PARALLEL STOCHASTIC ASYNCHRONOUS COORDINATE DESCENT: TIGHT BOUNDS ON THE POSSIBLE PARALLELISM*

YUN KUEN CHEUNG[†], RICHARD COLE[‡], AND YIXIN TAO[§]

Abstract. Several works have shown linear speedup is achieved by an asynchronous parallel implementation of stochastic coordinate descent so long as there is not too much parallelism. More specifically, it is known that if all updates are of similar duration, then linear speedup is possible with up to $\Theta(L_{\max}\sqrt{n}/L_{\overline{\text{tes}}})$ processors, where L_{\max} and $L_{\overline{\text{tes}}}$ are suitable Lipschitz parameters. This paper shows the bound is tight for almost all possible values of these parameters.

Key words. stochastic asynchronous coordinate descent, parallelism bound

AMS subject classifications. 90C25, 68W10, 60G50, 68Q99

 1. Introduction. Very large scale optimization problems have arisen in many areas such as machine learning. A natural approach for solving these huge problems is to employ parallel and more specifically asynchronous parallel algorithms. As common wisdom suggests, when a small number of processors are used in these algorithms, (linear) speedup can be achieved; but when too many processors are involved and when they are not properly coordinated, there may be undesirable outcomes. (Recall that speedup is defined as the ratio of the algorithm's execution time on a single core and its parallel execution time. Linear speedup means the speedup is linear in the number of cores.)¹ Typically, the bound on the number of processors is implicit and is expressed in terms of the maximum number q of basic iterations, namely single coordinate updates, that can overlap. In many scenarios q will be a small multiple of the number of processors or cores. In many earlier works the notation τ was used instead of q.²

This paper considers asynchronous implementations of stochastic coordinate descent (SCD) applied to smooth convex functions $f: \mathbb{R}^n \to \mathbb{R}$. Several recent works [3, 2, 1] quantify how large q can be in this context while guaranteeing linear speedup.³ More precisely, they showed: if at any time at most $q \leq \tilde{q}$ updates can overlap, then linear speedup is guaranteed. The goal in these works was to demonstrate as large a value of \tilde{q} as possible. Note that these results provide *lower* bounds on the actual value of \tilde{q} .

^{*}Submitted to the editors DATE. Authors are listed in alphabetical order.

Funding: Yun Kuen Cheung would like to acknowledge Singapore NRF 2018 Fellowship NRF-NRFF2018-07 and MOE AcRF Tier 2 Grant 2016-T2-1-170. The work of Richard Cole and Yixin Tao was supported in part by NSF Grants CCF-1527568 and CCF-1909538.

[†]Royal Holloway University of London (yunkuen.cheung@rhul.ac.uk, http://cs.rhul.ac.uk/~cheung/) Part of the work done while this author was a postdoctoral fellow at Max-Planck Institute for Informatics, Saarland Informatics Campus and Singapore University of Technology and Design, and also during two visits to the Courant Institute, NYU in the summers of 2017 and 2018.

[‡]Courant Institute, NYU (cole@cs.nyu.edu, https://cs.nyu.edu/cole/)

[§]Courant Institute, NYU (yt851@nyu.edu, https://tomtao26.github.io/). Yixin Tao is now a postdoc at LSE.

¹In optimization, we compare the times of two executions that achieve a given level of accuracy. ²We chose the notation q to emphasize its likely similarity to p, the number of processors.

³Many of these results hold for *composite* functions, i.e., functions of the form $f(\mathbf{x}) = g(\mathbf{x}) + \sum_{k=1}^{n} \Psi_k(x_k)$, where $g: \mathbb{R}^n \to \mathbb{R}$ is a convex function with a continuous gradient, and each $\Psi_k: \mathbb{R} \to \mathbb{R}$ is a univariate convex function, but may be non-smooth.

32

36

38

39

40

41

42

43

44

45

The best existing lower bound is $\tilde{q} = \Omega(\sqrt{n}L_{\text{max}}/L_{\overline{\text{res}}})$, where L_{max} and $L_{\overline{\text{res}}}$ are Lipschitz parameters defined in Section 2.⁴ Intuitively, one can view this as a lower bound on the possible parallelism supporting linear speedup; for if there are p processors at hand, and if the durations of the updates vary by at most a factor of d, then $q \leq (d+1)(p-1)$, so achieving linear speedup with up to q updates overlapping implies linear speedup occurs with p = 1 + q/(d+1) processors.

We present the first work concerning the inverse problem: to identify a value \overline{q} , such that if $q \geq \overline{q}$, then this can lead to an undesirable outcome.

Main result: The lower bound of $\Omega(\sqrt{n}L_{\max}/L_{\overline{\text{res}}})$ is asymptotically tight. (An undesirable outcome can occur if $\bar{q} \geq \bar{c} \cdot \sqrt{n}L_{\max}/L_{\overline{\text{res}}}$ for a sufficiently large constant \bar{c} .)

We will present an adversarial family of functions with the following property: if q exceeds $\Theta(\sqrt{n}L_{\max}/L_{\overline{res}})$ significantly, then there is an asynchronous schedule for which with high probability very rapid divergence occurs for a very long time. The function family uses a dimensionless parameter $\epsilon = \Theta(\frac{L_{\overline{res}}}{L_{\max}\sqrt{n}})$, which is approximately the inverse of the possible parallelism.

We use the following function family, $f_{\epsilon}: \mathbb{R}^n \to \mathbb{R}$:

$$f_{\epsilon}(\mathbf{x}) = \frac{1-\epsilon}{2} \cdot \sum_{i=1}^{n} (x_i)^2 + \frac{\epsilon}{2} \cdot \left(\sum_{i=1}^{n} x_i\right)^2,$$

for any ϵ satisfying $4/n \le \epsilon < 1$. As we shall see, $L_{\rm max} = 1$ and $L_{\overline{\rm res}} = \Omega(\epsilon \sqrt{n})$ for this function family; thus the existing lower bound on the parallelism achieving linear speedup is $\Omega(\sqrt{n}L_{\rm max}/L_{\overline{\rm res}}) = \Omega(1/\epsilon)$. To obtain bounds using arbitrary values of $L_{\rm max}$ one can simply multiply f_{ϵ} by $L_{\rm max}$, which also increases $L_{\overline{\rm res}}$ by an $L_{\rm max}$ factor.

Next, we discuss more precisely how we achieve this result. Recall that the performance of a sequential SCD algorithm is expressed in terms of its convergence rate. On strongly convex functions, it has a linear convergence rate, meaning that each update reduces the expected value of the difference $f(\mathbf{x}) - f^*$ by at least an $(1 - \alpha/n)$ multiplicative factor, for some constant $\alpha > 0$, where f^* denotes the minimum value of the function. Consequently,

$$\mathbb{E}\left[f(\mathbf{x}^t) - f^*\right] \leq \left(1 - \frac{\alpha}{n}\right)^t \cdot \left(f(\mathbf{x}^0) - f^*\right).$$

For our proposed function f_{ϵ} , which is strongly convex, we will show that for a suitable initial point \mathbf{x}^0 , for some constant $\alpha' \geq \alpha$,

$$\mathbb{E}\left[f_{\epsilon}(\mathbf{x}^{t}) - f_{\epsilon}^{*}\right] \geq \left(1 - \frac{\alpha'}{n}\right)^{t} \cdot \left(f_{\epsilon}(\mathbf{x}^{0}) - f_{\epsilon}^{*}\right),$$

and hence sequential SCD achieves no more than a linear convergence rate in general. For the function f_{ϵ} , we will show that $\alpha = \frac{1}{3}$ and $\alpha' = 2$.

To achieve linear speedup with a parallel algorithm means that the same convergence rate holds, up to constant factor, i.e., the α might be reduced by a constant factor $c \geq 1$, but no more:

66 (1.4)
$$\mathbb{E}\left[f(\mathbf{x}^t) - f^*\right] \leq \left(1 - \frac{\alpha}{cn}\right)^t \cdot \left(f(\mathbf{x}^0) - f^*\right),$$

⁴Here, we focus on the case where the step-size used in the asychronous SCD algorithm is $1/L_{\rm max}$; we will discuss the cases with smaller step-sizes later in the introduction.

where t is now the overall number of iterations performed by the various cores.⁵ This means that to guarantee a particular accuracy, the total number of iterations for a parallel execution is no more than a constant multiple of the number of iterations needed on a single core (so long as $\alpha \leq n/2$).

Prior work has shown that linear speedup is achieved when $q \leq \tilde{c}\sqrt{n}L_{\text{max}}/L_{\overline{\text{res}}}$ for some constant $\tilde{c} > 0$. To achieve our main result, we show that for the function family f_{ϵ} , when $q \geq \bar{c}\sqrt{n}L_{\text{max}}/L_{\overline{\text{res}}}$ for some constant $\bar{c} > \tilde{c}$, as an adversary, it is possible to pick asynchronous schedules such that for all $t \leq n^{10}$ (or more generally, for any constant $\hat{c} \geq 1$, for all $t \leq n^{\hat{c}}$),

$$\mathbb{E}\left[f_{\epsilon}(\mathbf{x}^{t}) - f_{\epsilon}^{*}\right] \geq \Omega(4^{t/q}) \cdot (f_{\epsilon}(\mathbf{x}^{0}) - f_{\epsilon}^{*})$$

for large enough n (in general, when $n = \Omega(c^4)$). This indicates that when q is too large, linear speedup cannot be achieved in worst-case scenarios.

The above upper bound on q holds when the step-size is $1/L_{\rm max}$. One might wonder what would happen if we reduced the step-size to $1/\Gamma$ for some $\Gamma \geq L_{\rm max}$. Would the permissible parallelism bound increase significantly, thereby improving the overall speedup? In fact, we show that the upper bound on q increases to at most $\mathcal{O}(\Gamma/L_{\rm res})$, for any $L_{\rm max} \leq \Gamma \leq \mathcal{O}(L_{\rm res}\sqrt{n})$. Since this increase is by a factor of at most $\mathcal{O}(\Gamma/L_{\rm max})$, but the step-size is reduced by a factor of $\Gamma/L_{\rm max}$, the overall speedup cannot improve by more than a constant factor. This upper bound is also asymptotically tight, as there were matching lower bounds for these choices of step-sizes [1].

Prior Work and Asynchrony Models. First, note that the bound on q is ensuring the asynchrony is bounded, and so we call it q-bounded asynchrony. Some requirement of this sort is unavoidable, otherwise there could be updates of arbitrarily long duration, which, when they commit, could undo an arbitrary amount of progress.

In addition to the q-bounded asynchrony assumption, we need to specify how the asynchronous environment affects the read operations. There are two models concerning how coordinates are read in asynchronous environments, namely "consistent" and "inconsistent" reads. Our upper bound applies to both models. Next, we discuss their differences.

Liu et al. [3] gave the first bound on the parallel performance of asynchronous SCD on convex functions, showing linear speedup when $q = O(\sqrt{n}L_{\text{max}}/L_{\text{res}})$ assuming a consistent read model, where L_{res} is another Lipschitz parameter defined in Section 2. We note that $L_{\text{res}} = L_{\overline{\text{res}}}$ for the function family f_{ϵ} we will be analyzing in this paper. In fact, as we shall see, $L_{\overline{\text{res}}}$ is equal to L_{res} on all quadratic functions f, i.e., f is of the form $\mathbf{x}^{\mathsf{T}}\mathbf{A}\mathbf{x} + \mathbf{b}^{\mathsf{T}}\mathbf{x} + \text{constant}$, where \mathbf{A} is an $n \times n$ matrix, and \mathbf{b} is an n-vector. In the consistent read model, all the coordinate values a processor reads when performing a single update on one coordinate may be out of date, but they must have been simultaneously current at some moment. To make this more precise, we view the updates as committing at integer times $t = 1, 2, \ldots$, and we write \mathbf{x}^t to be the value of \mathbf{x} after the update at time t. The consistent reads model requires the vector of \mathbf{x} values used by the time t update to be of the form $\mathbf{x}^{t-\tau}$ for some $\tau \geq 1$.

Consistent reads create a substantial constraint on the asynchrony, and so subsequent works sought to avoid this assumption. To this end, Liu and Wright [2] proposed

⁵Liu et al. [3] and Liu and Wright [2] call this near-linear speedup, reserving the term linear speedup for the case when c=1.

the inconsistent reads model. Allowing inconsistent reads means that the $\tilde{\mathbf{x}}$ values 111 used by the time t update can be any collection of the form $(x_1^{t-\tau_1}, \cdots, x_n^{t-\tau_n})$, where the τ_i 's can be distinct; the q-bounded asynchrony assumption implies that $1 \le \tau_i \le q$ 113 for each j. Liu and Wright showed that linear speedup (including the more general case of composite functions) can be achieved for $q = O(n^{1/4}\sqrt{L_{\text{max}}}/\sqrt{L_{\text{res}}})$, i.e., the 115 square root of the previous bound. There remained several constraints on the pos-116 sible asynchrony, in addition to the q-bounded asynchrony, as pointed out by Mania 117 et al. [4] and subsequently by Sun et al. [5]. The latter works also gave analyses 118 removing some or all of these constraints, but at the cost of reducing the bound on 119 q. Finally, Cheung, Cole and Tao [1] gave an analysis achieving linear speedup for 120 $q = O(\sqrt{n}L_{\text{max}}/L_{\overline{\text{res}}})$, again for composite functions, with the only constraint being 121 122 the q-bounded asynchrony.

In this paper we show the bounds in [3] and [1] are asymptotically tight for almost all possible values of L_{max} , L_{res} , and $L_{\overline{\text{res}}}$ for the function family (1.1).

2. **Notation.** Let \mathbf{e}_j denote the *n*-vector in which the *j*-th entry is 1 and every other entry is 0.

DEFINITION 2.1. For any coordinates j, k, the function f is L_{jk} -Lipschitz-smooth if for any $\mathbf{x} \in \mathbb{R}^n$ and $r \in \mathbb{R}$, $|\nabla_k f(\mathbf{x} + r \cdot \mathbf{e}_j) - \nabla_k f(\mathbf{x})| \le L_{jk} \cdot |r|$; it is L_{res} -Lipschitz-smooth if, for all j, $||\nabla f(\mathbf{x} + r \cdot \mathbf{e}_j) - \nabla f(\mathbf{x})|| \le L_{\text{res}} \cdot |r|$. Finally, $L_{\text{max}} := \max_j L_{jj}$ and $L_{\overline{\text{res}}} := \max_k \left(\sum_{j=1}^n (L_{kj})^2\right)^{1/2}$.

Observe that for the function f_{ϵ} we are considering,

132 (2.1)
$$\nabla_j f_{\epsilon}(\mathbf{x}) = (1 - \epsilon) \cdot x_j + \epsilon \cdot \sum_{i=1}^n x_i.$$

Consequently, $L_{\text{max}} = 1$, and $L_{\overline{\text{res}}} = \sqrt{1 + (n-1)\epsilon^2} = \Theta(\epsilon\sqrt{n})$, for $\epsilon = \Omega(1/\sqrt{n})$.

134 The difference between $L_{\rm res}$ and $L_{\overline{\rm res}}$. In general, $L_{\overline{\rm res}} \geq L_{\rm res}$. $L_{\rm res} = L_{\overline{\rm res}}$ when 135 the rates of change of the gradient are constant, as for example in quadratic functions 136 such as $\mathbf{x}^{\mathsf{T}}\mathbf{A}\mathbf{x} + \mathbf{b}^{\mathsf{T}}\mathbf{x} + c$. We refer the reader to [1] for a discussion of why $L_{\overline{\rm res}}$ is 137 needed in general for the analysis in [1].

Next, we define strong convexity.

DEFINITION 2.2. Let $f: \mathbb{R}^n \to \mathbb{R}$ be a convex function. f is strongly convex with parameter $\mu > 0$, if for all x, y, $f(y) - f(x) \ge \langle \nabla f(x), y - x \rangle + \frac{1}{2}\mu||y - x||^2$.

A simple calculation shows that the parameter μ for our function f_{ϵ} has value (1 - ϵ); see Lemma A.1 in Appendix A.

143 The update rule. Recall that in a standard coordinate descent, be it sequential 144 or parallel and synchronous, the update rule, applied to coordinate j, first computes 145 the accurate gradient $\nabla_j f(\mathbf{x}^{t-1})$, and then performs the update given below.

$$x_j^t \leftarrow x_j^{t-1} - \frac{\nabla_j f(\mathbf{x}^{t-1})}{\Gamma}$$

and for all $k \neq j$, $x_k^t \leftarrow x_k^{t-1}$, where $\Gamma \geq L_{\text{max}}$ is a parameter controlling the step size.

However, in an asynchronous environment, an updating processor might retrieve 149 outdated information $\tilde{\mathbf{x}}$ instead of \mathbf{x}^{t-1} , so the gradient the processor computes will 150 be $\nabla_i f(\tilde{\mathbf{x}})$, instead of the accurate value $\nabla_i f(\tilde{\mathbf{x}}^{t-1})$. Hence the update rule is in the 151 asynchronous environment is

153 (2.3)
$$x_j^t \leftarrow x_j^{t-1} - \frac{\nabla_j f(\tilde{\mathbf{x}})}{\Gamma}.$$

154

156

157

158

159

160

161

162

163

165

167

168

169

170

171

172 173

174

175

179 180

We want to show that for any fixed constants $c_1, c_2 \geq 1, f_{\epsilon}(\mathbf{x}^t) - f_{\epsilon}^*$ is rapidly growing for $t \leq q \cdot n_1^c$ with probability at least $1 - 1/n^{c_2}$. 155

2.1. The Stochastic Asynchronous Coordinate Descent (SACD) Algo**rithm.** The coordinate descent process starts at an initial point $\mathbf{x}^0 = (x_1^0, x_2^0, \cdots, x_n^0)$. Multiple processors then iteratively update the coordinate values, and for our analysis we assume that at each time, there is exactly one coordinate value being updated, which we can do, as we are choosing the asynchronous schedule.

Algorithm 2.1 SACD Algorithm for Smooth Functions.

Input: The initial point $\mathbf{x}^0 = (x_1^0, x_2^0, \cdots, x_n^0)$.

Multiple processors use a shared memory. Each processor iteratively repeats the following four-step procedure, with no global coordination among them:

- **Step 1:** Choose a coordinate $j \in \{1, 2, \dots, n\}$ uniformly at random.
- **Step 2:** Retrieve coordinate values \tilde{x} from the shared memory.
- Step 3: Compute the gradient $\nabla_j f(\tilde{\mathbf{x}})$.
- **Step 4:** Update coordinate j using rule (2.3) by atomic addition.
- 2.2. Asynchrony Assumptions, Basic Set-up and Terminology Used in the Construction. In our construction, updates are made in phases. In each phase, q updates are made in parallel by q processors. As in Step 1 of Algorithm 2.1, each processor chooses a coordinate to update. Since these processors are uncoordinated, they choose coordinates randomly and independently, and thus it is possible that some of them choose the same coordinate within a phase. Then, for every coordinate, each processor retrieves the value that was up to date at the end of the previous phase.

Note that $\mathbf{x}^{q\ell}$ denote the up-to-date values immediately after the ℓ -th phase. The starting point \mathbf{x}^0 can be viewed as the up-to-date values following the (non-existent) 0-th phase. For each $\ell \geq 0$, in the $(\ell + 1)$ -st phase, each of the q processors performs the following sequence of steps:

- it picks a coordinate k randomly and independently to update;
- it retrieves $\mathbf{x}^{q\ell}$ and then computes $\nabla_k f(\mathbf{x}^{q\ell})$;
- it increases the value of x_k in the main memory by $-\nabla_k f(\mathbf{x}^{q\ell})/\Gamma$, using an atomic addition operation.

Note that in a single phase a particular coordinate might be chosen two or more 176 times. If coordinate k is chosen b times in the $(\ell+1)$ -st phase, the value of x_k after 177 the $(\ell + 1)$ -st phase is $x_k^{q\ell} - b \cdot \nabla_k f(\mathbf{x}^{q\ell}) / \Gamma$. 178

3. The Result. Cheung, Cole and Tao [1] showed the following upper bound on the performance of asynchronous stochastic coordinate descent on strongly convex functions, where \mathbf{x}^* is a minimum point for the convex function.

THEOREM 3.1. i. Suppose the asynchronous updating is run for exactly t updates. Also suppose that $q \leq \min\left\{\frac{\sqrt{n}}{270}, \frac{\Gamma\sqrt{n}}{270L_{\text{res}}}\right\}$. If f is strongly convex with parameter μ ,

184 then

185

189

190

191

192

193 194

195

196

197

198

200201

202

203

204

205

$$\mathbb{E}\left[f(\mathbf{x}^t) - f(\mathbf{x}^*)\right] \le \left(1 - \frac{1}{3n} \cdot \frac{\mu}{\Gamma}\right)^t \left(f(\mathbf{x}^0) - f(\mathbf{x}^*)\right).$$

186 ii. This result holds for all $q \leq \frac{\Gamma\sqrt{n}}{12L_{\rm res}}$ if the asynchronous schedule obeys the Strong Common Value assumption.

The Strong Common Value assumption, specified in [1], captures a substantial class of asynchronous schedules including the consistent read constraint from [3], and all the schedules to which the analysis of Liu and Wright [2] apply, and more. The complementary construction in this paper observes the consistent reads condition, which is the most restrictive of these constraints. Consequently, our goal is to show that for $q = \Omega(\frac{\Gamma\sqrt{n}}{L_{\rm res}})$ convergence with linear speed-up is not guaranteed. In fact, we will show that there is no convergence for most values of q obeying this constraint. To achieve this it suffices to show there is a single asynchronous schedule observing the consistent read condition for which there is no convergence.

First, in Appendix A, we will prove the following theorem, which shows that the sequential stochastic coordinate descent on our function f_{ϵ} when starting at the point $\mathbf{x}^{0} = (-1, +1, -1, +1, \ldots)$ achieves a convergence rate that is at most a constant factor faster than the convergence rate given in Theorem 3.1. Recall that $\mu = 1 - \epsilon$ for f_{ϵ} .

PROPOSITION 3.2. Let \mathbf{x}^0 be the point with even-indexed coordinates equal to +1 and odd-indexed coordinates equal to -1. Suppose the sequential stochastic coordinate descent is run for t steps on function f_{ϵ} starting at point x^0 . Then

$$\mathbb{E}\left[f_{\epsilon}(\mathbf{x}^{t}) - f_{\epsilon}(\mathbf{x}^{*})\right] \ge \left(1 - \frac{2}{n} \cdot \frac{1 - \epsilon}{\Gamma}\right)^{t} \left(f_{\epsilon}(\mathbf{x}^{0}) - f_{\epsilon}(\mathbf{x}^{*})\right).$$

This result shows that for the function f_{ϵ} , the speedup guaranteed by Theorem 3.1 is indeed a linear speedup of the performance of the sequential SCD.

Now, we come to our main results.

Our analyses look at phases of q successive updates when the coordinate descent

210 in Algorithm 4 is applied to the function $f \equiv L_{\text{max}} \cdot f_{\epsilon}$ in (1.1) with $\epsilon = \sqrt{\frac{\left(\frac{L_{\text{res}}}{L_{\text{max}}}\right)^2 - 1}{n-1}}$.

The first result states that the expected value of $f(\mathbf{x}) - f^*$ grows by at least a factor of 4 from phase to phase.

THEOREM 3.3. Suppose that $L_{\text{res}} \geq L_{\text{max}} (1+8/n)$ and $q \geq \frac{4\Gamma\sqrt{n}}{\sqrt{L_{\text{res}}^2 - L_{\text{max}}^2}}$. Then

there is an asynchronous schedule for which the coordinate descent diverges when applied to the function f. Specifically, for every t = rq, with $r \ge 1$ an integer,

216
$$\mathbb{E}\left[f(\mathbf{x}^{t}) - f^{*}\right] \geq 4^{(t/q)-1} \cdot (f(\mathbf{x}^{0}) - f^{*}).$$

217 If
$$L_{\rm res} \geq 2L_{\rm max}$$
 (i.e. $\epsilon \geq \sqrt{3/(n-1)}$), the constraint on q becomes $q \geq \frac{8\Gamma\sqrt{n}}{\sqrt{3}L_{\rm res}}$.

When $L_{\rm res} \geq 2L_{\rm max}$, Theorem 3.3 states that once q exceeds the bound in The-219 orem 3.1(ii) by a constant factor, in order to have any possibility of convergence, let alone linear convergence, Γ has to increase at the same rate as q. Note that the upper bound on the rate of convergence decreases linearly with Γ , so this says that increasing q, the number of parallel updates, beyond this bound cannot increase the parallel runtime by more than a constant factor at best.

However, one might wonder if this divergence in expectation is a low probability event. Our next result shows that for most values of q < n it is in fact a high probability event.

THEOREM 3.4. Let $c_1, c_2 \ge 1$ be constants, let $c = c_1 + c_2$, and suppose that

$$\frac{n}{e} \ge q \ge \max \left\{ \frac{8\Gamma}{\epsilon L_{\max}}, \frac{10(c+1)}{\epsilon} \frac{\log n}{\log \frac{n}{q}} \right\}.$$

218

220

223

224

225

226

227

237

239

240

243

247

248

For each $c_1, c_2 \geq 1$, with probability at least $1 - 1/n^{c_2}$, there is an asynchronous 229 schedule for which the coordinate descent diverges for at least $q \cdot n^{c_1}$ updates when 230 applied to the function f. Specifically, for every t = rq with $1 \le r \le n^{c_1}$ an integer, 231

232
$$f(\mathbf{x}^t) - f^* \ge 4^{(t/q)-1} \cdot (f(\mathbf{x}^0) - f^*).$$

233 If
$$L_{\rm res} \geq 2L_{\rm max}$$
, the conditions $\frac{n}{e} \geq q \geq \max\left\{\frac{16\Gamma\sqrt{n}}{\sqrt{3}L_{\rm res}}, \frac{200(c+1)L_{\rm max}\sqrt{n}}{\sqrt{3}L_{\rm res}}\right\}$ and $n \geq 234 \left(\frac{544(c+1)}{\sqrt{3}} \cdot \frac{L_{\rm max}}{L_{\rm res}}\right)^{5/2}$ suffice.

Thus the comments in the paragraph after Theorem 3.3 apply here too, so long as 235 236 the stated bounds on n hold.

Finally, one might wonder what happens to the value $f(\mathbf{x}) - f^*$ during a phase. Could it be small at any time? As it happens, with our schedule the value oscillates a lot, but the next and final result shows that with high probability it remains large at all times.

241 Theorem 3.5. Let $c_1, c_2 \geq 1$ be constants, let $c = c_1 + c_2$, and suppose that $L_{\text{res}} \ge 2L_{\text{max}}, \ n \ge 400c^2(c+2)^2, \ and$ 242

$$\frac{n}{e} \ge q \ge \max \left\{ \frac{8\Gamma \sqrt{n}}{\sqrt{3}L_{\rm res}} + \frac{4c\Gamma}{L_{\rm max}} \ , \ 20c(c+2) \right\}.$$

For each $c_1, c_2 \geq 1$, with probability at least $1 - 1/n^{c_2}$, there is an asynchronous 244 schedule for which the coordinate descent diverges for at least $q \cdot n^{c_1}$ updates when 245 applied to the function f. Specifically, for every $1 \cdot n^{c_1}$, 246

$$f(\mathbf{x}^t) - f^* \geq \frac{1}{n^2} \cdot 4^{\lceil (t/q) \rceil - 1} \cdot (f(\mathbf{x}^0) - f^*).$$

Once more, the comments in the paragraph after Theorem 3.3 apply.

The three theorems incorporate general choices of Γ , c_1 and c_2 . In Theorems 3.4 249 and 3.5, if we pick c_1, c_2 to be large enough constants (for example, both are 10 — 250 which corresponds to a $1/n^{10}$ failure probability over the course of $q \cdot n^{10}$ updates), the 251 conditions on n and q reduce to $n \ge \Theta(1)$ and $n/e \ge q \ge \Theta(\Gamma \sqrt{n}/L_{\rm res})$, respectively. 252

259

263

269270

4. Analysis. We consider running SACD with q processors on the function

$$f(\mathbf{x}) = L_{\text{max}} \cdot \left[\frac{1 - \epsilon}{2} \cdot \sum_{i=1}^{n} (x_i)^2 + \frac{\epsilon}{2} \cdot \left(\sum_{i=1}^{n} x_i \right)^2 \right],$$

- for any ϵ satisfying $4/n \le \epsilon < 1$. We prove each theorem in turn.
- We choose the initial point to be the all ones point, $\mathbf{x}^0 = (+1, +1, \dots, +1)$. For $\ell \geq 0$, let \mathbf{y}^{ℓ} denote the up-to-date values of the coordinates immediately after the ℓ -th phase, i.e. $\mathbf{y}^{\ell} = \mathbf{x}^{q\ell}$. Also, let

$$G^{\ell} \triangleq \left| \sum_{k=1}^{n} y_{k}^{\ell} \right|$$
 and $M^{\ell} \triangleq \max_{k=1,2,\cdots,n} |y_{k}^{\ell}|.$

- 4.1. Proof of Theorem 3.3. The key claim, implying exponential growth in the value of $f(\mathbf{x})$ at the end of each successive phase of q updates, is given by the following lemma.
 - LEMMA 4.1. If $\epsilon \geq \frac{4}{n}$ and $q \geq \frac{4\Gamma}{\epsilon L_{\max}}$ then for all ℓ , $\mathbb{E}\left[G^{\ell}\right] \geq 2^{\ell}G^{0}$.
- 264 *Proof.* The proof is by induction. The claim clearly holds for $\ell = 0$.
- Now suppose the result holds for some $\ell \geq 0$. Without loss of generality, we assume $\sum_{j=1}^{n} y_{j}^{\ell} > 0$; the other case will be symmetric. In the $(\ell + 1)$ -st phase, there are q updates. Each update picks a coordinate k; by Equation (2.1), the update reduces its value by

$$\frac{L_{\max}}{\Gamma} \cdot \left[(1 - \epsilon) y_k^{\ell} + \epsilon \sum_{j=1}^n y_j^{\ell} \right] = \frac{L_{\max}}{\Gamma} \cdot \left[(1 - \epsilon) y_k^{\ell} + \epsilon G^{\ell} \right].$$

- The expected reduction due to q updates is at least
- 272 $\frac{qL_{\max}}{\Gamma} \cdot \left[\epsilon G^{\ell} \frac{1}{n} (1 \epsilon) G^{\ell} \right] \ge \frac{3q\epsilon L_{\max}}{4\Gamma} \cdot G^{\ell} \quad (\text{as } n\epsilon \ge 4)$ $\ge 3G^{\ell} \qquad (\text{as } q \ge \frac{4\Gamma}{\epsilon L_{\max}}).$
- 275 Thus, $\mathbb{E}\left[G^{\ell+1}|G^{\ell}\right] \geq \left|\mathbb{E}\left[\sum_{k=1}^{n}y_{k}^{\ell+1}|G^{\ell}\right]\right| \geq (3-1)G^{\ell} = 2G^{\ell}$; therefore $\mathbb{E}\left[G^{\ell+1}\right] \geq 2$ 76 $2 \cdot \mathbb{E}\left[G^{\ell}\right]$, demonstrating the inductive claim.
- 277 Proof of Theorem 3.3. We begin by showing that $f(\mathbf{y}^0) f^* \leq \epsilon \cdot (G^0)^2$, assuming that $\epsilon \geq 4/n$, which we justify in the next paragraph. To see this, note that $G^0 = n$ and $f(\mathbf{y}^0) f^* = \frac{1}{2}(1-\epsilon)n + \frac{1}{2}\epsilon n^2$. As $\epsilon \geq 4/n$, we see that $f(\mathbf{y}^0) f^* \leq \frac{1}{8} \cdot \epsilon n^2 + \frac{1}{2}\epsilon n^2 \leq \epsilon (G^0)^2$. Next, immediately after Phase ℓ , $f(\mathbf{y}^\ell) f^* \geq \frac{\epsilon}{2} \cdot (G^\ell)^2$. Then, by the Cauchy-281 Schwarz inequality and Lemma 4.1, $\mathbb{E}\left[f(\mathbf{y}^\ell) f^*\right] \geq \mathbb{E}\left[\frac{\epsilon}{2} \cdot (G^\ell)^2\right] \geq \frac{\epsilon}{2} \cdot \mathbb{E}\left[G^\ell\right]^2 \geq 4^{\ell-1}\epsilon \cdot (G^0)^2 \geq 4^{\ell-1}(f(\mathbf{y}^0) f^*)$.
- It remains to show that $\epsilon \geq 4/n$. Recall that $\epsilon^2 L_{\rm max}^2 = (L_{\rm res}^2 L_{\rm max}^2)/(n-1)$. Together, these imply $(L_{\rm res}^2 L_{\rm max}^2) \geq 16(n-1)/n^2 \cdot L_{\rm max}^2$. It suffices to have $L_{\rm res} \geq 16(n-1)/n^2 \cdot L_{\rm max}^2$.
- $L_{\max}(1+8/n)$, proving the theorem.
- 4.2. Proof of Theorem 3.4. The idea of the construction is to show that if we can bound M^{ℓ} by $\frac{1}{4}\epsilon G^{\ell}$, then we can show the bound $G^{\ell} \geq 2^{\ell}G_0$ holds absolutely and not just in expectation. We will show that $M^{\ell} \leq \frac{1}{4}\epsilon G^{\ell}$ with high probability.

Our construction uses a parameter b which we will specify later. We will need that 289 290 in each phase, each coordinate is chosen at most b times. In Lemma 4.2 and Corollary 4.3 below, we show that this occurs with high probability when b is suitably 291

- LEMMA 4.2. In one phase, the probability that each coordinate is chosen at most b 293 times by the q processors is at least $1 - n(eq/(b+1))^{b+1}/n^{b+1}$, where e = 2.71828...294
- *Proof.* Let k be one of the coordinates, and let $\mathcal{E}(k)$ denote the event that coordi-295 nate k is chosen more than b times. Note that 296
- 297 $\mathcal{E}(k) =$ {coordinate k is chosen by all the processors with labels in S }. $S:S \subset \{1,2,\cdots,q\}$ |S|=b+1
- Thus, by the union bound, the probability that $\mathcal{E}(k)$ holds is at most $\binom{q}{b+1} \cdot \left(\frac{1}{n}\right)^{b+1}$, which is at most $(eq/(b+1))^{b+1}/n^{b+1}$ by a well-known formula for bounding binomial 298
- 299 coefficients. 300
- By the union bound again, the probability that $\bigcup_{k=1}^n \mathcal{E}(k)$ holds is at most n(eq/(b+301 $(1)^{b+1}/n^{b+1}$. The lemma follows. 302
- COROLLARY 4.3. [Events \mathcal{E}_1 and \mathcal{E}'_1] Suppose that $q \leq \frac{n}{e}$, $c_1, c_2 \geq 1$, and $b \geq e-1$; 303
- also, let $\Lambda = \frac{\log n}{\log(n/q)}$. Then the probability that in each of the first n^{c_1} phases each 304
- coordinate is chosen at most b times by the q processors is at least $1 n^{c_1 + 1 (b+1)/\Lambda}$. 305 Furthermore, 306
- a. [Event \mathcal{E}_1] if $b \geq (c_1 + c_2 + 1)/\Lambda$, the probability is at least $1 1/n^{c_2}$; and 307
- b. [Event \mathcal{E}'_1] if $n \geq 2$ and $b \geq (c_1 + c_2 + 2)/\Lambda$, the probability is at least $1 1/n^{c_2+1} \geq 1$
- $1 1/2n^{c_2}$. 309
- Next, we show the inductive bounds on M^{ℓ} and G^{ℓ} . Recall that $c = c_1 + c_2$. 310
- LEMMA 4.4. For any fixed $c_1, c_2 \geq 1$, conditioned on \mathcal{E}_1 , if $\epsilon \geq \frac{4}{n}$ and $\frac{n}{e} \geq q \geq$ 311
- $\max\left\{\frac{8\Gamma}{\epsilon L_{\max}}, \frac{10(c+1)}{\epsilon} \cdot \frac{\log n}{\log \frac{n}{q}}\right\}, \ \ then \ \ for \ \ \ell = 0, 1, 2, \dots, n^{c_1}, \ \ G^{\ell} \ \geq \ 2^{\ell}G^0 \ \ and \ \ M^{\ell} \ \leq \frac{1}{2} \left(\frac{10(c+1)}{\epsilon} + \frac{\log n}{2}\right)$ 313
- *Proof.* Suppose the lemma holds for some $\ell \geq 0$, and without loss of generality, 314
- assume that $\sum_{j=1}^{n} y_{j}^{\ell}$ is positive; the other case will be symmetric. As in Lemma 4.1, 315
- consider an update in the $(\ell+1)$ -st phase that picks coordinate k; the update reduces 316
- 317 its value by

$$\frac{318}{319} \qquad \frac{L_{\max}}{\Gamma} \cdot \left[(1-\epsilon) y_k^\ell + \epsilon G^\ell \right] \ \geq \ \frac{L_{\max}}{\Gamma} \cdot \left[-(1-\epsilon) M^\ell + \epsilon G^\ell \right] \ \geq \ \frac{3\epsilon L_{\max}}{4\Gamma} G^\ell.$$

Thus, immediately after the $(\ell+1)$ -st phase, $\sum_{j=1}^n y_j^{\ell+1} \leq \sum_{j=1}^n y_j^{\ell} - q \cdot \frac{3\epsilon L_{\max}}{4\Gamma} G^{\ell} = 0$ 320

321
$$\left(1 - \frac{3\epsilon q L_{\max}}{4\Gamma}\right) G^{\ell}$$
. If $q \ge \frac{4\Gamma}{3\epsilon L_{\max}}$, then

322 (4.1)
$$G^{\ell+1} = \left| \sum_{j=1}^{n} y_j^{\ell+1} \right| \ge \left(\frac{3\epsilon q L_{\text{max}}}{4\Gamma} - 1 \right) G^{\ell}.$$

On the other hand, for each coordinate k, if it is chosen by b_k processors in the 323 $(\ell+1)$ -st phase, then the value of $y_k^{\ell+1}$ is 324

$$y_k^{\ell} - b_k \cdot \frac{L_{\max}}{\Gamma} \cdot \left[(1 - \epsilon) y_k^{\ell} + \epsilon G^{\ell} \right].$$

353

For each q we consider, we apply Corollary 4.3(a). Note that $\frac{1}{\Lambda} = \frac{\log n}{\log \frac{n}{2}}$. Conditioned 326 on \mathcal{E}_1 , $b_k \leq \frac{1}{\Lambda}(c+1)$. Also, since $|y_k^{\ell}| \leq M^{\ell} \leq \epsilon G^{\ell}/4$ and $L_{\max} \leq \Gamma$,

$$|y_k^{\ell+1}| \leq \frac{\epsilon}{4} \cdot G^{\ell} + \frac{1}{\Lambda} (c+1) \cdot \frac{L_{\text{max}}}{\Gamma} \cdot \left(\frac{\epsilon}{4} \cdot G^{\ell} + \epsilon G^{\ell}\right)$$

$$= \left(\frac{1}{4} + \frac{5 \cdot \frac{1}{\Lambda} (c+1) L_{\text{max}}}{4\Gamma}\right) \epsilon G^{\ell}.$$

Given Equations (4.1) and (4.2), to guarantee that $G^{\ell+1} \geq 2G^{\ell}$, having $\frac{3\epsilon q L_{\text{max}}}{4\Gamma}$ 331 $1 \ge 2$ suffices. We satisfy this by imposing $q \ge \frac{8\Gamma}{\epsilon L_{\text{max}}}$, or equivalently $\frac{\epsilon q L_{\text{max}}}{4\Gamma} \ge 2$, which is slightly stronger than needed, but will help improve the next constraint on 332 333 q. To guarantee that $M^{\ell+1} \leq \epsilon G^{\ell+1}/4$, the following condition suffices: 334

335 (4.3)
$$\frac{\epsilon}{4} \left[\frac{3\epsilon q L_{\text{max}}}{4\Gamma} - 1 \right] G^{\ell} \ge \left(\frac{1}{4} + \frac{5 \cdot \frac{1}{\Lambda} (c+1) L_{\text{max}}}{4\Gamma} \right) \epsilon G^{\ell}.$$

As $\frac{\epsilon q L_{\max}}{4\Gamma} \ge 2$, we see that $\frac{2\epsilon q L_{\max}}{4\Gamma} \ge \frac{5\cdot (c+1)L_{\max}}{\Lambda\Gamma}$ suffices; i.e. $q \ge \frac{10(c+1)}{\epsilon} \frac{\log n}{\log \frac{n}{q}}$ 336 337

Proof of Theorem 3.4. This theorem follows from Lemma 4.4, which requires that Event \mathcal{E}_1 hold. So the following conditions suffice: 339

340
$$\epsilon \ge \frac{4}{n} \text{ and } \frac{n}{e} \ge q \ge \max \left\{ \frac{8\Gamma}{\epsilon L_{\max}}, \frac{10(c+1)}{\epsilon} \cdot \frac{\log n}{\log \frac{n}{q}} \right\}.$$

If $L_{\rm res} \geq 2L_{\rm max}$ (implying $\frac{1}{\epsilon} \leq \frac{2L_{\rm max}\sqrt{n}}{\sqrt{3}L_{\rm res}}$), the final condition becomes 341

$$q \ge \max \left\{ \frac{16\Gamma\sqrt{n}}{\sqrt{3}L_{\text{res}}}, \frac{20(c+1)\sqrt{n}L_{\text{max}}}{\sqrt{3}L_{\text{res}}} \cdot \frac{\log n}{\log \frac{n}{q}} \right\}.$$

We focus on the second constraint, namely $q \geq \frac{20(c+1)\sqrt{n}L_{\max}}{\sqrt{3}L_{\text{res}}} \cdot \frac{\log n}{\log \frac{n}{q}}$. Let $\mathcal{C} =$ 344

 $\frac{20(c+1)\sqrt{n}L_{\max}}{\sqrt{3}L_{\text{res}}}$. The constraint becomes $q \geq C\frac{\log n}{\log \frac{n}{q}}$, or $q\log \frac{n}{q} \geq C\log n$, where the RHS is independent of q. Note that $q\log \frac{n}{q}$ is an increasing function for $1 \leq q \leq \frac{n}{e}$.

346

Therefore, if suffices to seek a $\hat{q} \geq 1$ such that $\hat{q} \log(n/\hat{q}) \geq \mathcal{C} \log n$; then $q \geq \hat{q}$ implies that the second constraint holds for $q \leq n/e$. 348

Consider $\hat{q} = 10\mathcal{C}$. Then $\hat{q} \log(n/\hat{q}) \geq \mathcal{C} \log n$ is equivalent to the inequality $9\mathcal{C} \log n \geq 10\mathcal{C} \log(10\mathcal{C})$, and hence to $n^9 \geq (10\mathcal{C})^{10}$. Substituting for \mathcal{C} , we obtain 349 350

$$351 \qquad n^9 \geq \left(\frac{200(c+1)}{\sqrt{3}} \cdot \frac{L_{\text{max}}}{L_{\text{res}}} \cdot \sqrt{n}\right)^{10}, \text{ or equivalently, } n \geq \left(\frac{200(c+1)}{\sqrt{3}} \cdot \frac{L_{\text{max}}}{L_{\text{res}}}\right)^{5/2}.$$

352 We also need
$$\hat{q} = 10\mathcal{C} \le n/e$$
; $n \ge \left(\frac{544(c+1)}{\sqrt{3}} \cdot \frac{L_{\text{max}}}{L_{\text{res}}}\right)^2$ suffices.

4.3. Proof of Theorem 3.5.

LEMMA 4.5. Conditioned on Event \mathcal{E}'_1 defined in Corollary 4.3, if $L_{res} \geq 2L_{max}$, 354 $n \ge 400c^2(c+2)^2$, and 355

$$q \ge \max \left\{ \frac{8\Gamma\sqrt{n}}{\sqrt{3}L_{\text{res}}} + \frac{4c\Gamma}{L_{\text{max}}} , 20c(c+2) \right\}$$

then $M^{\ell} \leq \frac{1}{4c}G^{\ell}$ and $G^{\ell} \geq 2^{\ell}G^{0}$, for $1 \leq \ell \leq n^{c_1}$. 357

358

361

367

373

383

384

386

387

388

Proof. We follow the inductive argument in the proof of Lemma 4.4 closely. In the spirit of deriving (4.3), we apply Corollary 4.3(b) instead of Corollary 4.3(a), causing 359 the c+1 term to be replaced by a c+2 term. Then, it suffices that 360

$$\frac{1}{4c} \left[\frac{3\epsilon q L_{\max}}{4\Gamma} - 1 \right] G^{\ell} \ge \left(\frac{1}{4} + \frac{5 \cdot (c+2) L_{\max}}{4\Gamma} \cdot \frac{\log n}{\log \frac{n}{a}} \right) \epsilon G^{\ell}.$$

The above constraint is equivalent to $q \ge \frac{4\Gamma}{3\epsilon L_{\max}} + \frac{4\Gamma c}{3L_{\max}} + \frac{20c(c+2)}{3} \cdot \frac{\log n}{\log \frac{n}{q}}$, so it suffices 362

363 that
$$q \ge \max\left\{\frac{4\Gamma}{L_{\max}}\left(\frac{1}{\epsilon} + c\right), 10c(c+2) \cdot \frac{\log n}{\log \frac{n}{q}}\right\}$$
.

We focus on the second constraint. As argued in the proof of Theorem 3.4, it 364 suffices to find a lower bound \hat{q} on q, such that $\hat{q} \log \frac{n}{\hat{q}} \geq \mathcal{C} \log n$, where $\mathcal{C} \triangleq 10c(c+2)$. 365 $\hat{q} = 2\mathcal{C}$ suffices if $n \geq (2\mathcal{C})^2$. 366

Since
$$q \ge \frac{4\Gamma}{\epsilon L_{\text{max}}}$$
, following the derivation of (4.1) yields $G^{\ell+1} \ge 2G^{\ell}$.

LEMMA 4.6. [Event \mathcal{E}_2] Let $\mathcal{I} = [-\frac{1}{2n}G^{\ell}, \frac{1}{2n}G^{\ell}]$. Suppose all the conditions imposed in Lemma 4.5 hold. Then, with probability at least $1 - 1/n^{c_2}$, for $0 \le \ell < n^{c_1}$, 368 369 at every time during phase $\ell+1$, at least one coordinate has a value outside \mathcal{I} .

Proof. We condition on Event \mathcal{E}'_1 , which by Corollary 4.3, occurs with probability 371 at least $1 - 1/2n^{c_2}$. 372

As before, without loss of generality, we assume that $\sum_i y_i^{\ell} > 0$.

We start by showing that at the start of the $(\ell+1)$ -st phase at least 2c coordinates have values larger than $\frac{1}{2n}G^{\ell}$. Note that $\sum_{i:y_i^{\ell}\in\mathcal{I}}y_i^{\ell}\leq\frac{1}{2}G^{\ell}$. By Lemma 4.5, $M^{\ell}\leq$ 375 $\frac{1}{4c}G^{\ell}$. Thus, at the start of Phase $\ell+1$, at least $\frac{1}{2}G^{\ell}/\frac{1}{4c}G^{\ell}=2c$ coordinates have 376 value greater than $\frac{1}{2n}G^{\ell}$, proving this claim. 377

Next, we observe that an update to a coordinate with value in I changes its value 378 to be smaller than $-\frac{1}{2n}G^{\ell}$. For the largest value resulting from such an update is $\left[1-(1-\epsilon)\frac{L_{\max}}{\Gamma}\right]\frac{1}{2n}G^{\ell}-\frac{L_{\max}}{\Gamma}\epsilon G^{\ell}$. By assumption, $\frac{\Gamma}{L_{\max}}\leq \frac{1}{4}\epsilon q\leq \frac{1}{4}\epsilon n$; therefore $\frac{\epsilon L_{\max}}{\Gamma}\geq \frac{4}{n}$. We deduce the update value is at most $-\left(\frac{4}{n}-\frac{1}{2n}\right)G^{\ell}$, which demonstrates 380 381 382

In order for all the coordinates to end up in \mathcal{I} during Phase $\ell+1$, we need that there be no coordinates less than $-\frac{1}{2n}G^{\ell}$ initially, and that there be no updates to the coordinates in \mathcal{I} until all the other coordinates enter \mathcal{I} , assuming this is possible. This requires some $k \geq 2c$ updates to these at least 2c coordinates.

The probability that these updates happen first is at most

$$\frac{k!}{n^k} \le \frac{(2c)!}{n^{2c}} \le \left(\frac{2c}{n}\right)^{2c},$$

and this is bounded by $1/(2n)^c$ if $n \geq 8c^2$. Summed over all n^{c_1} phases, this gives a 389 total failure probability of (significantly) less than $1/2n^{c_2}$. 390

Taking into account the $1/2n^{c_2}$ probability that Event \mathcal{E}'_1 does not occur, we see 391 that the overall probability of Event \mathcal{E}_2 is at least $1 - 1/n^{c_2}$, as claimed. 392

Proof of Theorem 3.5. Conditioned on Event \mathcal{E}_2 , throughout phase $\ell+1$, at least 393 one coordinate has a value outside the interval \mathcal{I} defined in Lemma 4.6. Thus $f(\mathbf{x}) \geq$ $\left[\frac{1}{2n}G^{\ell}\right]^{2} \geq \frac{1}{4n^{2}} \cdot 4^{\ell} \left(G^{0}\right)^{2} \geq \frac{1}{4} \cdot 4^{\ell}.$

Note that $f(\mathbf{x}^0) = \frac{1-\epsilon}{2} \cdot n + \frac{\epsilon}{2} n^2 \le \epsilon n^2$ as, by assumption, $\epsilon > \frac{1}{n}$. We deduce that $f(\mathbf{x}) \ge \frac{1}{4n^2} \cdot 4^{\ell} \cdot f(\mathbf{x}^0)$ throughout this phase. 396 397

- 5. Discussion. We have shown a tight asymptotic upper bound on the possible 398 parallelism for achieving linear speedup when using asynchronous coordinate descent 399 for almost the whole range of $\frac{L_{\overline{res}}}{L_{\max}}$. This upper bound holds even for composite functions. Furthermore, it holds even in the somewhat restrictive consistent read 401 model, and thus it holds for the inconsistent read model too.
- **Acknowledgments.** We thank the referees for their thoughtful and incisive com-403 404 ments.
 - Appendix A. Missing Proofs.
- LEMMA A.1. The strong convexity parameter of f_{ϵ} is (1ϵ) . 406 Proof.

$$407 f_{\epsilon}(y) - f_{\epsilon}(x) - \langle \nabla f_{\epsilon}(x), y - x \rangle
408 \geq \frac{1 - \epsilon}{2} \sum_{i} y_{i}^{2} - x_{i}^{2} + \frac{\epsilon}{2} \left[\left(\sum_{i} y_{i} \right)^{2} - \left(\sum_{i} x_{i} \right)^{2} \right]
409 - (1 - \epsilon) \sum_{i} (x_{i}y_{i} - x_{i}^{2}) - \epsilon \left[\left(\sum_{i} x_{i} \right) \left(\sum_{i} y_{i} - \sum_{i} x_{i} \right) \right]
410 \geq \frac{1 - \epsilon}{2} \sum_{i} (y_{i}^{2} + x_{i}^{2} - 2x_{i}y_{i}) + \frac{\epsilon}{2} \left[\left(\sum_{i} y_{i} \right)^{2} + \left(\sum_{i} x_{i} \right)^{2} - 2 \sum_{i} x_{i} \sum_{i} y_{i} \right]
411 \geq \frac{1 - \epsilon}{2} \sum_{i} (y_{i} - x_{i})^{2} + \frac{\epsilon}{2} \left[\sum_{i} y_{i} - \sum_{i} x_{i} \right]^{2}
412 \geq \frac{1 - \epsilon}{2} ||y - x||^{2}. \qquad \Box$$

Proof of Proposition 3.2. Recall that $f_{\epsilon}(\mathbf{x}^*) = 0$. Also, recall that the SCD 414 process starts at the all ones point. Let C_1 denotes the even index coordinates, and C_{-1} those of odd index. Consider the following random variables:

417
$$S_1(t) := \sum_{j \in C_1} x_j^t \qquad S_{-1}(t) := \sum_{j \in C_{-1}} (-x_j^t) \qquad S(t) := S_1(t) + S_{-1}(t).$$

- Note that $S_1(0) = S_{-1}(0) = n/2$ and S(0) = n. 418
- Next, we derive a recurrence which, conditioned on S(t), computes the expected 419 value of S(t+1) - S(t). Recall that if coordinate j is chosen to be updated at time 420 421
- t+1, then

422
$$x_j^{t+1} - x_j^t = -\frac{\nabla_j f(\mathbf{x}^t)}{\Gamma} = -\frac{1 - \epsilon}{\Gamma} \cdot x_j^t - \frac{\epsilon}{\Gamma} \cdot [S_1(t) - S_{-1}(t)].$$

423 Thus,

424
$$\mathbb{E}\left[S(t+1) - S(t) \mid S(t)\right] = \frac{1}{n} \left[\sum_{j \in C_1} \left(-\frac{1-\epsilon}{\Gamma} \cdot x_j^t - \frac{\epsilon}{\Gamma} \cdot [S_1(t) - S_{-1}(t)]\right) + \sum_{j \in C_{-1}} \left(\frac{1-\epsilon}{\Gamma} \cdot x_j^t + \frac{\epsilon}{\Gamma} \cdot [S_1(t) - S_{-1}(t)]\right)\right]$$
425
$$= \frac{1}{n} \left(-\frac{1-\epsilon}{\Gamma} \cdot S_1(t) - \frac{1-\epsilon}{\Gamma} \cdot S_{-1}(t)\right)$$
426
$$= -\frac{1-\epsilon}{n\Gamma} S(t).$$

- The second equality above holds because $|C_1| = |C_{-1}|$, leading to cancellation of the
- 430 terms $\pm \frac{\epsilon}{\Gamma} \cdot [S_1(t) S_{-1}(t)]$. Thus, $\mathbb{E}[S(t+1) \mid S(t)] = S(t) \cdot (1 \frac{1-\epsilon}{n\Gamma})$. Iterating this
- 431 recurrence yields

432 (A.1)
$$\mathbb{E}\left[S(t)\right] = n \cdot \left(1 - \frac{1 - \epsilon}{n\Gamma}\right)^{t}.$$

- Next, observe that for any fixed $S_1(t), S_{-1}(t)$, by the Power-Mean Inequality,
- 434 $\sum_{j \in C_1} (x_j^t)^2 \ge \frac{2|S_1(t)|^2}{n}$ and $\sum_{j \in C_{-1}} (x_j^t)^2 \ge \frac{2|S_{-1}(t)|^2}{n}$. On the other hand, for any
- 435 fixed S(t), which equals to $S_1(t) + S_{-1}(t)$, the sum $\frac{|S_1(t)|^2}{n} + \frac{|S_{-1}(t)|^2}{n}$ is minimized
- 436 when $S_1(t) = S_{-1}(t) = S(t)/2$. Thus,

437
$$\mathbb{E}\left[f_{\epsilon}(\mathbf{x}^{t})\right] \geq \mathbb{E}\left[\frac{1-\epsilon}{2}\sum_{j=1}^{n}(x_{j}^{t})^{2}\right] \geq (1-\epsilon)\cdot\mathbb{E}\left[\frac{|S_{1}(t)|^{2}}{n} + \frac{|S_{-1}(t)|^{2}}{n}\right]$$

$$\geq \frac{1-\epsilon}{2n}\cdot\mathbb{E}\left[S(t)^{2}\right].$$

Finally, we complete the proof by using the Cauchy-Schwarz inequality and Equation (A.1):

442
$$\frac{1-\epsilon}{2n} \cdot \mathbb{E}\left[S(t)^{2}\right] \geq \frac{1-\epsilon}{2n} \cdot \mathbb{E}\left[S(t)\right]^{2} = \frac{(1-\epsilon)n}{2} \cdot \left(1 - \frac{1-\epsilon}{n\Gamma}\right)^{2t}$$

$$\geq \frac{(1-\epsilon)n}{2} \cdot \left(1 - \frac{2(1-\epsilon)}{n\Gamma}\right)^{t}$$

$$= f_{\epsilon}(\mathbf{x}^{0}) \cdot \left(1 - \frac{2(1-\epsilon)}{n\Gamma}\right)^{t}.$$

446 REFERENCES

- 447 [1] Y. K. CHEUNG, R. COLE AND Y. TAO, Fully asynchronous stochastic coordinate descent: A
 448 tight lower bound on the parallelism achieving linear speedup, Math. Program. Series A, (to
 449 appear).
- 450 [2] J. LIU AND S. J. WRIGHT, Asynchronous stochastic coordinate descent: Parallelism and convergence properties, SIAM J. Optim., 25 (2015), pp. 351–376.
- 452 [3] J. LIU, S. J. WRIGHT, C. RÉ, V. BITTORF AND S. SRIDHAR, An asynchronous parallel stochastic coordinate descent algorithm, J. Mach. Learn. Res., 16 (2015), pp. 285–322.

- 454 [4] H. Mania, X. Pan, D. S. Papailiopoulos, B. Recht, K. Ramchandran and M. I. Jordan, 455 Perturbed iterate analysis for asynchronous stochastic optimization, SIAM J. Optim., 27 456 (2017), pp. 2202–2229.
- 457 [5] T. Sun, R. Hannah and W. Yin, Asynchronous coordinate descent under more realistic assumptions, Advances in Neural Information Processing Systems, (2017), pp. 6182–6190.