{\mu n} \stackrel{(62)}{\leq} \max_{i,j} \frac{\sqrt{\rho''_i} \frac{\mathbf{L}_{i,j}}{\mu n}}{\sqrt{\frac{\mathbf{L}_{i,j}}{\mu n}} + 1} \leq \max_{i,j} \sqrt{\rho''_i} \sqrt{\frac{\mathbf{L}_{i,j}}{\mu n}} \stackrel{(63)}{\leq} \left( \frac{d}{\tau} + \frac{\nu_2}{\tau} \sqrt{\frac{L_{\max}}{\mu n}} \right) \sqrt{\frac{L_{\max}}{\mu n}}. \quad (65)$$

Let  $\nu$  and  $\nu_2$  are  $\mathcal{O}(1)$ . Denote  $\omega = \frac{d}{\tau}$ ,  $\kappa_i = \frac{L_i}{\mu}$  and  $\kappa_{\max} = \max_{i \in [n]} \kappa_i$ . Then with this notation we have

$$\begin{aligned} \frac{L}{\mu} &\leq \frac{\nu}{n} \kappa_{\max} = \mathcal{O}\left(\frac{\kappa_{\max}}{n}\right) \\ \omega_{\max} &\leq \omega + \frac{\nu_2}{\tau} \sqrt{\frac{\kappa_{\max}}{n}} = \omega \left(1 + \frac{\nu_2}{d} \sqrt{\frac{\kappa_{\max}}{n}}\right) = \mathcal{O}\left(\omega \left(1 + \frac{\sqrt{\kappa_{\max}}}{d\sqrt{n}}\right)\right) \\ \frac{\mathcal{L}_{\max}}{\mu n} &\leq \left(\omega + \frac{\nu_2}{\tau} \sqrt{\frac{\kappa_{\max}}{n}}\right) \sqrt{\frac{\kappa_{\max}}{n}} = \mathcal{O}\left(\omega \left(1 + \frac{\sqrt{\kappa_{\max}}}{d\sqrt{n}}\right) \frac{\sqrt{\kappa_{\max}}}{\sqrt{n}}\right) \end{aligned} \quad (66)$$

Then, in case of  $nL \leq \tilde{\mathcal{L}}_{\max}$ , we have

$$\omega_{\max} + \sqrt{\omega_{\max} \frac{\tilde{\mathcal{L}}_{\max}}{\mu n}} = \mathcal{O}\left(\omega \left(1 + \frac{\sqrt{\kappa_{\max}}}{d\sqrt{n}}\right) \left(1 + \left(\frac{\kappa_{\max}}{n}\right)^{1/4}\right)\right),$$

which should be compared with  $\mathcal{O}\left(\omega \left(1 + \sqrt{\frac{\kappa_{\max}}{n}}\right)\right)$  [Li et al., 2020]. If  $\kappa_{\max} = \mathcal{O}(nd^2)$ , then we get  $\mathcal{O}(\sqrt{d})$  speedup factor. If  $nL > \tilde{\mathcal{L}}_{\max}$ , then

$$\begin{aligned} \omega_{\max} + \sqrt{\frac{L}{\mu}} + \sqrt{\omega_{\max} \sqrt{\frac{\tilde{\mathcal{L}}_{\max}}{\mu n}} \sqrt{\frac{L}{\mu}}} \\ = \mathcal{O}\left(\omega \left(1 + \frac{\sqrt{\kappa_{\max}}}{d\sqrt{n}}\right) + \sqrt{\frac{\kappa_{\max}}{n}} + \sqrt{\omega \left(1 + \frac{\sqrt{\kappa_{\max}}}{d\sqrt{n}}\right) \sqrt{\frac{\kappa_{\max}}{n}} \sqrt{\omega \left(1 + \frac{\sqrt{\kappa_{\max}}}{d\sqrt{n}}\right) \sqrt{\frac{\kappa_{\max}}{n}}}}\right) \\ = \mathcal{O}\left(\omega \left(1 + \frac{\sqrt{\kappa_{\max}}}{d\sqrt{n}}\right) + \sqrt{\frac{\kappa_{\max}}{n}} + \left[\omega \left(1 + \frac{\sqrt{\kappa_{\max}}}{d\sqrt{n}}\right) \sqrt{\frac{\kappa_{\max}}{n}}\right]^{3/4}\right), \end{aligned}$$

which should be compared with  $\omega + \kappa_{\max} + \omega^{3/4} n^{1/4} \sqrt{\frac{\kappa_{\max}}{n}}$  [Li et al., 2020]. If  $\kappa_{\max} = \mathcal{O}(nd^2)$ , then we get  $\mathcal{O}(\sqrt{n})$  times smaller second term and  $\mathcal{O}((nd)^{1/4})$  times smaller third term.  $\square$

## I Variance Reduction: ISEGA+

In this part we apply our redesign to another variance reduced method called ISEGA [Mishchenko et al., 2020, Hanzely and Richtárik, 2019b]. At the core of ISEGA, the mechanism for variance reduction is based on SEGAs method [Hanzely et al., 2018]. The key difference between ISEGA and DIANA is that ISEGA updates the control variates  $h$  more aggressively using projection instead of the mere  $\alpha$ -step towards the projection used in DIANA. Adapting our matrix-smoothness-aware sparsification to ISEGA, we define the update rule of control vectors  $h_i^k$  as follows (**for now assume  $\mathbf{L}_i$  is invertible**)

$$\begin{aligned} h_i^{k+1} &= \arg \min_{h \in \text{Range}(\mathbf{L}_i)} \|h - h_i^k\|_{\mathbf{L}_i^\dagger}^2 \\ &\quad \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^k) = \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} h \\ &= h_i^k + \mathbf{L}_i \mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i^k \left( \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i \mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i^k \right)^\dagger \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k) \\ &= h_i^k + \mathbf{L}_i^{1/2} \mathbf{C}_i^k (\mathbf{C}_i^k \mathbf{C}_i^k)^\dagger \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k) \\ &= h_i^k + \mathbf{L}_i^{1/2} \mathbf{Diag}(\mathbf{P}_i) \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k). \end{aligned}$$

Note that the update rule in DIANA+ has the form

$$h_i^{k+1} = h_i^k + \alpha \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k)$$

for some fixed scalar  $\alpha > 0$ , and thus is more conservative. Note that we choose the gradient estimator to be the same  $g_i^k = h_i^k + \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k)$ . The method is presented as Algorithm 7.

---

### Algorithm 7 ISEGA+

---

- 1: **Input:** Initial point  $x^0 \in \mathbb{R}^d$ , initial shifts  $h_i^0 \in \mathbb{R}^d$ , current point  $x^k$ , step size parameter  $\gamma$  and  $\alpha$ , sketch  $\mathbf{C}_i^k$  and  $\bar{\mathbf{C}}_i^k := \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2}$ , current shifts  $h_1^k, \dots, h_n^k$  and  $h^k := \frac{1}{n} \sum_{i=1}^n h_i^k$ .
  - 2: **on** each node
  - 3:   get  $x^k$  from the server
  - 4:   send sparse update  $\Delta_i^k = \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k)$
  - 5:    $g_i^k = h_i^k + \mathbf{L}_i^{1/2} \Delta_i^k$
  - 6:    $h_i^{k+1} = h_i^k + \mathbf{L}_i^{1/2} \mathbf{Diag}(\mathbf{P}_i) \Delta_i^k$
  - 7: **on** server
  - 8:   get sparse updates  $\Delta_i^k$  from each node
  - 9:    $g^k = \frac{1}{n} \sum_{i=1}^n g_i^k = h^k + \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i^{1/2} \Delta_i^k$
  - 10:    $x^{k+1} = \text{prox}_{\gamma R}(x^k - \gamma g^k)$
  - 11:    $h^{k+1} = \frac{1}{n} \sum_{i=1}^n h_i^{k+1} = h^k + \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i^{1/2} \mathbf{Diag}(\mathbf{P}_i) \Delta_i^k$
- 

Note that we can not obtain the convergence rate of ISEGA+ directly from the framework of Gorbunov et al. [2020a]. Instead, to get the tight convergence rate, we shall cast it as an instance of GJS method [Hanzely and Richtárik, 2019b]. Theorem 22 provides the result – we can see that the worst case complexity is identical to DIANA+. In terms of the practical performance, we expect ISEGA+ to outperform DIANA+ due to the more aggressive update rule of control variates.

**Theorem 22.** Suppose that  $\gamma \leq \frac{1}{\frac{4\tilde{\mathcal{L}}_{\max}}{n} + 2L + \mu(\omega_{\max} + 1)}$ . Then, we have

$$\mathbb{E}[\Psi^k] \leq (1 - \gamma\mu)\Psi^0,$$

where

$$\Psi^k := \|x^k - x^*\|^2 + \frac{\gamma}{2n} \sum_{i=1}^n \|\phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*)\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1}}^2$$

and  $\phi_i^k := \mathbf{L}_i^{\dagger 1/2} h_i^k$ . Consequently, the overall complexity of ISEGA+ is

$$\tilde{\mathcal{O}} \left( \omega_{\max} + \frac{L}{\mu} + \frac{\tilde{\mathcal{L}}_{\max}}{n\mu} \right).$$



*Proof.* The proof can be seen as a special case of the generalized Jacobian sketching theory of Hanzely and Richtárik [2019b]. For the sake of clarity, we provide a specialized proof here.

Note first that by (46), we have

$$\mathbb{E} [\|g^k - \nabla f(x^*)\|^2] \leq 2 \left( L + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) D_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|^2.$$

Similarly, we have

$$\begin{aligned} & \mathbb{E} \left[ \left\| \phi_i^{k+1} - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1}}^2 \right] \\ &= \mathbb{E} \left[ \left\| \phi_i^k + \mathbf{Diag}(\mathbf{P}_i) \mathbf{C}_i^k (\mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^k) - \phi_i^k) - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1}}^2 \right] \\ &= \mathbb{E} \left[ \left\| (\mathbf{I} - \mathbf{Diag}(\mathbf{P}_i) \mathbf{C}_i^k) (\phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*)) + \mathbf{Diag}(\mathbf{P}_i) \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - \nabla f_i(x^*)) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1}}^2 \right] \\ &= \mathbb{E} \left[ \left\| (\mathbf{I} - \mathbf{Diag}(\mathbf{P}_i) \mathbf{C}_i^k) \mathbf{Diag}(\mathbf{P}_i)^{-\frac{1}{2}} (\phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*)) + \mathbf{Diag}(\mathbf{P}_i)^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - \nabla f_i(x^*)) \right\|^2 \right] \\ &= \mathbb{E} \left[ \left\| (\mathbf{I} - \mathbf{Diag}(\mathbf{P}_i) \mathbf{C}_i^k) \mathbf{Diag}(\mathbf{P}_i)^{-\frac{1}{2}} (\phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*)) \right\|^2 \right] + \mathbb{E} \left[ \left\| \mathbf{Diag}(\mathbf{P}_i)^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - \nabla f_i(x^*)) \right\|^2 \right] \\ &= \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1} - \mathbf{I}}^2 + \left\| \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - \nabla f_i(x^*)) \right\|^2 \\ &\leq \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1} - \mathbf{I}}^2 + 2D_{f_i}(x^k, x^*) \end{aligned}$$

and therefore

$$\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \left\| \phi_i^{k+1} - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1}}^2 \right] \leq \frac{1}{n} \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1} - \mathbf{I}}^2 + 2D_f(x^k, x^*) \quad (67)$$

Following the classical analysis of SGD (i.e., proof of Lemma C.1 of Gorbunov et al. [2020a]), we get

$$\begin{aligned} \mathbb{E} [\|x^{k+1} - x^*\|^2] &= (1 - \gamma\mu) \|x^k - x^*\|^2 - 2\gamma D_f(x^k, x^*) + \gamma^2 \mathbb{E} [\|g^k - \nabla f(x^*)\|^2] \\ &\leq (1 - \gamma\mu) \|x^k - x^*\|^2 - 2\gamma \left( 1 - \gamma \left( L + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) \right) D_f(x^k, x^*) \\ &\quad + \frac{2\tilde{\mathcal{L}}_{\max}\gamma^2}{n^2} \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|^2. \end{aligned}$$

Adding  $\frac{\gamma}{2}$ -multiple of (67) to the above, we get

$$\begin{aligned} & \mathbb{E} [\|x^{k+1} - x^*\|^2] + \frac{\gamma}{2} \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \left\| \phi_i^{k+1} - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1}}^2 \right] \\ &\leq (1 - \gamma\mu) \|x^k - x^*\|^2 - 2\gamma \left( \frac{1}{2} - \gamma \left( L + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) \right) D_f(x^k, x^*) \\ &\quad + \frac{2\tilde{\mathcal{L}}_{\max}\gamma^2}{n^2} \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|^2 + \frac{\gamma}{2n} \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1}}^2 \quad (68) \end{aligned}$$

Next, note that we have

$$\begin{aligned} & \frac{2\tilde{\mathcal{L}}_{\max}\gamma^2}{n^2} \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|^2 + \frac{\gamma}{2n} \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1} - \mathbf{I}}^2 \\ &\leq \frac{(1 - \gamma\mu)\gamma}{2n} \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1}}^2 \quad (69) \end{aligned}$$

since it is equivalent to

$$\frac{4\tilde{\mathcal{L}}_{\max}\gamma}{n} \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|^2 + \gamma \mu \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|_{\mathbf{Diag}(\mathbf{P}_i)^{-1}}^2 \leq \sum_{i=1}^n \left\| \phi_i^k - \mathbf{L}_i^{\dagger 1/2} \nabla f_i(x^*) \right\|^2,$$

which holds since  $\gamma \leq \frac{1}{\frac{4\tilde{\mathcal{L}}_{\max}}{n} + \mu(\omega_{\max} + 1)}$ .

To finish the proof, it remains to plug (69) into (68), use that  $\gamma \leq \frac{1}{\frac{4\tilde{\mathcal{L}}_{\max}}{n} + 2L}$  and unroll the recurrence.

□

## J Variance Reduction with Bi-directional Compression: DIANA++

In this method, the master server applies compression in its turn with sketch  $\mathbf{C}$  independently. Thus, we maintain an additional control vector  $H^k$ , which helps to reduce the variance coming from the master's sparsification. Moreover, nodes keep track of  $H^k$  just like the central server.

---

### Algorithm 8 DIANA++

---

- 1: **Input:** Initial point  $x^0 \in \mathbb{R}^d$ , initial shifts  $h_i^0 \in \text{Range}(\mathbf{L}_i)$ ,  $H^0 \in \text{Range}(\mathbf{L})$ , current point  $x^k$ , step size parameter  $\gamma, \alpha$  and  $\beta$ , sketch  $\mathbf{C}_i^k$  and  $\bar{\mathbf{C}}_i^k := \mathbf{L}_i^{1/2} \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2}$ , current shifts  $h_1^k, \dots, h_n^k, H^k$  and  $h^k := \frac{1}{n} \sum_{i=1}^n h_i^k$ .
  - 2: **on** each node
  - 3:   **send** sparse update  $\Delta_i^k = \mathbf{C}_i^k \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k)$
  - 4:    $\bar{\Delta}_i^k = \mathbf{L}_i^{1/2} \Delta_i^k$ ,  $g_i^k = h_i^k + \bar{\Delta}_i^k$ ,  $h_i^{k+1} = h_i^k + \alpha \bar{\Delta}_i^k$
  - 5: **on** server
  - 6:   **get** sparse updates  $\Delta_i^k$  from each node
  - 7:    $\bar{\Delta}^k = \frac{1}{n} \sum_{i=1}^n \bar{\Delta}_i^k = \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i^{1/2} \Delta_i^k$
  - 8:    $g^k = \bar{\Delta}^k + h^k = \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{C}}_i^k (\nabla f_i(x^k) - h_i^k) + h^k$
  - 9:   **send** sparse update  $\delta^k = \mathbf{C}^k \mathbf{L}^{\dagger 1/2} (g^k - H^k)$
  - 10:    $\bar{\delta}^k = \mathbf{L}^{1/2} \delta^k$ ,  $\hat{g}^k = H^k + \bar{\delta}^k = H^k + \bar{\mathbf{C}}^k (g^k - H^k)$
  - 11:    $x^{k+1} = \text{prox}_{\gamma R}(x^k - \gamma \hat{g}^k)$
  - 12:    $h^{k+1} = h^k + \alpha \bar{\Delta}^k$
  - 13:    $H^{k+1} = H^k + \beta \bar{\delta}^k$
  - 14: **on** each node
  - 15:   **get**  $\delta^k$  from the server
  - 16:   reconstruct  $\bar{\delta}^k = \mathbf{L}^{1/2} \delta^k$ ,  $\hat{g}^k = H^k + \bar{\delta}^k = H^k + \bar{\mathbf{C}}^k (g^k - H^k)$
  - 17:    $x^{k+1} = \text{prox}_{\gamma R}(x^k - \gamma \hat{g}^k)$
  - 18:    $H^{k+1} = H^k + \beta \bar{\delta}^k$
- 

**Theorem 23.** *Let Assumptions 1 and 2 hold and assume that each node generates its own diagonal sketch  $\mathbf{C}_i$  independently from others. The master server, in its turn, generates  $\mathbf{C}$  independently from the nodes. Then, Algorithm 8 has the following iteration complexity*

$$\mathcal{O} \left( \frac{1}{\min(\alpha - \beta\theta', \beta)} + \frac{\alpha + \beta\theta + \beta\theta'}{\min(\alpha - \beta\theta', \beta)} \left( \frac{L}{\mu} + \frac{\tilde{\mathcal{L}}}{\mu} + \frac{\tilde{\mathcal{L}}'_{\max}}{n\mu} + \frac{\tilde{\mathcal{L}}_{\max}}{n\mu} \right) \right),$$

where we made the following notations

$$\theta := \frac{n\tilde{\mathcal{L}}}{\tilde{\mathcal{L}}_{\max} + 2\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}} \leq \frac{n}{2\tilde{\mathcal{L}}'_{\max}}, \quad \theta' := \frac{2\theta}{n} \tilde{\mathcal{L}}'_{\max} \leq 1 \in [0, 1]$$

$$\tilde{\mathcal{L}}'_{\max} := \max_{1 \leq i \leq n} \lambda_{\max} \left( \tilde{\mathbf{P}}_i \circ (\mathbf{L}_i^{1/2} \mathbf{L}^{\dagger} \mathbf{L}_i^{1/2}) \right), \quad \tilde{\mathcal{L}} := \lambda_{\max} \left( \tilde{\mathbf{P}} \circ \mathbf{L} \right)$$

with bounds  $\alpha \leq \frac{1}{1+\omega_{\max}} = \max_{i \in [n]} \max_{j \in [d]} \frac{1}{p_{i,j}}$  and  $\beta \leq \frac{1}{1+\omega} = \max_{j \in [d]} \frac{1}{p_j}$ .

**Remark 9.** *Note that, when master does not compress the messages, then we have  $\tilde{\mathbf{P}} = \mathbf{0}$ . This implies the same complexity we had for DIANA+ as quantities  $\tilde{\mathcal{L}}$ ,  $\theta$ ,  $\theta'$  are all become zeros.*

*Proof.* The proof follows the same structure as for DIANA+, with additional variance reduction process introduced for the master server. Analogously, we start bounding the following second moment:

$$\mathbb{E} [\|\hat{g}^k - \nabla f(x^*)\|^2] = \mathbb{E} [\|\hat{g}^k - g^k\|^2] + \mathbb{E} [\|g^k - \nabla f(x^*)\|^2]. \quad (70)$$

We can bound the second term as it was done in (46):

$$\mathbb{E} [\|g^k - \nabla f(x^*)\|^2] \leq 2 \left( L + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) D_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^{\dagger}}^2.$$

Then we decompose the first term  $\mathbb{E} [\|\hat{g}^k - g^k\|^2]$  into two as follows:

$$\begin{aligned}
\mathbb{E} [\|\hat{g}^k - g^k\|^2] &= \mathbb{E} [\|\bar{\mathbf{C}}^k(g^k - H^k) - (g^k - H^k)\|^2] \\
&= \|g^k - H^k\|_{\mathbb{E}[(\mathbf{I} - \bar{\mathbf{C}}^k)^\top (\mathbf{I} - \bar{\mathbf{C}}^k)]}^2 \\
&= \|g^k - H^k\|_{\mathbf{L}^\dagger^{1/2}(\tilde{\mathbf{P}} \circ \mathbf{L})\mathbf{L}^\dagger^{1/2}}^2 \\
&\leq \tilde{\mathcal{L}} \|g^k - H^k\|_{\mathbf{L}^\dagger}^2 \\
&\leq 2\tilde{\mathcal{L}} \|g^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + 2\tilde{\mathcal{L}} \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2.
\end{aligned} \tag{71}$$

To bound each of the two summands in (71), we derive the analogue of (39).

$$\begin{aligned}
&\mathbb{E} [\mathbf{L}_i^{1/2} (\bar{\mathbf{C}}_i - \mathbf{I})^\top \mathbf{L}^\dagger (\bar{\mathbf{C}}_i - \mathbf{I}) \mathbf{L}_i^{1/2}] \\
&= \mathbb{E} [\mathbf{L}_i^{1/2} (\mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i \mathbf{L}_i^{1/2} - \mathbf{I}) \mathbf{L}^\dagger (\mathbf{L}_i^{1/2} \mathbf{C}_i \mathbf{L}_i^{\dagger 1/2} - \mathbf{I}) \mathbf{L}_i^{1/2}] \\
&= \mathbb{E} [\mathbf{L}_i^{1/2} (\mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i (\mathbf{L}_i^{1/2} \mathbf{L}^\dagger \mathbf{L}_i^{1/2}) \mathbf{C}_i \mathbf{L}_i^{\dagger 1/2} - \mathbf{L}_i^{\dagger 1/2} \mathbf{C}_i \mathbf{L}_i^{1/2} \mathbf{L}^\dagger - \mathbf{L}^\dagger \mathbf{L}_i^{1/2} \mathbf{C}_i \mathbf{L}_i^{\dagger 1/2} + \mathbf{L}^\dagger) \mathbf{L}_i^{1/2}] \\
&\stackrel{(36)}{=} \mathbf{L}_i^{1/2} (\mathbf{L}_i^{\dagger 1/2} (\bar{\mathbf{P}}_i \circ (\mathbf{L}_i^{1/2} \mathbf{L}^\dagger \mathbf{L}_i^{1/2})) \mathbf{L}_i^{\dagger 1/2} - \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} \mathbf{L}^\dagger - \mathbf{L}^\dagger \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} + \mathbf{L}^\dagger) \mathbf{L}_i^{1/2} \\
&= \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\bar{\mathbf{P}}_i \circ (\mathbf{L}_i^{1/2} \mathbf{L}^\dagger \mathbf{L}_i^{1/2})) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2} - \mathbf{L}_i^{1/2} \mathbf{L}^\dagger \mathbf{L}_i^{1/2} \\
&= \mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\tilde{\mathbf{P}}_i \circ (\mathbf{L}_i^{1/2} \mathbf{L}^\dagger \mathbf{L}_i^{1/2})) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2}.
\end{aligned} \tag{72}$$

Then we bound them as follows. First, we have

$$\begin{aligned}
\mathbb{E} [\|g^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2] &= \|\nabla f(x^k) - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + \mathbb{E} [\|g^k - \nabla f(x^k)\|_{\mathbf{L}^\dagger}^2] \\
&\leq 2D_f(x^k, x^*) + \mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{C}}_i^k (\nabla f_i(x^k) - h_i^k) + h_i^k - \nabla f_i(x^k) \right\|_{\mathbf{L}^\dagger}^2 \right] \\
&= 2D_f(x^k, x^*) + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \left\| (\bar{\mathbf{C}}_i^k - \mathbf{I}) \mathbf{L}_i^{1/2} r_i^k \right\|_{\mathbf{L}^\dagger}^2 \right] \\
&= 2D_f(x^k, x^*) + \frac{1}{n^2} \sum_{i=1}^n \|r_i^k\|_{\mathbf{L}_i^{1/2}(\bar{\mathbf{C}}_i^k - \mathbf{I})^\top \mathbf{L}^\dagger (\bar{\mathbf{C}}_i^k - \mathbf{I}) \mathbf{L}_i^{1/2}}^2 \\
&\stackrel{(72)}{=} 2D_f(x^k, x^*) + \frac{1}{n^2} \sum_{i=1}^n \|r_i^k\|_{\mathbf{L}_i^{1/2} \mathbf{L}_i^{\dagger 1/2} (\tilde{\mathbf{P}}_i \circ (\mathbf{L}_i^{1/2} \mathbf{L}^\dagger \mathbf{L}_i^{1/2})) \mathbf{L}_i^{\dagger 1/2} \mathbf{L}_i^{1/2}}^2 \\
&= 2D_f(x^k, x^*) + \frac{1}{n^2} \sum_{i=1}^n \left\| \mathbf{L}_i^{\dagger 1/2} (\nabla f_i(x^k) - h_i^k) \right\|_{\tilde{\mathbf{P}}_i \circ (\mathbf{L}_i^{1/2} \mathbf{L}^\dagger \mathbf{L}_i^{1/2})}^2 \\
&\leq 2D_f(x^k, x^*) + \frac{\tilde{\mathcal{L}}'_{\max}}{n^2} \sum_{i=1}^n \|\nabla f_i(x^k) - h_i^k\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq 2D_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}'_{\max}}{n^2} \sum_{i=1}^n \|\nabla f_i(x^k) - f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + \frac{2\tilde{\mathcal{L}}'_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \\
&\leq 2D_f(x^k, x^*) + \frac{4\tilde{\mathcal{L}}'_{\max}}{n} D_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}'_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \\
&= 2 \left( 1 + \frac{2\tilde{\mathcal{L}}'_{\max}}{n} \right) D_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}'_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2
\end{aligned} \tag{73}$$

Then, for the control vectors  $H^k$  at the master, we have

$$\begin{aligned}
& \mathbb{E}_k \left[ \|H^{k+1} - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 \right] \\
&= \mathbb{E}_k \left[ \|H^k - \nabla f(x^*) + \beta \bar{\delta}^k\|_{\mathbf{L}^\dagger}^2 \right] \\
&= \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + 2\beta \mathbb{E} [\langle H^k - \nabla f(x^*), g^k - H^k \rangle_{\mathbf{L}^\dagger}] + \beta^2 \mathbb{E}_k \left[ \|\bar{\mathbf{C}}^k (g^k - H^k)\|_{\mathbf{L}^\dagger}^2 \right] \\
&= \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + 2\beta \mathbb{E}_k [\langle H^k - \nabla f(x^*), g^k - H^k \rangle_{\mathbf{L}^\dagger}] + \beta^2 \mathbb{E}_k \left[ \|g^k - H^k\|_{\mathbb{E}[(\bar{\mathbf{C}}^k)^\top \mathbf{L}^\dagger \bar{\mathbf{C}}^k]}^2 \right] \\
&\leq \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + 2\beta \mathbb{E}_k [\langle H^k - \nabla f(x^*), g^k - H^k \rangle_{\mathbf{L}^\dagger}] + \beta^2 \mathbb{E}_k \left[ \|g^k - H^k\|_{\mathbf{L}^\dagger^{1/2} \mathbb{E}[(\mathbf{C}^k)^2] \mathbf{L}^\dagger^{1/2}}^2 \right] \\
&\leq \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + 2\beta \mathbb{E}_k [\langle H^k - \nabla f(x^*), g^k - H^k \rangle_{\mathbf{L}^\dagger}] + \beta^2 (1 + \omega) \mathbb{E}_k \left[ \|g^k - H^k\|_{\mathbf{L}^\dagger}^2 \right] \\
&\leq \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + 2\beta \mathbb{E}_k [\langle H^k - \nabla f(x^*), g^k - H^k \rangle_{\mathbf{L}^\dagger}] + \beta \mathbb{E}_k \left[ \|g^k - H^k\|_{\mathbf{L}^\dagger}^2 \right] \\
&= (1 - \beta) \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + \beta \mathbb{E}_k \left[ \|g^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 \right] \\
&\leq (1 - \beta) \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + 2\beta \left( 1 + \frac{2\tilde{\mathcal{L}}'_{\max}}{n} \right) D_f(x^k, x^*) + \frac{2\beta\tilde{\mathcal{L}}'_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2
\end{aligned}$$

Now, for some  $\theta$  (to be defined later), let

$$\sigma^k := \frac{1}{n} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + \theta \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2.$$

Then, we have

$$\begin{aligned}
& \mathbb{E} [\|\hat{g}^k - \nabla f(x^*)\|^2] \\
&\stackrel{(70)}{=} \mathbb{E} [\|\hat{g}^k - g^k\|^2] + \mathbb{E} [\|g^k - \nabla f(x^*)\|^2] \\
&\stackrel{(71)}{\leq} 2\tilde{\mathcal{L}} \mathbb{E} [\|g^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2] + 2\tilde{\mathcal{L}} \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + \mathbb{E} [\|g^k - \nabla f(x^*)\|^2] \\
&\stackrel{(73)}{\leq} 4\tilde{\mathcal{L}} \left( 1 + \frac{2\tilde{\mathcal{L}}'_{\max}}{n} \right) D_f(x^k, x^*) + \frac{4\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \\
&\quad + 2 \left( L + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) D_f(x^k, x^*) + \frac{2\tilde{\mathcal{L}}_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \\
&\quad + 2\tilde{\mathcal{L}} \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 \\
&= 2 \left( L + 2\tilde{\mathcal{L}} + \frac{4\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}}{n} + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) D_f(x^k, x^*) \\
&\quad + \left( \frac{4\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}}{n} + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) \frac{1}{n} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + 2\tilde{\mathcal{L}} \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 \\
&= 2 \left( L + 2\tilde{\mathcal{L}} + \frac{4\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}}{n} + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) D_f(x^k, x^*) + \left( \frac{4\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}}{n} + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \right) \sigma^k,
\end{aligned}$$

with the following choice of  $\theta$ :

$$\theta := \frac{n\tilde{\mathcal{L}}}{\tilde{\mathcal{L}}_{\max} + 2\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}} \leq \frac{n}{2\tilde{\mathcal{L}}'_{\max}}, \quad \theta' := \frac{2\theta}{n} \tilde{\mathcal{L}}'_{\max} \leq 1.$$

For the control vectors  $h_i^k$  and  $H^k$ , we deduce

$$\begin{aligned}
& \mathbb{E} [\sigma^{k+1}] \\
& \leq (1 - \alpha) \frac{1}{n} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + 2\alpha D_f(x^k, x^*) \\
& \quad + (1 - \beta)\theta \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 + 2\beta\theta \left(1 + \frac{2\tilde{\mathcal{L}}'_{\max}}{n}\right) D_f(x^k, x^*) + \frac{2\beta\theta\tilde{\mathcal{L}}'_{\max}}{n^2} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 \\
& = \left(1 - \alpha + \frac{2\beta\theta\tilde{\mathcal{L}}'_{\max}}{n}\right) \frac{1}{n} \sum_{i=1}^n \|h_i^k - \nabla f_i(x^*)\|_{\mathbf{L}_i^\dagger}^2 + (1 - \beta)\theta \|H^k - \nabla f(x^*)\|_{\mathbf{L}^\dagger}^2 \\
& \quad + 2 \left(\alpha + \beta\theta \left(1 + \frac{2\tilde{\mathcal{L}}'_{\max}}{n}\right)\right) D_f(x^k, x^*) \\
& \leq \max \left(1 - \alpha + \frac{2\beta\theta\tilde{\mathcal{L}}'_{\max}}{n}, 1 - \beta\right) \sigma^k + 2 \left(\alpha + \beta\theta \left(1 + \frac{2\tilde{\mathcal{L}}'_{\max}}{n}\right)\right) D_f(x^k, x^*) \\
& = \max(1 - \alpha + \beta\theta', 1 - \beta) \sigma^k + 2(\alpha + \beta\theta + \beta\theta') D_f(x^k, x^*).
\end{aligned}$$

Thus the constants from [Gorbunov et al., 2020a] are as follows

$$\begin{aligned}
A &= L + 2\tilde{\mathcal{L}} + \frac{4\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}}{n} + \frac{2\tilde{\mathcal{L}}_{\max}}{n} \\
B &= \frac{4\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}}{n} + \frac{2\tilde{\mathcal{L}}_{\max}}{n} = \frac{2\tilde{\mathcal{L}}}{\theta} \\
C &= \alpha + \beta\theta + \beta\theta' \\
\rho &= \min(\alpha - \beta\theta', \beta).
\end{aligned}$$

Let  $M = \frac{2B}{\rho}$ , and note that  $B\theta = 2\tilde{\mathcal{L}}$  and  $B\theta' = \frac{4\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}}{n}$ . Then

$$\begin{aligned}
A + CM &= A + 2B \frac{\alpha + \beta\theta + \beta\theta'}{\min(\alpha - \beta\theta', \beta)} \\
&= \mathcal{O} \left( \frac{\alpha + \beta\theta + \beta\theta'}{\min(\alpha - \beta\theta', \beta)} \left( L + \tilde{\mathcal{L}} + \frac{\tilde{\mathcal{L}}\tilde{\mathcal{L}}'_{\max}}{n} + \frac{\tilde{\mathcal{L}}_{\max}}{n} \right) \right). \\
1 + \frac{B}{M} - \rho &= 1 - \frac{\rho}{2} = 1 - \frac{1}{2} \min(\alpha - \beta\theta', \beta).
\end{aligned}$$

□