

# An Asynchronous Parallel Stochastic Coordinate Descent Algorithm

**Ji Liu**

**Stephen J. Wright**

*Department of Computer Sciences  
University of Wisconsin-Madison  
Madison, WI 53706-1685*

JL.LIU.UWISC@GMAIL.COM

SWRIGHT@CS.WISC.EDU

**Christopher Ré**

*Department of Computer Science  
Stanford University  
353 Serra Mall  
Stanford, CA 94305-9025*

CHRISMRE@CS.STANFORD.EDU

**Victor Bittorf**

**Srikrishna Sridhar**

*Department of Computer Sciences  
University of Wisconsin-Madison  
Madison, WI 53706-1685*

BITTORF@CS.WISC.EDU

SRIKRIS@CS.WISC.EDU

**Editor:** Leon Bottou

## Abstract

We describe an asynchronous parallel stochastic coordinate descent algorithm for minimizing smooth unconstrained or separably constrained functions. The method achieves a linear convergence rate on functions that satisfy an essential strong convexity property and a sublinear rate  $(1/K)$  on general convex functions. Near-linear speedup on a multicore system can be expected if the number of processors is  $O(n^{1/2})$  in unconstrained optimization and  $O(n^{1/4})$  in the separable-constrained case, where  $n$  is the number of variables. We describe results from implementation on 40-core processors.

**Keywords:** asynchronous parallel optimization, stochastic coordinate descent

## 1. Introduction

Consider the convex optimization problem

$$\min_{x \in \Omega} f(x), \quad (1)$$

where  $\Omega \subset \mathbb{R}^n$  is a closed convex set and  $f$  is a smooth convex mapping from an open neighborhood of  $\Omega$  to  $\mathbb{R}$ . We consider two particular cases of  $\Omega$  in this paper: the unconstrained case  $\Omega = \mathbb{R}^n$ , and the separable case

$$\Omega = \Omega_1 \times \Omega_2 \times \dots \times \Omega_n, \quad (2)$$

where each  $\Omega_i$ ,  $i = 1, 2, \dots, n$  is a closed subinterval of the real line.

Formulations of the type (1,2) arise in many data analysis and machine learning problems, for example, support vector machines (linear or nonlinear dual formulation) (Cortes and Vapnik, 1995), LASSO (after decomposing  $x$  into positive and negative parts) (Tibshirani, 1996), and logistic regression. Algorithms based on gradient and approximate or partial gradient information have proved effective in these settings. We mention in particular gradient projection and its accelerated variants (Nesterov, 2004), accelerated proximal gradient methods for regularized objectives (Beck and Teboulle, 2009), and stochastic gradient methods (Nemirovski et al., 2009; Shamir and Zhang, 2013). These methods are inherently serial, in that each iteration depends on the result of the previous iteration. Recently, parallel multicore versions of stochastic gradient and stochastic coordinate descent have been described for problems involving large data sets; see for example Niu et al. (2011); Richtárik and Takáč (2012b); Avron et al. (2014).

This paper proposes an asynchronous stochastic coordinate descent (ASYSCD) algorithm for convex optimization. Each step of ASYSCD chooses an index  $i \in \{1, 2, \dots, n\}$  and subtracts a short, constant, positive multiple of the  $i$ th partial gradient  $\nabla_i f(x) := \partial f / \partial x_i$  from the  $i$ th component of  $x$ . When separable constraints (2) are present, the update is “clipped” to maintain feasibility with respect to  $\Omega_i$ . Updates take place in parallel across the cores of a multicore system, without any attempt to synchronize computation between cores. We assume that there is a bound  $\tau$  on the age of the updates, that is, no more than  $\tau$  updates to  $x$  occur between the time at which a processor reads  $x$  (and uses it to evaluate one element of the gradient) and the time at which this processor makes its update to a single element of  $x$ . (A similar model of parallel asynchronous computation was used in HOGWILD! (Niu et al., 2011).) Our implementation, described in Section 6, is a little more complex than this simple model would suggest, as it is tailored to the architecture of the Intel Xeon machine that we use for experiments.

We show that linear convergence can be attained if an “essential strong convexity” property (3) holds, while sublinear convergence at a “ $1/K$ ” rate can be proved for general convex functions. Our analysis also defines a sufficient condition for near-linear speedup in the number of cores used. This condition relates the value of delay parameter  $\tau$  (which relates to the number of cores / threads used in the computation) to the problem dimension  $n$ . A parameter that quantifies the cross-coordinate interactions in  $\nabla f$  also appears in this relationship. When the Hessian of  $f$  is nearly diagonal, the minimization problem can almost be separated along the coordinate axes, so higher degrees of parallelism are possible.

We review related work in Section 2. Section 3 specifies the proposed algorithm. Convergence results for unconstrained and constrained cases are described in Sections 4 and 5, respectively, with proofs given in the appendix. Computational experience is reported in Section 6. We discuss several variants of ASYSCD in Section 7. Some conclusions are given in Section 8.

## 1.1 Notation and Assumption

We use the following notation.

- $e_i \in \mathbb{R}^n$  denotes the  $i$ th natural basis vector  $(0, \dots, 0, 1, 0, \dots, 0)^T$  with the “1” in the  $i$ th position.
- $\|\cdot\|$  denotes the Euclidean norm  $\|\cdot\|_2$ .

- $S \subset \Omega$  denotes the set on which  $f$  attains its optimal value, which is denoted by  $f^*$ .
- $\mathcal{P}_S(\cdot)$  and  $\mathcal{P}_\Omega(\cdot)$  denote Euclidean projection onto  $S$  and  $\Omega$ , respectively.
- We use  $x_i$  for the  $i$ th element of  $x$ , and  $\nabla_i f(x)$  for the  $i$ th element of the gradient vector  $\nabla f(x)$ .
- We define the following *essential strong convexity* condition for a convex function  $f$  with respect to the optimal set  $S$ , with parameter  $l > 0$ :

$$f(x) - f(y) \geq \langle \nabla f(y), x - y \rangle + \frac{l}{2} \|x - y\|^2 \quad \text{for all } x, y \in \Omega \text{ with } \mathcal{P}_S(x) = \mathcal{P}_S(y). \quad (3)$$

This condition is significantly weaker than the usual strong convexity condition, which requires the inequality to hold for *all*  $x, y \in \Omega$ . In particular, it allows for non-singleton solution sets  $S$ , provided that  $f$  increases at a uniformly quadratic rate with distance from  $S$ . (This property is noted for convex quadratic  $f$  in which the Hessian is rank deficient.) Other examples of essentially strongly convex functions that are not strongly convex include:

- $f(Ax)$  with arbitrary linear transformation  $A$ , where  $f(\cdot)$  is strongly convex;
- $f(x) = \max(a^T x - b, 0)^2$ , for  $a \neq 0$ .
- Define  $L_{\text{res}}$  as the *restricted Lipschitz constant* for  $\nabla f$ , where the “restriction” is to the coordinate directions: We have

$$\|\nabla f(x) - \nabla f(x + te_i)\| \leq L_{\text{res}} |t|, \quad \text{for all } i = 1, 2, \dots, n \text{ and } t \in \mathbb{R}, \text{ with } x, x + te_i \in \Omega.$$

- Define  $L_i$  as the *coordinate Lipschitz constant* for  $\nabla f$  in the  $i$ th coordinate direction: We have

$$f(x + te_i) - f(x) \leq \langle \nabla_i f(x), t \rangle + \frac{L_i}{2} t^2, \quad \text{for } i \in \{1, 2, \dots, n\}, \text{ and } x, x + te_i \in \Omega,$$

or equivalently

$$|\nabla_i f(x) - \nabla_i f(x + te_i)| \leq L_i |t|.$$

- $L_{\text{max}} := \max_{i=1,2,\dots,n} L_i$ .

Note that  $L_{\text{res}} \geq L_{\text{max}}$ .

We use  $\{x_j\}_{j=0,1,2,\dots}$  to denote the sequence of iterates generated by the algorithm from starting point  $x_0$ . Throughout the paper, we make the following assumption.

### Assumption 1

- The optimal solution set  $S$  of (1) is nonempty.
- The radius of the iterate set  $\{x_j\}_{j=0,1,2,\dots}$  defined by

$$R := \sup_{j=0,1,2,\dots} \|x_j - \mathcal{P}_S(x_j)\|$$

is bounded, that is,  $R < +\infty$ .

## 1.2 Lipschitz Constants

The nonstandard Lipschitz constants  $L_{\text{res}}$ ,  $L_{\text{max}}$ , and  $L_i$ ,  $i = 1, 2, \dots, n$  defined above are crucial in the analysis of our method. Besides bounding the nonlinearity of  $f$  along various directions, these quantities capture the interactions between the various components in the gradient  $\nabla f$ , as quantified in the off-diagonal terms of the Hessian  $\nabla^2 f(x)$  — although the stated conditions do not require this matrix to exist.

We have noted already that  $L_{\text{res}}/L_{\text{max}} \geq 1$ . Let us consider upper bounds on this ratio under certain conditions. When  $f$  is twice continuously differentiable, we have

$$L_i = \sup_{x \in \Omega} \max_{i=1,2,\dots,n} [\nabla^2 f(x)]_{ii}.$$

Since  $\nabla^2 f(x) \succeq 0$  for  $x \in \Omega$ , we have that

$$|[\nabla^2 f(x)]_{ij}| \leq \sqrt{L_i L_j} \leq L_{\text{max}}, \quad \forall i, j = 1, 2, \dots, n.$$

Thus  $L_{\text{res}}$ , which is a bound on the largest column norm for  $\nabla^2 f(x)$  over all  $x \in \Omega$ , is bounded by  $\sqrt{n}L_{\text{max}}$ , so that

$$\frac{L_{\text{res}}}{L_{\text{max}}} \leq \sqrt{n}.$$

If the Hessian is structurally sparse, having at most  $p$  nonzeros per row/column, the same argument leads to  $L_{\text{res}}/L_{\text{max}} \leq \sqrt{p}$ .

If  $f(x)$  is a convex quadratic with Hessian  $Q$ , we have

$$L_{\text{max}} = \max_i Q_{ii}, \quad L_{\text{res}} = \max_i \|Q_{\cdot i}\|_2,$$

where  $Q_{\cdot i}$  denotes the  $i$ th column of  $Q$ . If  $Q$  is diagonally dominant, we have for any column  $i$  that

$$\|Q_{\cdot i}\|_2 \leq Q_{ii} + \|[Q_{ji}]_{j \neq i}\|_2 \leq Q_{ii} + \sum_{j \neq i} |Q_{ji}| \leq 2Q_{ii},$$

which, by taking the maximum of both sides, implies that  $L_{\text{res}}/L_{\text{max}} \leq 2$  in this case.

Finally, consider the objective  $f(x) = \frac{1}{2}\|Ax - b\|^2$  and assume that  $A \in \mathbb{R}^{m \times n}$  is a random matrix whose entries are i.i.d from  $\mathcal{N}(0, 1)$ . The diagonals of the Hessian are  $A_{\cdot i}^T A_{\cdot i}$  (where  $A_{\cdot i}$  is the  $i$ th column of  $A$ ), which have expected value  $m$ , so we can expect  $L_{\text{max}}$  to be not less than  $m$ . Recalling that  $L_{\text{res}}$  is the maximum column norm of  $A^T A$ , we have

$$\begin{aligned} \mathbb{E}(\|A^T A_{\cdot i}\|) &\leq \mathbb{E}(|A_{\cdot i}^T A_{\cdot i}|) + \mathbb{E}(\|[A_{\cdot j}^T A_{\cdot i}]_{j \neq i}\|) \\ &= m + \mathbb{E} \sqrt{\sum_{j \neq i} |A_{\cdot j}^T A_{\cdot i}|^2} \\ &\leq m + \sqrt{\sum_{j \neq i} \mathbb{E} |A_{\cdot j}^T A_{\cdot i}|^2} \\ &= m + \sqrt{(n-1)m}, \end{aligned}$$

where the second inequality uses Jensen's inequality and the final equality uses

$$\mathbb{E}(|A_{\cdot j}^T A_{\cdot i}|^2) = \mathbb{E}(A_{\cdot j}^T \mathbb{E}(A_{\cdot i} A_{\cdot i}^T) A_{\cdot j}) = \mathbb{E}(A_{\cdot j}^T I A_{\cdot j}) = \mathbb{E}(A_{\cdot j}^T A_{\cdot j}) = m.$$

We can thus estimate the upper bound on  $L_{\text{res}}/L_{\text{max}}$  roughly by  $1 + \sqrt{n/m}$  for this case.

## 2. Related Work

This section reviews some related work on coordinate relaxation and stochastic gradient algorithms.

Among *cyclic coordinate descent algorithms*, Tseng (2001) proved the convergence of a *block coordinate descent method* for nondifferentiable functions with certain conditions. Local and global linear convergence were established under additional assumptions, by Luo and Tseng (1992) and Wang and Lin (2014), respectively. Global linear (sublinear) convergence rate for strongly (weakly) convex optimization was proved by Beck and Tetruashvili (2013). Block-coordinate approaches based on proximal-linear subproblems are described by Tseng and Yun (2009, 2010). Wright (2012) uses acceleration on reduced spaces (corresponding to the optimal manifold) to improve the local convergence properties of this approach.

*Stochastic coordinate descent* is almost identical to cyclic coordinate descent except selecting coordinates in a random manner. Nesterov (2012) studied the convergence rate for a stochastic block coordinate descent method for unconstrained and separably constrained convex smooth optimization, proving linear convergence for the strongly convex case and a sublinear  $1/K$  rate for the convex case. Extensions to minimization of composite functions are described by Richtárik and Takáč (2012a) and Lu and Xiao (2013).

*Synchronous parallel methods* distribute the workload and data among multiple processors, and coordinate the computation among processors. Ferris and Mangasarian (1994) proposed to distribute variables among multiple processors and optimize concurrently over each subset. The synchronization step searches the affine hull formed by the current iterate and the points found by each processor. Similar ideas appeared in (Mangasarian, 1995), with a different synchronization step. Goldfarb and Ma (2012) considered a multiple splitting algorithm for functions of the form  $f(x) = \sum_{k=1}^N f_k(x)$  in which  $N$  models are optimized separately and concurrently, then combined in an synchronization step. The alternating direction method-of-multiplier (ADMM) framework (Boyd et al., 2011) can also be implemented in parallel. This approach dissects the problem into multiple subproblems (possibly after replication of primal variables) and optimizes concurrently, then synchronizes to update multiplier estimates. Duchi et al. (2012) described a subgradient dual-averaging algorithm for partially separable objectives, with subgradient evaluations distributed between cores and combined in ways that reflect the structure of the objective. Parallel stochastic gradient approaches have received broad attention; see Agarwal and Duchi (2011) for an approach that allows delays between evaluation and update, and Cotter et al. (2011) for a *minibatch stochastic gradient approach* with Nesterov acceleration. Shalev-Shwartz and Zhang (2013) proposed an accelerated stochastic dual coordinate ascent method.

Among *synchronous parallel methods for (block) coordinate descent*, Richtárik and Takáč (2012b) described a method of this type for convex composite optimization problems. All processors update randomly selected coordinates or blocks, concurrently and synchronously, at each iteration. Speedup depends on the sparsity of the data matrix that defines the loss functions. Several variants that select blocks greedily are considered by Scherrer et al. (2012) and Peng et al. (2013). Yang (2013) studied the parallel stochastic dual coordinate ascent method and emphasized the balance between computation and communication.

We turn now to *asynchronous parallel methods*. Bertsekas and Tsitsiklis (1989) introduced an asynchronous parallel implementation for general fixed point problems  $x = q(x)$  over a separable convex closed feasible region. (The optimization problem (1) can be formulated in this way by defining  $q(x) := \mathcal{P}_\Omega[(I - \alpha \nabla f)(x)]$  for some fixed  $\alpha > 0$ .) Their analysis allows inconsistent reads for  $x$ , that is, the coordinates of the read  $x$  have different “ages.” Linear convergence is established if all ages are bounded and  $\nabla^2 f(x)$  satisfies a diagonal dominance condition guaranteeing that the iteration  $x = q(x)$  is a maximum-norm contraction mapping for sufficient small  $\alpha$ . However, this condition is strong — stronger, in fact, than the strong convexity condition. For convex quadratic optimization  $f(x) = \frac{1}{2}x^T A x + b x$ , the contraction condition requires diagonal dominance of the Hessian:  $A_{ii} > \sum_{i \neq j} |A_{ij}|$  for all  $i = 1, 2, \dots, n$ . By comparison, ASYSCD guarantees linear convergence rate under the essential strong convexity condition (3), though we do not allow inconsistent read. (We require the vector  $x$  used for each evaluation of  $\nabla_i f(x)$  to have existed at a certain point in time.)

HOGWILD! (Niu et al., 2011) is a lock-free, asynchronous parallel implementation of a stochastic-gradient method, targeted to a multicore computational model similar to the one considered here. Its analysis assumes consistent reading of  $x$ , and it is implemented without locking or coordination between processors. Under certain conditions, convergence of HOGWILD! approximately matches the sublinear  $1/K$  rate of its serial counterpart, which is the constant-steplength stochastic gradient method analyzed in Nemirovski et al. (2009).

We also note recent work by Avron et al. (2014), who proposed an asynchronous linear solver to solve  $Ax = b$  where  $A$  is a symmetric positive definite matrix, proving a linear convergence rate. Both inconsistent- and consistent-read cases are analyzed in this paper, with the convergence result for inconsistent read being slightly weaker.

### 3. Algorithm

In ASYSCD, multiple processors have access to a shared data structure for the vector  $x$ , and each processor is able to compute a randomly chosen element of the gradient vector  $\nabla f(x)$ . Each processor repeatedly runs the following coordinate descent process (the steplength parameter  $\gamma$  is discussed further in the next section):

- R: Choose an index  $i \in \{1, 2, \dots, n\}$  at random, read  $x$ , and evaluate  $\nabla_i f(x)$ ;
- U: Update component  $i$  of the shared  $x$  by taking a step of length  $\gamma/L_{\max}$  in the direction  $-\nabla_i f(x)$ .

Since these processors are being run concurrently and without synchronization,  $x$  may change between the time at which it is read (in step R) and the time at which it is updated (step U). We capture the system-wide behavior of ASYSCD in Algorithm 1. There is a global counter  $j$  for the total number of updates;  $x_j$  denotes the state of  $x$  after  $j$  updates. The index  $i(j) \in \{1, 2, \dots, n\}$  denotes the component updated at step  $j$ .  $k(j)$  denotes the  $x$ -iterate at which the update applied at iteration  $j$  was calculated. Obviously, we have  $k(j) \leq j$ , but we assume that the delay between the time of evaluation and updating is bounded uniformly by a positive integer  $\tau$ , that is,  $j - k(j) \leq \tau$  for all  $j$ . The value of  $\tau$  captures the essential parallelism in the method, as it indicates the number of processors that are involved in the computation.

---

**Algorithm 1** Asynchronous Stochastic Coordinate Descent Algorithm  $x_{K+1} = \text{ASYSCD}(x_0, \gamma, K)$

---

**Require:**  $x_0 \in \Omega$ ,  $\gamma$ , and  $K$

**Ensure:**  $x_{K+1}$

- 1: Initialize  $j \leftarrow 0$ ;
  - 2: **while**  $j \leq K$  **do**
  - 3:   Choose  $i(j)$  from  $\{1, \dots, n\}$  with equal probability;
  - 4:    $x_{j+1} \leftarrow \mathcal{P}_\Omega \left( x_j - \frac{\gamma}{L_{\max}} e_{i(j)} \nabla_{i(j)} f(x_{k(j)}) \right)$ ;
  - 5:    $j \leftarrow j + 1$ ;
  - 6: **end while**
- 

The projection operation  $P_\Omega$  onto the feasible set is not needed in the case of unconstrained optimization. For separable constraints (2), it requires a simple clipping operation on the  $i(j)$  component of  $x$ .

We note several differences with earlier asynchronous approaches. Unlike the asynchronous scheme in Bertsekas and Tsitsiklis (1989, Section 6.1), the *latest* value of  $x$  is updated at each step, not an earlier iterate. Although our model of computation is similar to HOGWILD! (Niu et al., 2011), the algorithm differs in that each iteration of ASYSCD evaluates a single component of the gradient exactly, while HOGWILD! computes only a (usually crude) estimate of the full gradient. Our analysis of ASYSCD below is comprehensively different from that of Niu et al. (2011), and we obtain stronger convergence results.

#### 4. Unconstrained Smooth Convex Case

This section presents results about convergence of ASYSCD in the unconstrained case  $\Omega = \mathbb{R}^n$ . The theorem encompasses both the linear rate for essentially strongly convex  $f$  and the sublinear rate for general convex  $f$ . The result depends strongly on the delay parameter  $\tau$ . (Proofs of results in this section appear in Appendix A.) In Algorithm 1, the indices  $i(j)$ ,  $j = 0, 1, 2, \dots$  are random variables. We denote the expectation over all random variables as  $\mathbb{E}$ , the conditional expectation in term of  $i(j)$  given  $i(0), i(1), \dots, i(j-1)$  as  $\mathbb{E}_{i(j)}$ .

A crucial issue in ASYSCD is the choice of steplength parameter  $\gamma$ . This choice involves a tradeoff: We would like  $\gamma$  to be long enough that significant progress is made at each step, but not so long that the gradient information computed at step  $k(j)$  is stale and irrelevant by the time the update is applied at step  $j$ . We enforce this tradeoff by means of a bound on the ratio of expected squared norms on  $\nabla f$  at successive iterates; specifically,

$$\rho^{-1} \leq \frac{\mathbb{E} \|\nabla f(x_{j+1})\|^2}{\mathbb{E} \|\nabla f(x_j)\|^2} \leq \rho, \quad (4)$$

where  $\rho > 1$  is a user defined parameter. The analysis becomes a delicate balancing act in the choice of  $\rho$  and steplength  $\gamma$  between aggression and excessive conservatism. We find, however, that these values can be chosen to ensure steady convergence for the asynchronous

method at a *linear* rate, with rate constants that are almost consistent with vanilla short-step full-gradient descent.

**Theorem 1** *Suppose that  $\Omega = \mathbb{R}^n$  in (1) and that Assumption 1 is satisfied. For any  $\rho > 1$ , define the quantity  $\psi$  as follows:*

$$\psi := 1 + \frac{2\tau\rho^\tau L_{\text{res}}}{\sqrt{n}L_{\text{max}}}. \quad (5)$$

*Suppose that the steplength parameter  $\gamma > 0$  satisfies the following three upper bounds:*

$$\gamma \leq \frac{1}{\psi}, \quad (6a)$$

$$\gamma \leq \frac{(\rho - 1)\sqrt{n}L_{\text{max}}}{2\rho^{\tau+1}L_{\text{res}}}, \quad (6b)$$

$$\gamma \leq \frac{(\rho - 1)\sqrt{n}L_{\text{max}}}{L_{\text{res}}\rho^\tau(2 + \frac{L_{\text{res}}}{\sqrt{n}L_{\text{max}}})}. \quad (6c)$$

*Then we have that for any  $j \geq 0$  that*

$$\rho^{-1}\mathbb{E}(\|\nabla f(x_j)\|^2) \leq \mathbb{E}(\|\nabla f(x_{j+1})\|^2) \leq \rho\mathbb{E}(\|\nabla f(x_j)\|^2). \quad (7)$$

*Moreover, if the essentially strong convexity property (3) holds with  $l > 0$ , we have*

$$\mathbb{E}(f(x_j) - f^*) \leq \left(1 - \frac{2l\gamma}{nL_{\text{max}}} \left(1 - \frac{\psi}{2}\gamma\right)\right)^j (f(x_0) - f^*), \quad (8)$$

*while for general smooth convex functions  $f$ , we have*

$$\mathbb{E}(f(x_j) - f^*) \leq \frac{1}{(f(x_0) - f^*)^{-1} + j\gamma(1 - \frac{\psi}{2}\gamma)/(nL_{\text{max}}R^2)}. \quad (9)$$

This theorem demonstrates linear convergence (8) for ASYSCD in the unconstrained essentially strongly convex case. This result is better than that obtained for HOGWILD! (Niu et al., 2011), which guarantees only sublinear convergence under the stronger assumption of strict convexity.

The following corollary proposes an interesting particular choice of the parameters for which the convergence expressions become more comprehensible. The result requires a condition on the delay bound  $\tau$  in terms of  $n$  and the ratio  $L_{\text{max}}/L_{\text{res}}$ .

**Corollary 2** *Suppose that Assumption 1 holds, and that*

$$\tau + 1 \leq \frac{\sqrt{n}L_{\text{max}}}{2eL_{\text{res}}}. \quad (10)$$

*Then if we choose*

$$\rho = 1 + \frac{2eL_{\text{res}}}{\sqrt{n}L_{\text{max}}}, \quad (11)$$



define  $\psi$  by (5), and set  $\gamma = 1/\psi$ , we have for the essentially strongly convex case (3) with  $l > 0$  that

$$\mathbb{E}(f(x_j) - f^*) \leq \left(1 - \frac{l}{2nL_{\max}}\right)^j (f(x_0) - f^*), \quad (12)$$

while for the case of general convex  $f$ , we have

$$\mathbb{E}(f(x_j) - f^*) \leq \frac{1}{(f(x_0) - f^*)^{-1} + j/(4nL_{\max}R^2)}. \quad (13)$$

We note that the linear rate (12) is broadly consistent with the linear rate for the classical steepest descent method applied to strongly convex functions, which has a rate constant of  $(1 - 2l/L)$ , where  $L$  is the standard Lipschitz constant for  $\nabla f$ . If we assume (not unreasonably) that  $n$  steps of stochastic coordinate descent cost roughly the same as one step of steepest descent, and note from (12) that  $n$  steps of stochastic coordinate descent would achieve a reduction factor of about  $(1 - l/(2L_{\max}))$ , a standard argument would suggest that stochastic coordinate descent would require about  $4L_{\max}/L$  times more computation. (Note that  $L_{\max}/L \in [1/n, 1]$ .) The stochastic approach may gain an advantage from the parallel implementation, however. Steepest descent requires synchronization and careful division of gradient evaluations, whereas the stochastic approach can be implemented in an asynchronous fashion.

For the general convex case, (13) defines a sublinear rate, whose relationship with the rate of the steepest descent for general convex optimization is similar to the previous paragraph.

As noted in Section 1, the parameter  $\tau$  is closely related to the number of cores that can be involved in the computation, without degrading the convergence performance of the algorithm. In other words, if the number of cores is small enough such that (10) holds, the convergence expressions (12), (13) do not depend on the number of cores, implying that linear speedup can be expected. A small value for the ratio  $L_{\text{res}}/L_{\max}$  (not much greater than 1) implies a greater degree of potential parallelism. As we note at the end of Section 1, this ratio tends to be small in some important applications — a situation that would allow  $O(\sqrt{n})$  cores to be used with near-linear speedup.

We conclude this section with a high-probability estimate for convergence of the sequence of function values.

**Theorem 3** *Suppose that the assumptions of Corollary 2 hold, including the definitions of  $\rho$  and  $\psi$ . Then for any  $\epsilon \in (0, f(x_0) - f^*)$  and  $\eta \in (0, 1)$ , we have that*

$$\mathbb{P}(f(x_j) - f^* \leq \epsilon) \geq 1 - \eta, \quad (14)$$

*provided that either of the following sufficient conditions hold for the index  $j$ . In the essentially strongly convex case (3) with  $l > 0$ , it suffices to have*

$$j \geq \frac{2nL_{\max}}{l} \left\lceil \log \frac{f(x_0) - f^*}{\epsilon\eta} \right\rceil, \quad (15)$$

*while in the general convex case, a sufficient condition is*

$$j \geq 4nL_{\max}R^2 \left( \frac{1}{\epsilon\eta} - \frac{1}{f(x_0) - f^*} \right). \quad (16)$$

## 5. Constrained Smooth Convex Case

This section considers the case of separable constraints (2). We show results about convergence rates and high-probability complexity estimates, analogous to those of the previous section. Proofs appear in Appendix B.

As in the unconstrained case, the steplength  $\gamma$  should be chosen to ensure steady progress while ensuring that update information does not become too stale. Because constraints are present, the ratio (4) is no longer appropriate. We use instead a ratio of squares of expected differences in successive primal iterates:

$$\frac{\mathbb{E}\|x_{j-1} - \bar{x}_j\|^2}{\mathbb{E}\|x_j - \bar{x}_{j+1}\|^2}, \quad (17)$$

where  $\bar{x}_{j+1}$  is the hypothesized full update obtained by applying the single-component update to *every* component of  $x_j$ , that is,

$$\bar{x}_{j+1} := \arg \min_{x \in \Omega} \langle \nabla f(x_{k(j)}), x - x_j \rangle + \frac{L_{\max}}{2\gamma} \|x - x_j\|^2.$$

In the unconstrained case  $\Omega = \mathbb{R}^n$ , the ratio (17) reduces to

$$\frac{\mathbb{E}\|\nabla f(x_{k(j-1)})\|^2}{\mathbb{E}\|\nabla f(x_{k(j)})\|^2},$$

which is evidently related to (4), but not identical.

We have the following result concerning convergence of the expected error to zero.

**Theorem 4** *Suppose that  $\Omega$  has the form (2), that Assumption 1 is satisfied, and that  $n \geq 5$ . Let  $\rho$  be a constant with  $\rho > (1 - 2/\sqrt{n})^{-1}$ , and define the quantity  $\psi$  as follows:*

$$\psi := 1 + \frac{L_{\text{res}}\tau\rho^\tau}{\sqrt{n}L_{\max}} \left( 2 + \frac{L_{\max}}{\sqrt{n}L_{\text{res}}} + \frac{2\tau}{n} \right). \quad (18)$$

*Suppose that the steplength parameter  $\gamma > 0$  satisfies the following two upper bounds:*

$$\gamma \leq \frac{1}{\psi}, \quad \gamma \leq \left( 1 - \frac{1}{\rho} - \frac{2}{\sqrt{n}} \right) \frac{\sqrt{n}L_{\max}}{4L_{\text{res}}\tau\rho^\tau}. \quad (19)$$

*Then we have*

$$\mathbb{E}\|x_{j-1} - \bar{x}_j\|^2 \leq \rho \mathbb{E}\|x_j - \bar{x}_{j+1}\|^2, \quad j = 1, 2, \dots \quad (20)$$

*If the essential strong convexity property (3) holds with  $l > 0$ , we have for  $j = 1, 2, \dots$  that*

$$\begin{aligned} & \mathbb{E}\|x_j - \mathcal{P}_S(x_j)\|^2 + \frac{2\gamma}{L_{\max}} (\mathbb{E}f(x_j) - f^*) \\ & \leq \left( 1 - \frac{l}{n(l + \gamma^{-1}L_{\max})} \right)^j \left( R^2 + \frac{2\gamma}{L_{\max}} (f(x_0) - f^*) \right). \end{aligned} \quad (21)$$

*For general smooth convex function  $f$ , we have*

$$\mathbb{E}f(x_j) - f^* \leq \frac{n(R^2L_{\max} + 2\gamma(f(x_0) - f^*))}{2\gamma(n + j)}. \quad (22)$$

Similarly to the unconstrained case, the following corollary proposes an interesting particular choice for the parameters for which the convergence expressions become more comprehensible. The result requires a condition on the delay bound  $\tau$  in terms of  $n$  and the ratio  $L_{\max}/L_{\text{res}}$ .

**Corollary 5** *Suppose that Assumption 1 holds, that  $\tau \geq 1$  and  $n \geq 5$ , and that*

$$\tau(\tau + 1) \leq \frac{\sqrt{n}L_{\max}}{4eL_{\text{res}}}. \quad (23)$$

*If we choose*

$$\rho = 1 + \frac{4e\tau L_{\text{res}}}{\sqrt{n}L_{\max}}, \quad (24)$$

*then the steplength  $\gamma = 1/2$  will satisfy the bounds (19). In addition, for the essentially strongly convex case (3) with  $l > 0$ , we have for  $j = 1, 2, \dots$  that*

$$\mathbb{E}(f(x_j) - f^*) \leq \left(1 - \frac{l}{n(l + 2L_{\max})}\right)^j (L_{\max}R^2 + f(x_0) - f^*), \quad (25)$$

*while for the case of general convex  $f$ , we have*

$$\mathbb{E}(f(x_j) - f^*) \leq \frac{n(L_{\max}R^2 + f(x_0) - f^*)}{j + n}. \quad (26)$$

Similarly to Section 4, and provided  $\tau$  satisfies (23), the convergence rate is not affected appreciably by the delay bound  $\tau$ , and near-linear speedup can be expected for multicore implementations when (23) holds. This condition is more restrictive than (10) in the unconstrained case, but still holds in many problems for interesting values of  $\tau$ . When  $L_{\text{res}}/L_{\max}$  is bounded independently of dimension, the maximal number of cores allowed is of the order of  $n^{1/4}$ , which is smaller than the  $O(n^{1/2})$  value obtained for the unconstrained case.

We conclude this section with another high-probability bound, whose proof tracks that of Theorem 3.

**Theorem 6** *Suppose that the conditions of Corollary 5 hold, including the definitions of  $\rho$  and  $\psi$ . Then for  $\epsilon > 0$  and  $\eta \in (0, 1)$ , we have that*

$$\mathbb{P}(f(x_j) - f^* \leq \epsilon) \geq 1 - \eta,$$

*provided that one of the following conditions holds: In the essentially strongly convex case (3) with  $l > 0$ , we require*

$$j \geq \frac{n(l + 2L_{\max})}{l} \left\lceil \log \frac{L_{\max}R^2 + f(x_0) - f^*}{\epsilon\eta} \right\rceil,$$

*while in the general convex case, it suffices that*

$$j \geq \frac{n(L_{\max}R^2 + f(x_0) - f^*)}{\epsilon\eta} - n.$$

## 6. Experiments

We illustrate the behavior of two variants of the stochastic coordinate descent approach on test problems constructed from several data sets. Our interests are in the efficiency of multicore implementations (by comparison with a single-threaded implementation) and in performance relative to alternative solvers for the same problems.

All our test problems have the form (1), with either  $\Omega = \mathbb{R}^n$  or  $\Omega$  separable as in (2). The objective  $f$  is quadratic, that is,

$$f(x) = \frac{1}{2}x^T Qx + c^T x,$$

with  $Q$  symmetric positive definite.

Our implementation of ASYSCD is called DIMM-WITTED (or DW for short). It runs on various numbers of threads, from 1 to 40, each thread assigned to a single core in our 40-core Intel Xeon architecture. Cores on the Xeon architecture are arranged into four sockets — ten cores per socket, with each socket having its own memory. Non-uniform memory access (NUMA) means that memory accesses to local memory (on the same socket as the core) are less expensive than accesses to memory on another socket. In our DW implementation, we assign each socket an equal-sized “slice” of  $Q$ , a row submatrix. The components of  $x$  are partitioned between cores, each core being responsible for updating its own partition of  $x$  (though it can read the components of  $x$  from other cores). The components of  $x$  assigned to the cores correspond to the rows of  $Q$  assigned to that core’s socket. Computation is grouped into “epochs,” where an epoch is defined to be the period of computation during which each component of  $x$  is updated exactly once. We use the parameter  $p$  to denote the number of epochs that are executed between reordering (shuffling) of the coordinates of  $x$ . We investigate both shuffling after every epoch ( $p = 1$ ) and after every tenth epoch ( $p = 10$ ). Access to  $x$  is lock-free, and updates are performed asynchronously. This update scheme does not implement exactly the “sampling with replacement” scheme analyzed in previous sections, but can be viewed as a high performance, practical adaptation of the ASYSCD method.

To do each coordinate descent update, a thread must read the latest value of  $x$ . Most components are already in the cache for that core, so that it only needs to fetch those components recently changed. When a thread writes to  $x_i$ , the hardware ensures that this  $x_i$  is simultaneously removed from other cores, signaling that they must fetch the updated version before proceeding with their respective computations.

Although DW is not a precise implementation of ASYSCD, it largely achieves the consistent-read condition that is assumed by the analysis. Inconsistent read happens on a core only if the following three conditions are satisfied simultaneously:

- A core does not finish reading recently changed coordinates of  $x$  (note that it needs to read no more than  $\tau$  coordinates);
- Among these recently changed coordinates, modifications take place both to coordinates that *have been read* and that are *still to be read* by this core;
- Modification of the already-read coordinates happens earlier than the modification of the still-unread coordinates.

Inconsistent read will occur only if at least two coordinates of  $x$  are modified twice during a stretch of approximately  $\tau$  updates to  $x$  (that is, iterations of Algorithm 1). For the DW implementation, inconsistent read would require repeated updating of a particular component in a stretch of approximately  $\tau$  iterations that straddles two epochs. This event would be rare, for typical values of  $n$  and  $\tau$ . Of course, one can avoid the inconsistent read issue altogether by changing the shuffling rule slightly, enforcing the requirement that no coordinate can be modified twice in a span of  $\tau$  iterations. From the practical perspective, this change does not improve performance, and detracts from the simplicity of the approach. From the theoretical perspective, however, the analysis for the inconsistent-read model would be interesting and meaningful, and we plan to study this topic in future work.

The first test problem **QP** is an unconstrained, regularized least squares problem constructed with synthetic data. It has the form

$$\min_{x \in \mathbb{R}^n} f(x) := \frac{1}{2} \|Ax - b\|^2 + \frac{\alpha}{2} \|x\|^2. \quad (27)$$

All elements of  $A \in \mathbb{R}^{m \times n}$ , the true model  $\tilde{x} \in \mathbb{R}^n$ , and the observation noise vector  $\delta \in \mathbb{R}^m$  are generated in i.i.d. fashion from the Gaussian distribution  $\mathcal{N}(0, 1)$ , following which each column in  $A$  is scaled to have a Euclidean norm of 1. The observation  $b \in \mathbb{R}^m$  is constructed from  $A\tilde{x} + \delta\|A\tilde{x}\|/(5m)$ . We choose  $m = 6000$ ,  $n = 20000$ , and  $\alpha = 0.5$ . We therefore have  $L_{\max} = 1 + \alpha = 1.5$  and

$$\frac{L_{\text{res}}}{L_{\max}} \approx \frac{1 + \sqrt{n/m} + \alpha}{1 + \alpha} \approx 2.2.$$

This problem is diagonally dominant, and the condition (10) is satisfied when delay parameter  $\tau$  is less than about 95. In Algorithm 1, we set the steplength parameter  $\gamma$  to 1, and we choose initial iterate to be  $x_0 = \mathbf{0}$ . We measure convergence of the residual norm  $\|\nabla f(x)\|$ .

Our second problem **QPc** is a bound-constrained version of (27):

$$\min_{x \in \mathbb{R}_+^n} f(x) := \frac{1}{2} (x - \tilde{x})^T (A^T A + \alpha I) (x - \tilde{x}). \quad (28)$$

The methodology for generating  $A$  and  $\tilde{x}$  and for choosing the values of  $m$ ,  $n$ ,  $\gamma$ , and  $x_0$  is the same as for (27). We measure convergence via the residual  $\|x - \mathcal{P}_\Omega(x - \nabla f(x))\|$ , where  $\Omega$  is the nonnegative orthant  $\mathbb{R}_+^n$ . At the solution of (28), about half the components of  $x$  are at their lower bound of 0.

Our third and fourth problems are quadratic penalty functions for linear programming relaxations of vertex cover problems on large graphs. The vertex cover problem for an undirected graph with edge set  $E$  and vertex set  $V$  can be written as a binary linear program:

$$\min_{y \in \{0,1\}^{|V|}} \sum_{v \in V} y_v \quad \text{subject to } y_u + y_v \geq 1, \quad \forall (u, v) \in E.$$

By relaxing each binary constraint to the interval  $[0, 1]$ , introducing slack variables for the cover inequalities, we obtain a problem of the form

$$\min_{y_v \in [0,1], s_{uv} \in [0,1]} \sum_{v \in V} y_v \quad \text{subject to } y_u + y_v - s_{uv} = 0, \quad \forall (u, v) \in E.$$

This has the form

$$\min_{x \in [0,1]^n} c^T x \quad \text{subject to } Ax = b,$$

for  $n = |V| + |E|$ . The test problem (29) is a regularized quadratic penalty reformulation of this linear program for some penalty parameter  $\beta$ :

$$\min_{x \in [0,1]^n} c^T x + \frac{\beta}{2} \|Ax - b\|^2 + \frac{1}{2\beta} \|x\|^2, \quad (29)$$

with  $\beta = 5$ . Two test data sets **Amazon** and **DBLP** have dimensions  $n = 561050$  and  $n = 520891$ , respectively.

We tracked the behavior of the residual as a function of the number of epochs, when executed on different numbers of cores. Figure 1 shows convergence behavior for each of our four test problems on various numbers of cores with two different shuffling periods:  $p = 1$  and  $p = 10$ . We note the following points.

- The total amount of computation to achieve any level of precision appears to be almost independent of the number of cores, at least up to 40 cores. In this respect, the performance of the algorithm does not change appreciably as the number of cores is increased. Thus, any deviation from linear speedup is due not to degradation of convergence speed in the algorithm but rather to systems issues in the implementation.
- When we reshuffle after every epoch ( $p = 1$ ), convergence is slightly faster in synthetic unconstrained QP but slightly slower in **Amazon** and **DBLP** than when we do occasional reshuffling ( $p = 10$ ). Overall, the convergence rates with different shuffling periods are comparable in the sense of epochs. However, when the dimension of the variable is large, the shuffling operation becomes expensive, so we would recommend using a large value for  $p$  for large-dimensional problems.

Results for speedup on multicore implementations are shown in Figures 2 and 3 for DW with  $p = 10$ . Speedup is defined as follows:

$$\frac{\text{runtime a single core using DW}}{\text{runtime on } P \text{ cores}}.$$

Near-linear speedup can be observed for the two QP problems with synthetic data. For Problems 3 and 4, speedup is at most 12-14; there are few gains when the number of cores exceeds about 12. We believe that the degradation is due mostly to memory contention. Although these problems have high dimension, the matrix  $Q$  is very sparse (in contrast to the dense  $Q$  for the synthetic data set). Thus, the ratio of computation to data movement / memory access is much lower for these problems, making memory contention effects more significant.

Figures 2 and 3 also show results of a global-locking strategy for the parallel stochastic coordinate descent method, in which the vector  $x$  is locked by a core whenever it performs a read or update. The performance curve for this strategy hugs the horizontal axis; it is not competitive.

Wall clock times required for the four test problems on 1 and 40 cores, to reduce residuals below  $10^{-5}$  are shown in Table 1. (Similar speedups are noted when we use a convergence tolerance looser than  $10^{-5}$ .)

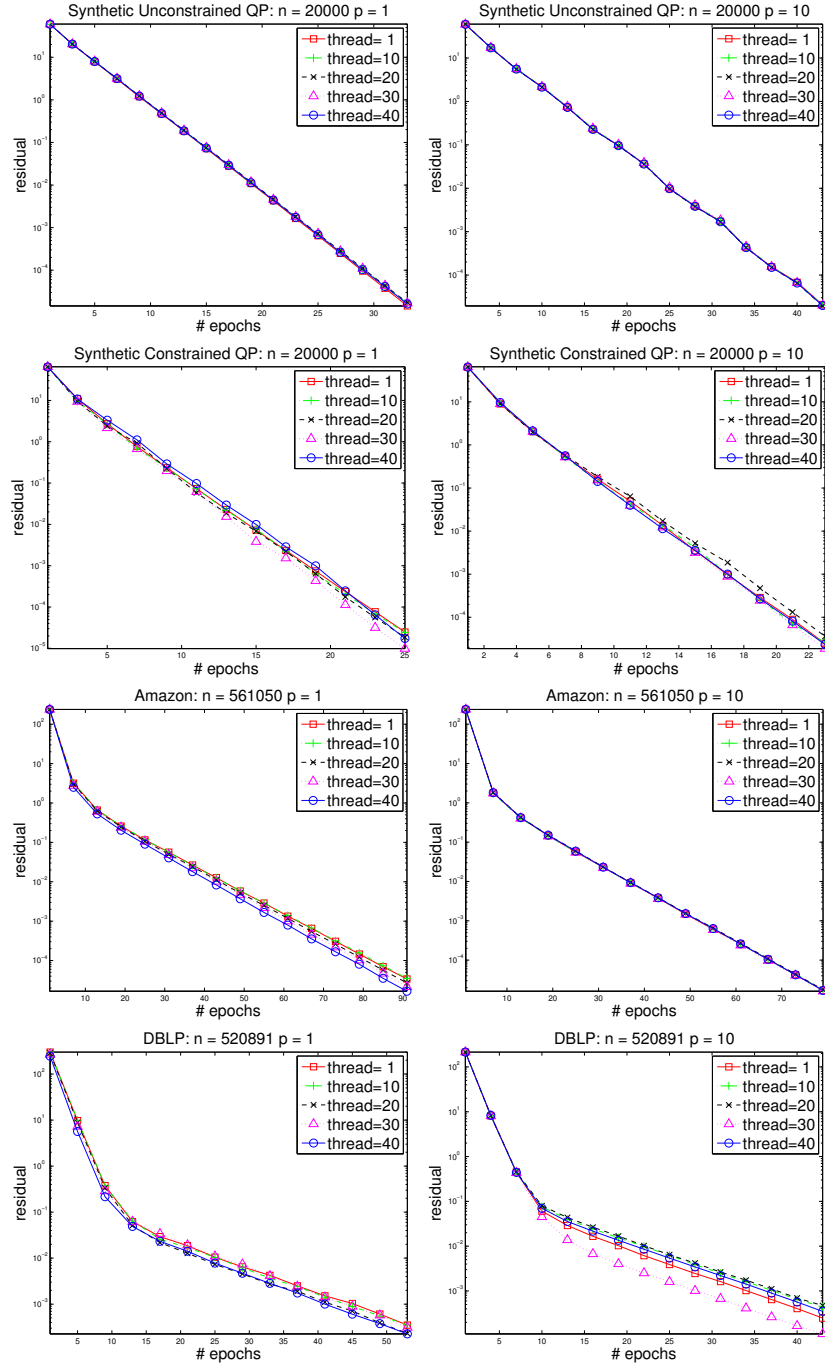


Figure 1: Residuals vs epoch number for the four test problems. Results are reported for variants in which indices are reshuffled after every epoch ( $p = 1$ ) and after every tenth epoch ( $p = 10$ ).

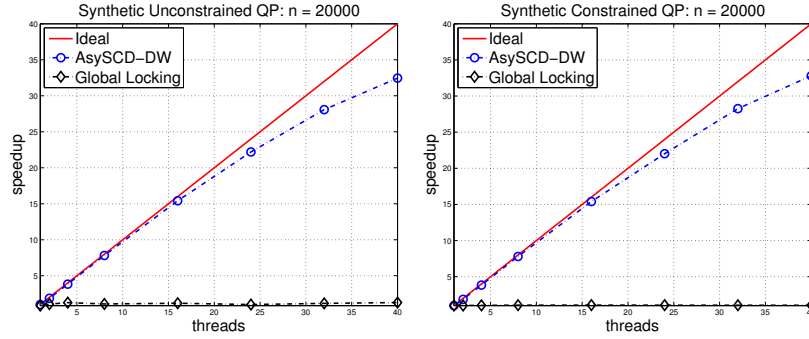


Figure 2: Test problems 1 and 2: Speedup of multicore implementations of DW on up to 40 cores of an Intel Xeon architecture. Ideal (linear) speedup curve is shown for reference, along with poor speedups obtained for a global-locking strategy.

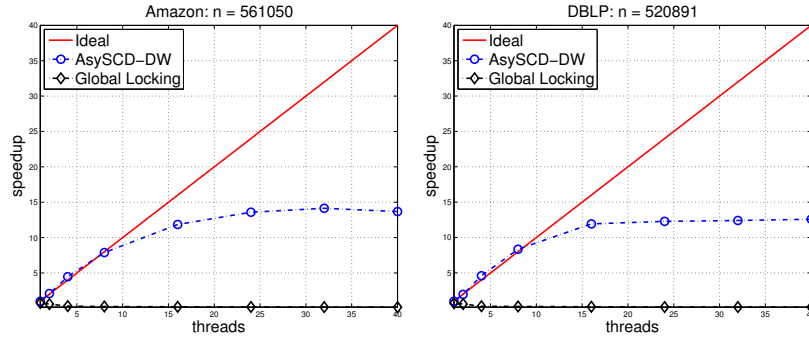


Figure 3: Test problems 3 and 4: Speedup of multicore implementations of DW on up to 40 cores of an Intel Xeon architecture. Ideal (linear) speedup curve is shown for reference, along with poor speedups obtained for a global-locking strategy.

Problem	1 core	40 cores
QP	98.4	3.03
QPc	59.7	1.82
Amazon	17.1	1.25
DBLP	11.5	.91

Table 1: Runtimes (seconds) for the four test problems on 1 and 40 cores.

All problems reported on above are essentially strongly convex. Similar speedup properties can be obtained in the weakly convex case as well. We show speedups for the QPc problem with  $\alpha = 0$ . Table 2 demonstrates similar speedup to the essentially strongly convex case shown in Figure 2.

Turning now to comparisons between ASYSCD and alternative algorithms, we start by considering the basic gradient descent method. We implement gradient descent in a



#cores	Time(sec)	Speedup
1	55.9	1
10	5.19	10.8
20	2.77	20.2
30	2.06	27.2
40	1.81	30.9

Table 2: Runtimes (seconds) and speedup for multicore implementations of DW on different number of cores for the weakly convex **QPc** problem (with  $\alpha = 0$ ) to achieve a residual below 0.06.

#cores	Time(sec)	Speedup
	SYNGD / ASYSCD	SYNGD / ASYSCD
1	121. / 27.1	0.22 / 1.00
10	11.4 / 2.57	2.38 / 10.5
20	6.00 / 1.36	4.51 / 19.9
30	4.44 / 1.01	6.10 / 26.8
40	3.91 / 0.88	6.93 / 30.8

Table 3: Efficiency comparison between SYNGD and ASYSCD for the **QP** problem. The running time and speedup are based on the residual achieving a tolerance of  $10^{-5}$ .

Dataset	# of	# of	Train time(sec)	
	Samples	Features	LIBSVM	ASYSCD
adult	32561	123	16.15	1.39
news	19996	1355191	214.48	7.22
rcv	20242	47236	40.33	16.06
reuters	8293	18930	1.63	0.81
w8a	49749	300	33.62	5.86

Table 4: Efficiency comparison between LIBSVM and ASYSCD for kernel SVM using 40 cores using homogeneous kernels ( $K(x_i, x_j) = (x_i^T x_j)^2$ ). The running time and speedup are calculated based on the “residual”  $10^{-3}$ . Here, to make both algorithms comparable, the “residual” is defined by  $\|x - \mathcal{P}_\Omega(x - \nabla f(x))\|_\infty$ .

parallel, synchronous fashion, distributing the gradient computation load on multiple cores and updates the variable  $x$  in parallel at each step. The resulting implementation is called SYNGD. Table 3 reports running time and speedup of both ASYSCD over SYNGD, showing a clear advantage for ASYSCD.

Next we compare ASYSCD to LIBSVM (Chang and Lin, 2011) a popular parallel solver for kernel support vector machines (SVM). Both algorithms are run on 40 cores to solve the dual formulation of kernel SVM, without an intercept term. All data sets used in 4 except

reuters were obtained from the LIBSVM data set repository.<sup>1</sup> The data set reuters is a sparse binary text classification data set constructed as a one-versus-all version of Reuters-2159.<sup>2</sup> Our comparisons, shown in Table 4, indicate that ASYSCD outperforms LIBSVM on these test sets.

## 7. Extension

The ASYSCD algorithm can be extended by partitioning the coordinates into blocks, and modifying Algorithm 1 to work with these blocks rather than with single coordinates. If  $L_i$ ,  $L_{\max}$ , and  $L_{\text{res}}$  are defined in the block sense, as follows:

$$\begin{aligned} \|\nabla f(x) - \nabla f(x + E_i t)\| &\leq L_{\text{res}} \|t\| \quad \forall x, i, t \in \mathbb{R}^{|i|}, \\ \|\nabla_i f(x) - \nabla_i f(x + E_i t)\| &\leq L_i \|t\| \quad \forall x, i, t \in \mathbb{R}^{|i|}, \\ L_{\max} &= \max_i L_i, \end{aligned}$$

where  $E_i$  is the projection from the  $i$ th block to  $\mathbb{R}^n$  and  $|i|$  denotes the number of components in block  $i$ , our analysis can be extended appropriately.

To make the ASYSCD algorithm more efficient, one can redefine the steplength in Algorithm 1 to be  $\frac{\gamma}{L_{i(j)}}$  rather than  $\frac{\gamma}{L_{\max}}$ . Our analysis can be applied to this variant by doing a change of variables to  $\tilde{x}$ , with  $x_i = \frac{L_i}{L_{\max}} \tilde{x}_i$  and defining  $L_i$ ,  $L_{\text{res}}$ , and  $L_{\max}$  in terms of  $\tilde{x}$ .

## 8. Conclusion

This paper proposes an asynchronous parallel stochastic coordinate descent algorithm for minimizing convex objectives, in the unconstrained and separable-constrained cases. Sub-linear convergence (at rate  $1/K$ ) is proved for general convex functions, with stronger linear convergence results for functions that satisfy an essential strong convexity property. Our analysis indicates the extent to which parallel implementations can be expected to yield near-linear speedup, in terms of a parameter that quantifies the cross-coordinate interactions in the gradient  $\nabla f$  and a parameter  $\tau$  that bounds the delay in updating. Our computational experience confirms the theory.

## Acknowledgments

This project is supported by NSF Grants DMS-0914524, DMS-1216318, and CCF-1356918; NSF CAREER Award IIS-1353606; ONR Awards N00014-13-1-0129 and N00014-12-1-0041; AFOSR Award FA9550-13-1-0138; a Sloan Research Fellowship; and grants from Oracle, Google, and ExxonMobil.

## Appendix A. Proofs for Unconstrained Case

This section contains convergence proofs for ASYSCD in the unconstrained case.

1. <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>
2. <http://www.daviddlewis.com/resources/testcollections/reuters21578/>

We start with a technical result, then move to the proofs of the three main results of Section 4.

**Lemma 7** *For any  $x$ , we have*

$$\|x - \mathcal{P}_S(x)\|^2 \|\nabla f(x)\|^2 \geq (f(x) - f^*)^2.$$

*If the essential strong convexity property (3) holds, we have*

$$\|\nabla f(x)\|^2 \geq 2l(f(x) - f^*).$$

**Proof** The first inequality is proved as follows:

$$f(x) - f^* \leq \langle \nabla f(x), x - \mathcal{P}_S(x) \rangle \leq \|\nabla f(x)\| \|\mathcal{P}_S(x) - x\|.$$

For the second bound, we have from the definition (3), setting  $y \leftarrow x$  and  $x \leftarrow \mathcal{P}_S(x)$ , that

$$\begin{aligned} f^* - f(x) &\geq \langle \nabla f(x), \mathcal{P}_S(x) - x \rangle + \frac{l}{2} \|x - \mathcal{P}_S(x)\|^2 \\ &= \frac{l}{2} \|\mathcal{P}_S(x) - x\|^2 + \frac{1}{l} \|\nabla f(x)\|^2 - \frac{1}{2l} \|\nabla f(x)\|^2, \end{aligned}$$

as required. ■

**Proof** (Theorem 1) We prove each of the two inequalities in (7) by induction. We start with the left-hand inequality. For all values of  $j$ , we have

$$\begin{aligned} &\mathbb{E} (\|\nabla f(x_j)\|^2 - \|\nabla f(x_{j+1})\|^2) \\ &= \mathbb{E} \langle \nabla f(x_j) + \nabla f(x_{j+1}), \nabla f(x_j) - \nabla f(x_{j+1}) \rangle \\ &= \mathbb{E} \langle 2\nabla f(x_j) + \nabla f(x_{j+1}) - \nabla f(x_j), \nabla f(x_j) - \nabla f(x_{j+1}) \rangle \\ &\leq 2\mathbb{E} \langle \nabla f(x_j), \nabla f(x_j) - \nabla f(x_{j+1}) \rangle \\ &\leq 2\mathbb{E} (\|\nabla f(x_j)\| \|\nabla f(x_j) - \nabla f(x_{j+1})\|) \\ &\leq 2L_{\text{res}} \mathbb{E} (\|\nabla f(x_j)\| \|x_j - x_{j+1}\|) \\ &\leq \frac{2L_{\text{res}}\gamma}{L_{\text{max}}} \mathbb{E} (\|\nabla f(x_j)\| \|\nabla_{i(j)} f(x_{k(j)})\|) \\ &\leq \frac{L_{\text{res}}\gamma}{L_{\text{max}}} \mathbb{E} (n^{-1/2} \|\nabla f(x_j)\|^2 + n^{1/2} \|\nabla_{i(j)} f(x_{k(j)})\|^2) \\ &= \frac{L_{\text{res}}\gamma}{L_{\text{max}}} \mathbb{E} (n^{-1/2} \|\nabla f(x_j)\|^2 + n^{1/2} \mathbb{E}_{i(j)} (\|\nabla_{i(j)} f(x_{k(j)})\|^2)) \\ &= \frac{L_{\text{res}}\gamma}{L_{\text{max}}} \mathbb{E} (n^{-1/2} \|\nabla f(x_j)\|^2 + n^{-1/2} \|\nabla f(x_{k(j)})\|^2) \\ &\leq \frac{L_{\text{res}}\gamma}{\sqrt{n}L_{\text{max}}} \mathbb{E} (\|\nabla f(x_j)\|^2 + \|\nabla f(x_{k(j)})\|^2). \end{aligned} \tag{30}$$

We can use this bound to show that the left-hand inequality in (7) holds for  $j = 0$ . By setting  $j = 0$  in (30) and noting that  $k(0) = 0$ , we obtain

$$\mathbb{E} (\|\nabla f(x_0)\|^2 - \|\nabla f(x_1)\|^2) \leq \frac{L_{\text{res}}\gamma}{\sqrt{n}L_{\text{max}}} 2\mathbb{E} (\|\nabla f(x_0)\|^2). \tag{31}$$

From (6b), we have

$$\frac{2L_{\text{res}}\gamma}{\sqrt{n}L_{\text{max}}} \leq \frac{\rho - 1}{\rho^\tau} \leq \frac{\rho - 1}{\rho} = 1 - \rho^{-1},$$

where the second inequality follows from  $\rho > 1$ . By substituting into (31), we obtain  $\rho^{-1}\mathbb{E}(\|\nabla f(x_0)\|^2) \leq \mathbb{E}(\|\nabla f(x_1)\|^2)$ , establishing the result for  $j = 1$ . For the inductive step, we use (30) again, assuming that the left-hand inequality in (7) holds up to stage  $j$ , and thus that

$$\mathbb{E}(\|\nabla f(x_{k(j)})\|^2) \leq \rho^\tau \mathbb{E}(\|\nabla f(x_j)\|^2),$$

provided that  $0 \leq j - k(j) \leq \tau$ , as assumed. By substituting into the right-hand side of (30) again, and using  $\rho > 1$ , we obtain

$$\mathbb{E}(\|\nabla f(x_j)\|^2 - \|\nabla f(x_{j+1})\|^2) \leq \frac{2L_{\text{res}}\gamma\rho^\tau}{\sqrt{n}L_{\text{max}}} \mathbb{E}(\|\nabla f(x_j)\|^2).$$

By substituting (6b) we conclude that the left-hand inequality in (7) holds for all  $j$ .

We now work on the right-hand inequality in (7). For all  $j$ , we have the following:

$$\begin{aligned} & \mathbb{E}(\|\nabla f(x_{j+1})\|^2 - \|\nabla f(x_j)\|^2) \\ &= \mathbb{E}\langle \nabla f(x_j) + \nabla f(x_{j+1}), \nabla f(x_{j+1}) - \nabla f(x_j) \rangle \\ &\leq \mathbb{E}(\|\nabla f(x_j) + \nabla f(x_{j+1})\| \|\nabla f(x_{j+1}) - \nabla f(x_j)\|) \\ &\leq L_{\text{res}} \mathbb{E}(\|\nabla f(x_j) + \nabla f(x_{j+1})\| \|x_j - x_{j+1}\|) \\ &\leq L_{\text{res}} \mathbb{E}((2\|\nabla f(x_j)\| + \|\nabla f(x_{j+1}) - \nabla f(x_j)\|) \|x_j - x_{j+1}\|) \\ &\leq L_{\text{res}} \mathbb{E}(2\|\nabla f(x_j)\| \|x_j - x_{j+1}\| + L_{\text{res}} \|x_j - x_{j+1}\|^2) \\ &\leq L_{\text{res}} \mathbb{E}\left(\frac{2\gamma}{L_{\text{max}}} \|\nabla f(x_j)\| \|\nabla_{i(j)} f(x_{k(j)})\| + \frac{L_{\text{res}}\gamma^2}{L_{\text{max}}^2} \|\nabla_{i(j)} f(x_{k(j)})\|^2\right) \\ &\leq L_{\text{res}} \mathbb{E}\left(\frac{\gamma}{L_{\text{max}}} (n^{-1/2} \|\nabla f(x_j)\|^2 + n^{1/2} \|\nabla_{i(j)} f(x_{k(j)})\|^2) + \frac{L_{\text{res}}\gamma^2}{L_{\text{max}}^2} \|\nabla_{i(j)} f(x_{k(j)})\|^2\right) \\ &= L_{\text{res}} \mathbb{E}\left(\frac{\gamma}{L_{\text{max}}} (n^{-1/2} \|\nabla f(x_j)\|^2 + n^{1/2} \mathbb{E}_{i(j)}(\|\nabla_{i(j)} f(x_{k(j)})\|^2)) + \right. \\ &\quad \left. \frac{L_{\text{res}}\gamma^2}{L_{\text{max}}^2} \mathbb{E}_{i(j)}(\|\nabla_{i(j)} f(x_{k(j)})\|^2)\right) \\ &= L_{\text{res}} \mathbb{E}\left(\frac{\gamma}{L_{\text{max}}} (n^{-1/2} \|\nabla f(x_j)\|^2 + n^{-1/2} \|\nabla f(x_{k(j)})\|^2) + \frac{L_{\text{res}}\gamma^2}{nL_{\text{max}}^2} \|\nabla f(x_{k(j)})\|^2\right) \\ &= \frac{\gamma L_{\text{res}}}{\sqrt{n}L_{\text{max}}} \mathbb{E}(\|\nabla f(x_j)\|^2 + \|\nabla f(x_{k(j)})\|^2) + \frac{\gamma^2 L_{\text{res}}^2}{nL_{\text{max}}^2} \mathbb{E}(\|\nabla f(x_{k(j)})\|^2) \\ &\leq \frac{\gamma L_{\text{res}}}{\sqrt{n}L_{\text{max}}} \mathbb{E}(\|\nabla f(x_j)\|^2) + \left(\frac{\gamma L_{\text{res}}}{\sqrt{n}L_{\text{max}}} + \frac{\gamma L_{\text{res}}^2}{nL_{\text{max}}^2}\right) \mathbb{E}(\|\nabla f(x_{k(j)})\|^2), \end{aligned} \tag{32}$$

where the last inequality is from the observation  $\gamma \leq 1$ . By setting  $j = 0$  in this bound, and noting that  $k(0) = 0$ , we obtain

$$\mathbb{E}(\|\nabla f(x_1)\|^2 - \|\nabla f(x_0)\|^2) \leq \left(\frac{2\gamma L_{\text{res}}}{\sqrt{n}L_{\text{max}}} + \frac{\gamma L_{\text{res}}^2}{nL_{\text{max}}^2}\right) \mathbb{E}(\|\nabla f(x_0)\|^2). \tag{33}$$

By using (6c), we have

$$\frac{2\gamma L_{\text{res}}}{\sqrt{n}L_{\text{max}}} + \frac{\gamma L_{\text{res}}^2}{nL_{\text{max}}^2} = \frac{L_{\text{res}}\gamma}{\sqrt{n}L_{\text{max}}} \left( 2 + \frac{L_{\text{res}}}{\sqrt{n}L_{\text{max}}} \right) \leq \frac{\rho - 1}{\rho^\tau} < \rho - 1,$$

where the last inequality follows from  $\rho > 1$ . By substituting into (33), we obtain  $\mathbb{E}(\|\nabla f(x_1)\|^2) \leq \rho \mathbb{E}(\|\nabla f(x_0)\|^2)$ , so the right-hand bound in (7) is established for  $j = 0$ . For the inductive step, we use (32) again, assuming that the right-hand inequality in (7) holds up to stage  $j$ , and thus that

$$\mathbb{E}(\|\nabla f(x_j)\|^2) \leq \rho^\tau \mathbb{E}(\|\nabla f(x_{k(j)})\|^2),$$

provided that  $0 \leq j - k(j) \leq \tau$ , as assumed. From (32) and the left-hand inequality in (7), we have by substituting this bound that

$$\mathbb{E}(\|\nabla f(x_{j+1})\|^2 - \|\nabla f(x_j)\|^2) \leq \left( \frac{2\gamma L_{\text{res}}\rho^\tau}{\sqrt{n}L_{\text{max}}} + \frac{\gamma L_{\text{res}}^2\rho^\tau}{nL_{\text{max}}^2} \right) \mathbb{E}(\|\nabla f(x_j)\|^2). \quad (34)$$

It follows immediately from (6c) that the term in parentheses in (34) is bounded above by  $\rho - 1$ . By substituting this bound into (34), we obtain  $\mathbb{E}(\|\nabla f(x_{j+1})\|^2) \leq \rho \mathbb{E}(\|\nabla f(x_j)\|^2)$ , as required.

At this point, we have shown that both inequalities in (7) are satisfied for all  $j$ .

Next we prove (8) and (9). Take the expectation of  $f(x_{j+1})$  in terms of  $i(j)$ :

$$\begin{aligned} \mathbb{E}_{i(j)} f(x_{j+1}) &= \mathbb{E}_{i(j)} f \left( x_j - \frac{\gamma}{L_{\text{max}}} e_{i(j)} \nabla_{i(j)} f(x_{k(j)}) \right) \\ &= \frac{1}{n} \sum_{i=1}^n f \left( x_j - \frac{\gamma}{L_{\text{max}}} e_i \nabla_i f(x_{k(j)}) \right) \\ &\leq \frac{1}{n} \sum_{i=1}^n f(x_j) - \frac{\gamma}{L_{\text{max}}} \langle \nabla f(x_j), e_i \nabla_i f(x_{k(j)}) \rangle + \frac{L_i}{2L_{\text{max}}^2} \gamma^2 \|\nabla_i f(x_{k(j)})\|^2 \\ &\leq f(x_j) - \frac{\gamma}{nL_{\text{max}}} \langle \nabla f(x_j), \nabla f(x_{k(j)}) \rangle + \frac{\gamma^2}{2nL_{\text{max}}} \|\nabla f(x_{k(j)})\|^2 \\ &= f(x_j) + \frac{\gamma}{nL_{\text{max}}} \underbrace{\langle \nabla f(x_{k(j)}) - \nabla f(x_j), \nabla f(x_{k(j)}) \rangle}_{T_1} \\ &\quad - \left( \frac{\gamma}{nL_{\text{max}}} - \frac{\gamma^2}{2nL_{\text{max}}} \right) \|\nabla f(x_{k(j)})\|^2. \end{aligned} \quad (35)$$

The second term  $T_1$  is caused by delay. If there is no the delay issue,  $T_1$  should be 0 because of  $\nabla f(x_j) = \nabla f(x_{k(j)})$ . We estimate the upper bound of  $\|\nabla f(x_{k(j)}) - \nabla f(x_j)\|$ :

$$\begin{aligned}
 \|\nabla f(x_{k(j)}) - \nabla f(x_j)\| &\leq \sum_{d=k(j)}^{j-1} \|\nabla f(x_{d+1}) - \nabla f(x_d)\| \\
 &\leq L_{\text{res}} \sum_{d=k(j)}^{j-1} \|x_{d+1} - x_d\| \\
 &= \frac{L_{\text{res}}\gamma}{L_{\text{max}}} \sum_{d=k(j)}^{j-1} \|\nabla_{i(d)} f(x_{k(d)})\|. \tag{36}
 \end{aligned}$$

Then  $\mathbb{E}(|T_1|)$  can be bounded by

$$\begin{aligned}
 \mathbb{E}(|T_1|) &\leq \mathbb{E}(\|\nabla f(x_{k(j)}) - \nabla f(x_j)\| \|\nabla f(x_{k(j)})\|) \\
 &\leq \frac{L_{\text{res}}\gamma}{L_{\text{max}}} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} \|\nabla_{i(d)} f(x_{k(d)})\| \|\nabla f(x_{k(j)})\| \right) \\
 &\leq \frac{L_{\text{res}}\gamma}{2L_{\text{max}}} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} n^{1/2} \|\nabla_{i(d)} f(x_{k(d)})\|^2 + n^{-1/2} \|\nabla f(x_{k(j)})\|^2 \right) \\
 &= \frac{L_{\text{res}}\gamma}{2L_{\text{max}}} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} n^{1/2} \mathbb{E}_{i(d)} (\|\nabla_{i(d)} f(x_{k(d)})\|^2) + n^{-1/2} \|\nabla f(x_{k(j)})\|^2 \right) \\
 &= \frac{L_{\text{res}}\gamma}{2L_{\text{max}}} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} n^{-1/2} \|\nabla f(x_{k(d)})\|^2 + n^{-1/2} \|\nabla f(x_{k(j)})\|^2 \right) \\
 &= \frac{L_{\text{res}}\gamma}{2\sqrt{n}L_{\text{max}}} \sum_{d=k(j)}^{j-1} \mathbb{E}(\|\nabla f(x_{k(d)})\|^2 + \|\nabla f(x_{k(j)})\|^2) \\
 &\leq \frac{\tau \rho^\tau L_{\text{res}}\gamma}{\sqrt{n}L_{\text{max}}} \mathbb{E}(\|\nabla f(x_{k(j)})\|^2) \tag{37}
 \end{aligned}$$

where the second line uses (36), and the final inequality uses the fact for  $d$  between  $k(j)$  and  $j-1$ ,  $k(d)$  lies in the range  $k(j) - \tau$  and  $j-1$ , so we have  $|k(d) - k(j)| \leq \tau$  for all  $d$ .

Taking expectation on both sides of (35) in terms of all random variables, together with (37), we obtain

$$\begin{aligned}
 &\mathbb{E}(f(x_{j+1}) - f^*) \\
 &\leq \mathbb{E}(f(x_j) - f^*) + \frac{\gamma}{nL_{\text{max}}} \mathbb{E}(|T_1|) - \left( \frac{\gamma}{nL_{\text{max}}} - \frac{\gamma^2}{2nL_{\text{max}}} \right) \mathbb{E}(\|\nabla f(x_{k(j)})\|^2) \\
 &\leq \mathbb{E}(f(x_j) - f^*) - \left( \frac{\gamma}{nL_{\text{max}}} - \frac{\tau \rho^\tau L_{\text{res}}\gamma^2}{n^{3/2}L_{\text{max}}^2} - \frac{\gamma^2}{2nL_{\text{max}}} \right) \mathbb{E}(\|\nabla f(x_{k(j)})\|^2) \\
 &= \mathbb{E}(f(x_j) - f^*) - \frac{\gamma}{nL_{\text{max}}} \left( 1 - \frac{\psi}{2}\gamma \right) \mathbb{E}(\|\nabla f(x_{k(j)})\|^2),
 \end{aligned}$$

which (because of (6a)) implies that  $\mathbb{E}(f(x_j) - f^*)$  is monotonically decreasing. From Lemma 7 and the assumption  $\|x_j - \mathcal{P}_S(x_j)\| \leq R$  for all  $j$ , we have

$$\begin{aligned} \|\nabla f(x_{k(j)})\|^2 &\geq \max \left\{ 2l(f(x_{k(j)}) - f^*), \frac{(f(x_{k(j)}) - f^*)^2}{\|x_{k(j)} - \mathcal{P}_S(x_{k(j)})\|^2} \right\} \\ &\geq \max \left\{ 2l(f(x_{k(j)}) - f^*), \frac{(f(x_{k(j)}) - f^*)^2}{R^2} \right\}, \end{aligned}$$

which implies

$$\begin{aligned} \mathbb{E}(\|\nabla f(x_{k(j)})\|^2) &\geq \max \left\{ 2l\mathbb{E}(f(x_{k(j)}) - f^*), \frac{\mathbb{E}(f(x_{k(j)}) - f^*)^2}{R^2} \right\} \\ &\geq \max \left\{ 2l\mathbb{E}(f(x_j) - f^*), \frac{\mathbb{E}(f(x_j) - f^*)^2}{R^2} \right\}. \end{aligned}$$

From the first upper bound  $\|\nabla f(x_{k(j)})\|^2 \geq 2l\mathbb{E}(f(x_j) - f^*)$ , we have

$$\begin{aligned} \mathbb{E}(f(x_{j+1}) - f^*) &\leq \mathbb{E}(f(x_j) - f^*) - \frac{\gamma}{nL_{\max}} \left( 1 - \frac{\psi}{2}\gamma \right) \mathbb{E}(\|\nabla f(x_{k(j)})\|^2) \\ &\leq \left( 1 - \frac{2l\gamma}{nL_{\max}} \left( 1 - \frac{\psi}{2}\gamma \right) \right) \mathbb{E}(f(x_j) - f^*), \end{aligned}$$

from which the linear convergence claim (8) follows by an obvious induction. From the other bound  $\|\nabla f(x_{k(j)})\|^2 \geq \frac{(f(x_{k(j)}) - f^*)^2}{R^2}$ , we have

$$\begin{aligned} \mathbb{E}(f(x_{j+1}) - f^*) &\leq \mathbb{E}(f(x_j) - f^*) - \frac{\gamma}{nL_{\max}} \left( 1 - \frac{\psi}{2}\gamma \right) \mathbb{E}(\|\nabla f(x_{k(j)})\|^2) \\ &\leq \mathbb{E}(f(x_j) - f^*) - \frac{\gamma}{nL_{\max}R^2} \left( 1 - \frac{\psi}{2}\gamma \right) \mathbb{E}((f(x_j) - f^*)^2) \\ &\leq \mathbb{E}(f(x_j) - f^*) - \frac{\gamma}{nL_{\max}R^2} \left( 1 - \frac{\psi}{2}\gamma \right) (\mathbb{E}(f(x_j) - f^*))^2, \end{aligned}$$

where the third line uses the Jensen's inequality  $\mathbb{E}(v^2) \geq (\mathbb{E}(v))^2$ . Defining

$$C := \frac{\gamma}{nL_{\max}R^2} \left( 1 - \frac{\psi}{2}\gamma \right),$$

we have

$$\begin{aligned} \mathbb{E}(f(x_{j+1}) - f^*) &\leq \mathbb{E}(f(x_j) - f^*) - C(\mathbb{E}(f(x_j) - f^*))^2 \\ \Rightarrow \frac{1}{\mathbb{E}(f(x_j) - f^*)} &\leq \frac{1}{\mathbb{E}(f(x_{j+1}) - f^*)} - C \frac{\mathbb{E}(f(x_j) - f^*)}{\mathbb{E}(f(x_{j+1}) - f^*)} \\ \Rightarrow \frac{1}{\mathbb{E}(f(x_{j+1}) - f^*)} - \frac{1}{\mathbb{E}(f(x_j) - f^*)} &\geq C \frac{\mathbb{E}(f(x_j) - f^*)}{\mathbb{E}(f(x_{j+1}) - f^*)} \geq C \\ \Rightarrow \frac{1}{\mathbb{E}(f(x_{j+1}) - f^*)} &\geq \frac{1}{f(x_0) - f^*} + C(j+1) \\ \Rightarrow \mathbb{E}(f(x_{j+1}) - f^*) &\leq \frac{1}{(f(x_0) - f^*)^{-1} + C(j+1)}, \end{aligned}$$

which completes the proof of the sublinear rate (9).  $\blacksquare$

**Proof** (Corollary 2) Note first that for  $\rho$  defined by (11), we have

$$\rho^\tau \leq \rho^{\tau+1} = \left( \left( 1 + \frac{2eL_{\text{res}}}{\sqrt{n}L_{\text{max}}} \right)^{\frac{\sqrt{n}L_{\text{max}}}{2eL_{\text{res}}}} \right)^{\frac{2eL_{\text{res}}(\tau+1)}{\sqrt{n}L_{\text{max}}}} \leq e^{\frac{2eL_{\text{res}}(\tau+1)}{\sqrt{n}L_{\text{max}}}} \leq e,$$

and thus from the definition of  $\psi$  (5) that

$$\psi = 1 + \frac{2\tau\rho^\tau L_{\text{res}}}{\sqrt{n}L_{\text{max}}} \leq 1 + \frac{2\tau e L_{\text{res}}}{\sqrt{n}L_{\text{max}}} \leq 2. \quad (38)$$

We show now that the steplength parameter choice  $\gamma = 1/\psi$  satisfies all the bounds in (6), by showing that the second and third bounds are implied by the first. For the second bound (6b), we have

$$\frac{(\rho - 1)\sqrt{n}L_{\text{max}}}{2\rho^{\tau+1}L_{\text{res}}} \geq \frac{(\rho - 1)\sqrt{n}L_{\text{max}}}{2eL_{\text{res}}} \geq 1 \geq \frac{1}{\psi},$$

where the second inequality follows from (11). For the third bound (6c), we have

$$\frac{(\rho - 1)\sqrt{n}L_{\text{max}}}{L_{\text{res}}\rho^\tau(2 + \frac{L_{\text{res}}}{\sqrt{n}L_{\text{max}}})} = \frac{2eL_{\text{res}}}{L_{\text{res}}\rho^\tau(2 + \frac{L_{\text{res}}}{\sqrt{n}L_{\text{max}}})} \geq \frac{2eL_{\text{res}}}{L_{\text{res}}e(2 + \frac{L_{\text{res}}}{\sqrt{n}L_{\text{max}}})} = \frac{2}{2 + \frac{L_{\text{res}}}{\sqrt{n}L_{\text{max}}}} \geq \frac{1}{\psi}.$$

We can thus set  $\gamma = 1/\psi$ , and by substituting this choice into (8) and using (38), we obtain (12). We obtain (13) by making the same substitution into (9).  $\blacksquare$

**Proof** (Theorem 3) From Markov's inequality, we have

$$\begin{aligned} \mathbb{P}(f(x_j) - f^* \geq \epsilon) &\leq \epsilon^{-1} \mathbb{E}(f(x_j) - f^*) \\ &\leq \epsilon^{-1} \left( 1 - \frac{l}{2nL_{\text{max}}} \right)^j (f(x_0) - f^*) \\ &\leq \epsilon^{-1} (1 - c)^{(1/c) \left| \log \frac{f(x_0) - f^*}{\eta\epsilon} \right|} (f(x_0) - f^*) \quad \text{with } c = l/(2nL_{\text{max}}) \\ &\leq \epsilon^{-1} (f(x_0) - f^*) e^{-\left| \log \frac{f(x_0) - f^*}{\eta\epsilon} \right|} \\ &= \eta e^{\log \frac{(f(x_0) - f^*)}{\eta\epsilon}} e^{-\left| \log \frac{f(x_0) - f^*}{\eta\epsilon} \right|} \\ &\leq \eta, \end{aligned}$$

where the second inequality applies (12), the third inequality uses the definition of  $j$  (15), and the second last inequality uses the inequality  $(1 - c)^{1/c} \leq e^{-1} \forall c \in (0, 1)$ , which proves the essentially strongly convex case. Similarly, the general convex case is proven by

$$\mathbb{P}(f(x_j) - f^* \geq \epsilon) \leq \epsilon^{-1} \mathbb{E}(f(x_j) - f^*) \leq \frac{f(x_0) - f^*}{\epsilon \left( 1 + j \frac{f(x_0) - f^*}{4nL_{\text{max}}R^2} \right)} \leq \eta,$$

where the second inequality uses (13) and the last inequality uses the definition of  $j$  (16).  $\blacksquare$



## Appendix B. Proofs for Constrained Case

We start by introducing notation and proving several preliminary results. Define

$$(\Delta_j)_{i(j)} := (x_j - x_{j+1})_{i(j)}, \quad (39)$$

and formulate the update in Step 4 of Algorithm 1 in the following way:

$$x_{j+1} = \arg \min_{x \in \Omega} \langle \nabla_{i(j)} f(x_{k(j)}), (x - x_j)_{i(j)} \rangle + \frac{L_{\max}}{2\gamma} \|x - x_j\|^2.$$

(Note that  $(x_{j+1})_i = (x_j)_i$  for  $i \neq i(j)$ .) From the optimality condition for this formulation, we have

$$\left\langle (x - x_{j+1})_{i(j)}, \nabla_{i(j)} f(x_{k(j)}) - \frac{L_{\max}}{\gamma} (\Delta_j)_{i(j)} \right\rangle \geq 0, \quad \text{for all } x \in \Omega.$$

This implies in particular that for all  $x \in \Omega$ , we have

$$\langle (\mathcal{P}_S(x) - x_{j+1})_{i(j)}, \nabla_{i(j)} f(x_{k(j)}) \rangle \geq \frac{L_{\max}}{\gamma} \langle (\mathcal{P}_S(x) - x_{j+1})_{i(j)}, (\Delta_j)_{i(j)} \rangle. \quad (40)$$

From the definition of  $L_{\max}$ , and using the notation (39), we have

$$f(x_{j+1}) \leq f(x_j) + \langle \nabla_{i(j)} f(x_j), -(\Delta_j)_{i(j)} \rangle + \frac{L_{\max}}{2} \|(\Delta_j)_{i(j)}\|^2,$$

which indicates that

$$\langle \nabla_{i(j)} f(x_j), (\Delta_j)_{i(j)} \rangle \leq f(x_j) - f(x_{j+1}) + \frac{L_{\max}}{2} \|(\Delta_j)_{i(j)}\|^2. \quad (41)$$

From optimality conditions for this definition, we have

$$\left\langle x - \bar{x}_{j+1}, \nabla f(x_{k(j)}) + \frac{L_{\max}}{\gamma} (\bar{x}_{j+1} - x_j) \right\rangle \geq 0 \quad \forall x \in \Omega. \quad (42)$$

We now define  $\Delta_j := x_j - \bar{x}_{j+1}$ , and note that this definition is consistent with  $(\Delta)_{i(j)}$  defined in (39). It can be seen that

$$\mathbb{E}_{i(j)}(\|x_{j+1} - x_j\|^2) = \frac{1}{n} \|\bar{x}_{j+1} - x_j\|^2.$$

We now proceed to prove the main results of Section 5.

**Proof** (Theorem 4) We prove (20) by induction. First, note that for any vectors  $a$  and  $b$ , we have

$$\|a\|^2 - \|b\|^2 = 2\|a\|^2 - (\|a\|^2 + \|b\|^2) \leq 2\|a\|^2 - 2\langle a, b \rangle \leq 2\langle a, a - b \rangle \leq 2\|a\|\|a - b\|,$$

Thus for all  $j$ , we have

$$\|x_{j-1} - \bar{x}_j\|^2 - \|x_j - \bar{x}_{j+1}\|^2 \leq 2\|x_{j-1} - \bar{x}_j\|\|x_j - \bar{x}_{j+1} - x_{j-1} + \bar{x}_j\|. \quad (43)$$

The second factor in the r.h.s. of (43) is bounded as follows:

$$\begin{aligned}
 & \|x_j - \bar{x}_{j+1} - x_{j-1} + \bar{x}_j\| \\
 &= \left\| x_j - \mathcal{P}_\Omega\left(x_j - \frac{\gamma}{L_{\max}} \nabla f(x_{k(j)})\right) - \left(x_{j-1} - \mathcal{P}_\Omega\left(x_{j-1} - \frac{\gamma}{L_{\max}} \nabla f(x_{k(j-1)})\right)\right) \right\| \\
 &\leq \left\| x_j - \frac{\gamma}{L_{\max}} \nabla f(x_{k(j)}) - \mathcal{P}_\Omega\left(x_j - \frac{\gamma}{L_{\max}} \nabla f(x_{k(j)})\right) - \left(x_{j-1} - \frac{\gamma}{L_{\max}} \nabla f(x_{k(j-1)})\right) \right. \\
 &\quad \left. - \mathcal{P}_\Omega\left(x_{j-1} - \frac{\gamma}{L_{\max}} \nabla f(x_{k(j-1)})\right) \right\| + \frac{\gamma}{L_{\max}} \|\nabla f(x_{k(j-1)}) - \nabla f(x_{k(j)})\| \\
 &\leq \left\| x_j - \frac{\gamma}{L_{\max}} \nabla f(x_{k(j)}) - x_{j-1} + \frac{\gamma}{L_{\max}} \nabla f(x_{k(j-1)}) \right\| \\
 &\quad + \frac{\gamma}{L_{\max}} \|\nabla f(x_{k(j-1)}) - \nabla f(x_{k(j)})\| \\
 &\leq \|x_j - x_{j-1}\| + 2 \frac{\gamma}{L_{\max}} \|\nabla f(x_{k(j)}) - \nabla f(x_{k(j-1)})\| \\
 &\leq \|x_j - x_{j-1}\| + 2 \frac{\gamma}{L_{\max}} \sum_{d=\min\{k(j-1), k(j)\}}^{\max\{k(j-1), k(j)\}-1} \|\nabla f(x_d) - \nabla f(x_{d+1})\| \\
 &\leq \|x_j - x_{j-1}\| + 2 \frac{\gamma L_{\text{res}}}{L_{\max}} \sum_{d=\min\{k(j-1), k(j)\}}^{\max\{k(j-1), k(j)\}-1} \|x_d - x_{d+1}\|, \tag{44}
 \end{aligned}$$

where the first inequality follows by adding and subtracting a term, and the second inequality uses the nonexpansive property of projection:

$$\|(z - \mathcal{P}_\Omega(z)) - (y - \mathcal{P}_\Omega(y))\| \leq \|z - y\|.$$

One can see that  $j - 1 - \tau \leq k(j - 1) \leq j - 1$  and  $j - \tau \leq k(j) \leq j$ , which implies that  $j - 1 - \tau \leq d \leq j - 1$  for each index  $d$  in the summation in (44). It also follows that

$$\max\{k(j - 1), k(j)\} - 1 - \min\{k(j - 1), k(j)\} \leq \tau. \tag{45}$$

We set  $j = 1$ , and note that  $k(0) = 0$  and  $k(1) \leq 1$ . Thus, in this case, we have that the lower and upper limits of the summation in (44) are 0 and 0, respectively. Thus, this summation is vacuous, and we have

$$\|x_1 - \bar{x}_2 + x_0 - \bar{x}_1\| \leq \left(1 + 2 \frac{\gamma L_{\text{res}}}{L_{\max}}\right) \|x_1 - x_0\|,$$

By substituting this bound in (43) and setting  $j = 1$ , we obtain

$$\mathbb{E}(\|x_0 - \bar{x}_1\|^2) - \mathbb{E}(\|x_1 - \bar{x}_2\|^2) \leq \left(2 + 4 \frac{\gamma L_{\text{res}}}{L_{\max}}\right) \mathbb{E}(\|x_1 - x_0\| \|\bar{x}_1 - x_0\|). \tag{46}$$

For any  $j$ , we have

$$\begin{aligned}
\mathbb{E}(\|x_j - x_{j-1}\| \|\bar{x}_j - x_{j-1}\|) &\leq \frac{1}{2} \mathbb{E}(n^{1/2} \|x_j - x_{j-1}\|^2 + n^{-1/2} \|\bar{x}_j - x_{j-1}\|^2) \\
&= \frac{1}{2} \mathbb{E}(n^{1/2} \mathbb{E}_{i(j-1)}(\|x_j - x_{j-1}\|^2) + n^{-1/2} \|\bar{x}_j - x_{j-1}\|^2) \\
&= \frac{1}{2} \mathbb{E}(n^{-1/2} \|\bar{x}_j - x_{j-1}\|^2 + n^{-1/2} \|\bar{x}_j - x_{j-1}\|^2) \\
&= n^{-1/2} \mathbb{E} \|\bar{x}_j - x_{j-1}\|^2.
\end{aligned} \tag{47}$$

Returning to (46), we have

$$\mathbb{E}(\|x_0 - \bar{x}_1\|^2) - \mathbb{E}(\|x_1 - \bar{x}_2\|^2) \leq 2n^{-1/2} \mathbb{E} \|\bar{x}_1 - x_0\|^2$$

which implies that

$$\mathbb{E}(\|x_0 - \bar{x}_1\|^2) \leq \left(1 - \frac{2}{\sqrt{n}} - \frac{4\gamma L_{\text{res}}}{\sqrt{n} L_{\text{max}}}\right)^{-1} \mathbb{E}(\|x_1 - \bar{x}_2\|^2) \leq \rho \mathbb{E}(\|x_1 - \bar{x}_2\|^2).$$

To see the last inequality above, we only need to verify that

$$\gamma \leq \left(1 - \rho^{-1} - \frac{2}{\sqrt{n}}\right) \frac{\sqrt{n} L_{\text{max}}}{4L_{\text{res}}}.$$

This proves that (20) holds for  $j = 1$ .

To take the inductive step, we assume that show that (20) holds up to index  $j - 1$ . We have for  $j - 1 - \tau \leq d \leq j - 2$  that

$$\begin{aligned}
\mathbb{E}(\|x_d - x_{d+1}\| \|\bar{x}_j - x_{j-1}\|) &\leq \frac{1}{2} \mathbb{E}(n^{1/2} \|x_d - x_{d+1}\|^2 + n^{-1/2} \|\bar{x}_j - x_{j-1}\|^2) \\
&= \frac{1}{2} \mathbb{E}(n^{1/2} \mathbb{E}_{i(d)}(\|x_d - x_{d+1}\|^2) + n^{-1/2} \|\bar{x}_j - x_{j-1}\|^2) \\
&= \frac{1}{2} \mathbb{E}(n^{-1/2} \|x_d - \bar{x}_{d+1}\|^2 + n^{-1/2} \|\bar{x}_j - x_{j-1}\|^2) \\
&\leq \frac{1}{2} \mathbb{E}(n^{-1/2} \rho^\tau \|x_{j-1} - \bar{x}_j\|^2 + n^{-1/2} \|\bar{x}_j - x_{j-1}\|^2) \\
&\leq \frac{\rho^\tau}{n^{1/2}} \mathbb{E}(\|\bar{x}_j - x_{j-1}\|^2),
\end{aligned} \tag{48}$$

where the second inequality uses the inductive hypothesis. By substituting (44) into (43) and taking expectation on both sides of (43), we obtain

$$\begin{aligned}
 & \mathbb{E}(\|x_{j-1} - \bar{x}_j\|^2) - \mathbb{E}(\|x_j - \bar{x}_{j+1}\|^2) \\
 & \leq 2\mathbb{E}(\|\bar{x}_j - x_{j-1}\| \|\bar{x}_j - \bar{x}_{j+1} + x_j - x_{j-1}\|) \\
 & \leq 2\mathbb{E} \left( \|\bar{x}_j - x_{j-1}\| \left( \|x_j - x_{j-1}\| + 2 \frac{\gamma L_{\text{res}}}{L_{\text{max}}} \sum_{d=\min\{k(j-1), k(j)\}}^{\max\{k(j-1), k(j)\}-1} \|x_d - x_{d+1}\| \right) \right) \\
 & = 2\mathbb{E}(\|\bar{x}_j - x_{j-1}\| \|x_j - x_{j-1}\|) + \\
 & \quad 4 \frac{\gamma L_{\text{res}}}{L_{\text{max}}} \sum_{d=\min\{k(j-1), k(j)\}}^{\max\{k(j-1), k(j)\}-1} \mathbb{E}(\|\bar{x}_j - x_{j-1}\| \|x_d - x_{d+1}\|) \\
 & \leq n^{-1/2} \left( 2 + \frac{4\gamma L_{\text{res}} \tau \rho^\tau}{L_{\text{max}}} \right) \mathbb{E}(\|x_{j-1} - \bar{x}_j\|^2),
 \end{aligned}$$

where the last line uses (45), (47), and (48). It follows that

$$\mathbb{E}(\|x_{j-1} - \bar{x}_j\|^2) \leq \left( 1 - n^{-1/2} \left( 2 + \frac{4\gamma L_{\text{res}} \tau \rho^\tau}{L_{\text{max}}} \right) \right)^{-1} \mathbb{E}(\|x_j - \bar{x}_{j+1}\|^2) \leq \rho \mathbb{E}(\|x_j - \bar{x}_{j+1}\|^2).$$

To see the last inequality, one only needs to verify that

$$\rho^{-1} \leq 1 - \frac{1}{\sqrt{n}} \left( 2 + \frac{4\gamma L_{\text{res}} \tau \rho^\tau}{L_{\text{max}}} \right) \Leftrightarrow \gamma \leq \left( 1 - \rho^{-1} - \frac{2}{\sqrt{n}} \right) \frac{\sqrt{n} L_{\text{max}}}{4L_{\text{res}} \tau \rho^\tau},$$

and the last inequality is true because of the upper bound of  $\gamma$  in (19). It proves (20).

Next we will show the expectation of objective is monotonically decreasing. We have using the definition (39) that

$$\begin{aligned}
 \mathbb{E}_{i(j)}(f(x_{j+1})) &= n^{-1} \sum_{i=1}^n f(x_j + (\Delta_j)_i) \\
 &\leq n^{-1} \sum_{i=1}^n \left[ f(x_j) + \langle \nabla_i f(x_j), (\bar{x}_{j+1} - x_j)_i \rangle + \frac{L_{\text{max}}}{2} \|(\bar{x}_{j+1} - x_j)_i\|^2 \right] \\
 &= f(x_j) + n^{-1} \left( \langle \nabla f(x_j), \bar{x}_{j+1} - x_j \rangle + \frac{L_{\text{max}}}{2} \|\bar{x}_{j+1} - x_j\|^2 \right) \\
 &= f(x_j) + \frac{1}{n} \left( \langle \nabla f(x_{k(j)}), \bar{x}_{j+1} - x_j \rangle + \frac{L_{\text{max}}}{2} \|\bar{x}_{j+1} - x_j\|^2 \right) + \frac{1}{n} \langle \nabla f(x_j) - \nabla f(x_{k(j)}), \bar{x}_{j+1} - x_j \rangle \\
 &\leq f(x_j) + \frac{1}{n} \left( \frac{L_{\text{max}}}{2} \|\bar{x}_{j+1} - x_j\|^2 - \frac{L_{\text{max}}}{\gamma} \|\bar{x}_{j+1} - x_j\|^2 \right) + \frac{1}{n} \langle \nabla f(x_j) - \nabla f(x_{k(j)}), \bar{x}_{j+1} - x_j \rangle \\
 &= f(x_j) - \left( \frac{1}{\gamma} - \frac{1}{2} \right) \frac{L_{\text{max}}}{n} \|\bar{x}_{j+1} - x_j\|^2 + \frac{1}{n} \langle \nabla f(x_j) - \nabla f(x_{k(j)}), \bar{x}_{j+1} - x_j \rangle, \tag{49}
 \end{aligned}$$

where the second inequality uses (42). Consider the expectation of the last term on the right-hand side of this expression. We have

$$\begin{aligned}
& \mathbb{E} \langle \nabla f(x_j) - \nabla f(x_{k(j)}), \bar{x}_{j+1} - x_j \rangle \\
& \leq \mathbb{E} \|\nabla f(x_j) - \nabla f(x_{k(j)})\| \|\bar{x}_{j+1} - x_j\| \\
& \leq \mathbb{E} \sum_{d=k(j)}^{j-1} \|\nabla f(x_d) - \nabla f(x_{d+1})\| \|\bar{x}_{j+1} - x_j\| \\
& \leq L_{\text{res}} \mathbb{E} \sum_{d=k(j)}^{j-1} \|x_d - x_{d+1}\| \|\bar{x}_{j+1} - x_j\| \\
& \leq \frac{L_{\text{res}}}{2} \mathbb{E} \sum_{d=k(j)}^{j-1} (n^{1/2} \|x_d - x_{d+1}\|^2 + n^{-1/2} \|\bar{x}_{j+1} - x_j\|^2) \\
& = \frac{L_{\text{res}}}{2} \mathbb{E} \sum_{d=k(j)}^{j-1} (n^{1/2} \mathbb{E}_{i(d)}(\|x_d - x_{d+1}\|^2) + n^{-1/2} \|\bar{x}_{j+1} - x_j\|^2) \\
& = \frac{L_{\text{res}}}{2} \mathbb{E} \sum_{d=k(j)}^{j-1} (n^{-1/2} \|x_d - \bar{x}_{d+1}\|^2 + n^{-1/2} \|\bar{x}_{j+1} - x_j\|^2) \\
& \leq \frac{L_{\text{res}}}{2n^{1/2}} \mathbb{E} \sum_{d=k(j)}^{j-1} (1 + \rho^\tau) \|\bar{x}_{j+1} - x_j\|^2 \\
& \leq \frac{L_{\text{res}} \tau \rho^\tau}{n^{1/2}} \mathbb{E} \|\bar{x}_{j+1} - x_j\|^2, \tag{50}
\end{aligned}$$

where the fifth inequality uses (20). By taking expectation on both sides of (49) and substituting (50), we have

$$\mathbb{E}(f(x_{j+1})) \leq \mathbb{E}(f(x_j)) - \frac{1}{n} \left( \left( \frac{1}{\gamma} - \frac{1}{2} \right) L_{\text{max}} - \frac{L_{\text{res}} \tau \rho^\tau}{n^{1/2}} \right) \mathbb{E} \|\bar{x}_{j+1} - x_j\|^2.$$

To see  $\left( \frac{1}{\gamma} - \frac{1}{2} \right) L_{\text{max}} - \frac{L_{\text{res}} \tau \rho^\tau}{n^{1/2}} \geq 0$ , we only need to verify

$$\gamma \leq \left( \frac{1}{2} + \frac{L_{\text{res}} \tau \rho^\tau}{\sqrt{n} L_{\text{max}}} \right)^{-1}$$

which is implied by the first upper bound of  $\gamma$  (19). Therefore, we have proved the monotonicity  $\mathbb{E}(f(x_{j+1})) \leq \mathbb{E}(f(x_j))$ .

Next we prove the sublinear convergence rate for the constrained smooth convex case in (22). We have

$$\begin{aligned}
 \|x_{j+1} - \mathcal{P}_S(x_{j+1})\|^2 &\leq \|x_{j+1} - \mathcal{P}_S(x_j)\|^2 \\
 &= \|x_j - (\Delta_j)_{i(j)}e_{i(j)} - \mathcal{P}_S(x_j)\|^2 \\
 &= \|x_j - \mathcal{P}_S(x_j)\|^2 + |(\Delta_j)_{i(j)}|^2 - 2(x_j - \mathcal{P}_S(x_j))_{i(j)}(\Delta_j)_{i(j)} \\
 &= \|x_j - \mathcal{P}_S(x_j)\|^2 - |(\Delta_j)_{i(j)}|^2 - 2((x_j - \mathcal{P}_S(x_j))_{i(j)} - (\Delta_j)_{i(j)})(\Delta_j)_{i(j)} \\
 &= \|x_j - \mathcal{P}_S(x_j)\|^2 - \|(\Delta_j)_{i(j)}\|^2 + 2(\mathcal{P}_S(x_j) - x_{j+1})_{i(j)}(\Delta_j)_{i(j)} \\
 &\leq \|x_j - \mathcal{P}_S(x_j)\|^2 - |(\Delta_j)_{i(j)}|^2 + \frac{2\gamma}{L_{\max}}(\mathcal{P}_S(x_j) - x_{j+1})_{i(j)}\nabla_{i(j)}f(x_{k(j)}) \\
 &= \|x_j - \mathcal{P}_S(x_j)\|^2 - |(\Delta_j)_{i(j)}|^2 + \frac{2\gamma}{L_{\max}}(\mathcal{P}_S(x_j) - x_j)_{i(j)}\nabla_{i(j)}f(x_{k(j)}) + \\
 &\quad \frac{2\gamma}{L_{\max}}((\Delta_j)_{i(j)}\nabla_{i(j)}f(x_j) + (\Delta_j)_{i(j)}(\nabla_{i(j)}f(x_{k(j)}) - \nabla_{i(j)}f(x_j))) \\
 &\leq \|x_j - \mathcal{P}_S(x_j)\|^2 - |(\Delta_j)_{i(j)}|^2 + \frac{2\gamma}{L_{\max}}(\mathcal{P}_S(x_j) - x_j)_{i(j)}\nabla_{i(j)}f(x_{k(j)}) + \\
 &\quad \frac{2\gamma}{L_{\max}}\left(f(x_j) - f(x_{j+1}) + \frac{L_{\max}}{2}|(\Delta_j)_{i(j)}|^2 \right. \\
 &\quad \left. + (\Delta_j)_{i(j)}(\nabla_{i(j)}f(x_{k(j)}) - \nabla_{i(j)}f(x_j))\right) \\
 &= \|x_j - \mathcal{P}_S(x_j)\|^2 - (1 - \gamma)|(\Delta_j)_{i(j)}|^2 + \frac{2\gamma}{L_{\max}}\underbrace{(\mathcal{P}_S(x_j) - x_j)_{i(j)}\nabla_{i(j)}f(x_{k(j)})}_{T_1} + \\
 &\quad \frac{2\gamma}{L_{\max}}(f(x_j) - f(x_{j+1})) + \frac{2\gamma}{L_{\max}}\underbrace{(\Delta_j)_{i(j)}(\nabla_{i(j)}f(x_{k(j)}) - \nabla_{i(j)}f(x_j))}_{T_2}, \tag{51}
 \end{aligned}$$

where the second inequality uses (40) and the third inequality uses (41). We now seek upper bounds on the quantities  $T_1$  and  $T_2$  in the expectation sense. For  $T_1$ , we have

$$\begin{aligned}
 \mathbb{E}(T_1) &= n^{-1} \mathbb{E} \langle \mathcal{P}_S(x_j) - x_j, \nabla f(x_{k(j)}) \rangle \\
 &= n^{-1} \mathbb{E} \langle \mathcal{P}_S(x_j) - x_{k(j)}, \nabla f(x_{k(j)}) \rangle + n^{-1} \mathbb{E} \sum_{d=k(j)}^{j-1} \langle x_d - x_{d+1}, \nabla f(x_{k(j)}) \rangle \\
 &= n^{-1} \mathbb{E} \langle \mathcal{P}_S(x_j) - x_{k(j)}, \nabla f(x_{k(j)}) \rangle \\
 &\quad + n^{-1} \mathbb{E} \sum_{d=k(j)}^{j-1} \langle x_d - x_{d+1}, \nabla f(x_d) \rangle + \langle x_d - x_{d+1}, \nabla f(x_{k(j)}) - \nabla f(x_d) \rangle \\
 &\leq n^{-1} \mathbb{E}(f^* - f(x_{k(j)})) + n^{-1} \mathbb{E} \sum_{d=k(j)}^{j-1} \left( f(x_d) - f(x_{d+1}) + \frac{L_{\max}}{2} \|x_d - x_{d+1}\|^2 \right) \\
 &\quad + n^{-1} \mathbb{E} \sum_{d=k(j)}^{j-1} \langle x_d - x_{d+1}, \nabla f(x_{k(j)}) - \nabla f(x_d) \rangle \\
 &= n^{-1} \mathbb{E}(f^* - f(x_j)) + \frac{L_{\max}}{2n} \mathbb{E} \sum_{d=k(j)}^{j-1} \|x_d - x_{d+1}\|^2 \\
 &\quad + n^{-1} \mathbb{E} \sum_{d=k(j)}^{j-1} \langle x_d - x_{d+1}, \nabla f(x_{k(j)}) - \nabla f(x_d) \rangle \\
 &= n^{-1} \mathbb{E}(f^* - f(x_j)) + \frac{L_{\max}}{2n^2} \mathbb{E} \sum_{d=k(j)}^{j-1} \|x_d - \bar{x}_{d+1}\|^2 \\
 &\quad + n^{-1} \mathbb{E} \sum_{d=k(j)}^{j-1} \langle x_d - x_{d+1}, \nabla f(x_{k(j)}) - \nabla f(x_d) \rangle \\
 &\leq n^{-1} \mathbb{E}(f^* - f(x_j)) + \frac{L_{\max} \tau \rho^\tau}{2n^2} \mathbb{E} \|x_j - \bar{x}_{j+1}\|^2 \\
 &\quad + n^{-1} \sum_{d=k(j)}^{j-1} \underbrace{\mathbb{E} \langle x_d - x_{d+1}, \nabla f(x_{k(j)}) - \nabla f(x_d) \rangle}_{T_3},
 \end{aligned}$$

where the last inequality uses (20). The upper bound of  $\mathbb{E}(T_3)$  is estimated by

$$\begin{aligned}
 \mathbb{E}(T_3) &= \mathbb{E}\langle x_d - x_{d+1}, \nabla f(x_{k(j)}) - \nabla f(x_d) \rangle \\
 &= \mathbb{E}(\mathbb{E}_{i(d)}\langle x_d - x_{d+1}, \nabla f(x_{k(j)}) - \nabla f(x_d) \rangle) \\
 &= n^{-1} \mathbb{E}\langle x_d - \bar{x}_{d+1}, \nabla f(x_{k(j)}) - \nabla f(x_d) \rangle \\
 &\leq n^{-1} \mathbb{E}\|x_d - \bar{x}_{d+1}\| \|\nabla f(x_{k(j)}) - \nabla f(x_d)\| \\
 &\leq n^{-1} \mathbb{E}(\|x_d - \bar{x}_{d+1}\| \sum_{t=k(j)}^{d-1} \|\nabla f(x_t) - \nabla f(x_{t+1})\|) \\
 &\leq \frac{L_{\text{res}}}{n} \sum_{t=k(j)}^{d-1} \mathbb{E}(\|x_d - \bar{x}_{d+1}\| \|x_t - x_{t+1}\|) \\
 &\leq \frac{L_{\text{res}}}{2n} \sum_{t=k(j)}^{d-1} \mathbb{E}(n^{-1/2} \|x_d - \bar{x}_{d+1}\|^2 + n^{1/2} \|x_t - x_{t+1}\|^2) \\
 &\leq \frac{L_{\text{res}}}{2n} \sum_{t=k(j)}^{d-1} \mathbb{E}(n^{-1/2} \|x_d - \bar{x}_{d+1}\|^2 + n^{-1/2} \|x_t - \bar{x}_{t+1}\|^2) \\
 &\leq \frac{L_{\text{res}} \rho^\tau}{n^{3/2}} \sum_{t=k(j)}^{d-1} \mathbb{E}(\|x_j - \bar{x}_{j+1}\|^2) \\
 &\leq \frac{L_{\text{res}} \tau \rho^\tau}{n^{3/2}} \mathbb{E}(\|x_j - \bar{x}_{j+1}\|^2).
 \end{aligned}$$

where the second last inequality uses (20). Therefore,  $\mathbb{E}(T_1)$  can be bounded by

$$\begin{aligned}
 \mathbb{E}(T_1) &= \mathbb{E}\langle (\mathcal{P}_S(x_j) - x_j)_{i(j)}, \nabla_{i(j)} f(x_{k(j)}) \rangle \\
 &\leq \frac{1}{n} \mathbb{E}(f^* - f(x_j)) + \frac{L_{\max} \tau \rho^\tau}{2n^2} \mathbb{E}\|x_j - \bar{x}_{j+1}\|^2 + \sum_{d=k(j)}^{j-1} \frac{L_{\text{res}} \tau \rho^\tau}{n^{5/2}} \mathbb{E}(\|x_j - \bar{x}_{j+1}\|^2) \\
 &= \frac{1}{n} \left( f^* - \mathbb{E}f(x_j) + \left( \frac{L_{\max} \tau \rho^\tau}{2n} + \frac{L_{\text{res}} \tau^2 \rho^\tau}{n^{3/2}} \right) \mathbb{E}(\|x_j - \bar{x}_{j+1}\|^2) \right). \tag{52}
 \end{aligned}$$



For  $T_2$ , we have

$$\begin{aligned}
\mathbb{E}(T_2) &= \mathbb{E}(\Delta_j)_{i(j)} (\nabla_{i(j)} f(x_{k(j)}) - \nabla_{i(j)} f(x_j)) \\
&= n^{-1} \mathbb{E} \langle \Delta_j, \nabla f(x_{k(j)}) - \nabla f(x_j) \rangle \\
&\leq n^{-1} \mathbb{E} (\|\Delta_j\| \|\nabla f(x_{k(j)}) - \nabla f(x_j)\|) \\
&\leq n^{-1} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} \|\Delta_j\| \|\nabla f(x_d) - \nabla f(x_{d+1})\| \right) \\
&\leq \frac{L_{\text{res}}}{n} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} \|\Delta_j\| \|x_d - x_{d+1}\| \right) \\
&= \frac{L_{\text{res}}}{2n} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} n^{-1/2} \|\Delta_j\|^2 + n^{1/2} \|x_d - x_{d+1}\|^2 \right) \\
&= \frac{L_{\text{res}}}{2n} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} n^{-1/2} \|x_j - \bar{x}_{j+1}\|^2 + n^{1/2} \mathbb{E}_{i(d)} \|x_d - x_{d+1}\|^2 \right) \\
&= \frac{L_{\text{res}}}{2n} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} n^{-1/2} \|x_j - \bar{x}_{j+1}\|^2 + n^{-1/2} \|x_d - \bar{x}_{d+1}\|^2 \right) \\
&= \frac{L_{\text{res}}}{2n^{3/2}} \mathbb{E} \left( \sum_{d=k(j)}^{j-1} \mathbb{E} \|x_j - \bar{x}_{j+1}\|^2 + \mathbb{E} \|x_d - \bar{x}_{d+1}\|^2 \right) \\
&\leq \frac{L_{\text{res}}(1 + \rho^\tau)}{2n^{3/2}} \sum_{d=k(j)}^{j-1} \mathbb{E} \|x_j - \bar{x}_{j+1}\|^2 \\
&\leq \frac{L_{\text{res}} \tau \rho^\tau}{n^{3/2}} \mathbb{E} \|x_j - \bar{x}_{j+1}\|^2, \tag{53}
\end{aligned}$$

where the second last inequality uses (20).

By taking the expectation on both sides of (51), using  $\mathbb{E}_{i(j)}(|(\Delta_j)_{i(j)}|^2) = n^{-1} \|x_j - \bar{x}_{j+1}\|^2$ , and substituting the upper bounds from (52) and (53), we obtain

$$\begin{aligned}
\mathbb{E} \|x_{j+1} - \mathcal{P}_S(x_{j+1})\|^2 &\leq \mathbb{E} \|x_j - \mathcal{P}_S(x_j)\|^2 \\
&\quad - \frac{1}{n} \left( 1 - \gamma - \frac{2\gamma L_{\text{res}} \tau \rho^\tau}{L_{\text{max}} n^{1/2}} - \frac{\gamma \tau \rho^\tau}{n} - \frac{2\gamma L_{\text{res}} \tau^2 \rho^\tau}{L_{\text{max}} n^{3/2}} \right) \mathbb{E} \|x_j - \bar{x}_{j+1}\|^2 \\
&\quad + \frac{2\gamma}{L_{\text{max}} n} (f^* - \mathbb{E} f(x_j)) + \frac{2\gamma}{L_{\text{max}}} (\mathbb{E} f(x_j) - \mathbb{E} f(x_{j+1})) \\
&\leq \mathbb{E} \|x_j - \mathcal{P}_S(x_j)\|^2 + \frac{2\gamma}{L_{\text{max}} n} (f^* - \mathbb{E} f(x_j)) + \frac{2\gamma}{L_{\text{max}}} (\mathbb{E} f(x_j) - \mathbb{E} f(x_{j+1})). \tag{54}
\end{aligned}$$

In the second inequality, we were able to drop the term involving  $\mathbb{E} \|x_j - \bar{x}_{j+1}\|^2$  by using the fact that

$$1 - \gamma - \frac{2\gamma L_{\text{res}} \tau \rho^\tau}{L_{\text{max}} n^{1/2}} - \frac{\gamma \tau \rho^\tau}{n} - \frac{2\gamma L_{\text{res}} \tau^2 \rho^\tau}{L_{\text{max}} n^{3/2}} = 1 - \gamma \psi \geq 0,$$

which follows from the definition (18) of  $\psi$  and from the first upper bound on  $\gamma$  in (19). It follows that

$$\begin{aligned}
 & \mathbb{E}\|x_{j+1} - \mathcal{P}_S(x_{j+1})\|^2 + \frac{2\gamma}{L_{\max}}(\mathbb{E}f(x_{j+1}) - f^*) \\
 & \leq \mathbb{E}\|x_j - \mathcal{P}_S(x_j)\|^2 + \frac{2\gamma}{L_{\max}}(\mathbb{E}f(x_j) - f^*) - \frac{2\gamma}{L_{\max}n}(\mathbb{E}f(x_j) - f^*) \\
 & \leq \|x_0 - \mathcal{P}_S(x_0)\|^2 + \frac{2\gamma}{L_{\max}}(f(x_0) - f^*) - \frac{2\gamma}{L_{\max}n} \sum_{t=0}^j (\mathbb{E}f(x_t) - f^*) \\
 & \leq R^2 + \frac{2\gamma}{L_{\max}}(f(x_0) - f^*) - \frac{2\gamma(j+1)}{L_{\max}n}(\mathbb{E}f(x_{j+1}) - f^*),
 \end{aligned} \tag{55}$$

where the second inequality follows by applying induction to the inequality

$$S_{j+1} \leq S_j - \frac{2\gamma}{L_{\max}n} \mathbb{E}(f(x_j) - f^*),$$

where

$$S_j := \mathbb{E}(\|x_j - \mathcal{P}_S(x_j)\|^2) + \frac{2\gamma}{L_{\max}} \mathbb{E}(f(x_j) - \mathcal{P}_S(x_j)),$$

and the last line uses the monotonicity of  $\mathbb{E}f(x_j)$  (proved above) and the assumed bound  $\|x_0 - \mathcal{P}_S(x_0)\| \leq R$ . It implies that

$$\begin{aligned}
 & \mathbb{E}\|x_{j+1} - \mathcal{P}_S(x_{j+1})\|^2 + \frac{2\gamma}{L_{\max}}(\mathbb{E}f(x_{j+1}) - f^*) + \frac{2\gamma(j+1)}{L_{\max}n}(\mathbb{E}f(x_{j+1}) - f^*) \\
 & \leq R^2 + \frac{2\gamma}{L_{\max}}(f(x_0) - f^*) \\
 \Rightarrow & \frac{2\gamma(n+j+1)}{L_{\max}n}(\mathbb{E}f(x_{j+1}) - f^*) \leq R^2 + \frac{2\gamma}{L_{\max}}(f(x_0) - f^*) \\
 \Rightarrow & \mathbb{E}f(x_{j+1}) - f^* \leq \frac{n(R^2 L_{\max} + 2\gamma(f(x_0) - f^*))}{2\gamma(n+j+1)}.
 \end{aligned}$$

This completes the proof of the sublinear convergence rate (22).

Finally, we prove the linear convergence rate (21) for the essentially strongly convex case. All bounds proven above hold, and we make use the following additional property:

$$f(x_j) - f^* \geq \langle \nabla f(\mathcal{P}_S(x_j)), x_j - \mathcal{P}_S(x_j) \rangle + \frac{l}{2} \|x_j - \mathcal{P}_S(x_j)\|^2 \geq \frac{l}{2} \|x_j - \mathcal{P}_S(x_j)\|^2,$$

due to feasibility of  $x_j$  and  $\langle \nabla f(\mathcal{P}_S(x_j)), x_j - \mathcal{P}_S(x_j) \rangle \geq 0$ . By using this result together with some elementary manipulation, we obtain

$$\begin{aligned}
 f(x_j) - f^* &= \left(1 - \frac{L_{\max}}{l\gamma + L_{\max}}\right) (f(x_j) - f^*) + \frac{L_{\max}}{l\gamma + L_{\max}} (f(x_j) - f^*) \\
 &\geq \left(1 - \frac{L_{\max}}{l\gamma + L_{\max}}\right) (f(x_j) - f^*) + \frac{L_{\max}l}{2(l\gamma + L_{\max})} \|x_j - \mathcal{P}_S(x_j)\|^2 \\
 &= \frac{L_{\max}l}{2(l\gamma + L_{\max})} \left( \|x_j - \mathcal{P}_S(x_j)\|^2 + \frac{2\gamma}{L_{\max}} (f(x_j) - f^*) \right).
 \end{aligned} \tag{56}$$

Recalling (55), we have

$$\begin{aligned} & \mathbb{E}\|x_{j+1} - \mathcal{P}_S(x_{j+1})\|^2 + \frac{2\gamma}{L_{\max}}(\mathbb{E}f(x_{j+1}) - f^*) \\ & \leq \mathbb{E}\|x_j - \mathcal{P}_S(x_j)\|^2 + \frac{2\gamma}{L_{\max}}(\mathbb{E}f(x_j) - f^*) - \frac{2\gamma}{L_{\max}n}(\mathbb{E}f(x_j) - f^*). \end{aligned} \quad (57)$$

By taking the expectation of both sides in (56) and substituting in the last term of (57), we obtain

$$\begin{aligned} & \mathbb{E}\|x_{j+1} - \mathcal{P}_S(x_{j+1})\|^2 + \frac{2\gamma}{L_{\max}}(\mathbb{E}f(x_{j+1}) - f^*) \\ & \leq \mathbb{E}\|x_j - \mathcal{P}_S(x_j)\|^2 + \frac{2\gamma}{L_{\max}}(\mathbb{E}f(x_j) - f^*) \\ & \quad - \frac{2\gamma}{L_{\max}n} \left( \frac{L_{\max}l}{2(l\gamma + L_{\max})} \left( \mathbb{E}\|x_j - \mathcal{P}_S(x_j)\|^2 + \frac{2\gamma}{L_{\max}}(\mathbb{E}f(x_j) - f^*) \right) \right) \\ & = \left( 1 - \frac{l}{n(l + \gamma^{-1}L_{\max})} \right) \left( \mathbb{E}\|x_j - \mathcal{P}_S(x_j)\|^2 + \frac{2\gamma}{L_{\max}}(\mathbb{E}f(x_j) - f^*) \right) \\ & \leq \left( 1 - \frac{l}{n(l + \gamma^{-1}L_{\max})} \right)^{j+1} \left( \|x_0 - \mathcal{P}_S(x_0)\|^2 + \frac{2\gamma}{L_{\max}}(f(x_0) - f^*) \right), \end{aligned}$$

which yields (21). ■

**Proof** (Corollary 5) To apply Theorem 4, we first show  $\rho > \left(1 - \frac{2}{\sqrt{n}}\right)^{-1}$ . Using the bound (23), together with  $L_{\text{res}}/L_{\max} \geq 1$ , we obtain

$$\begin{aligned} & \left(1 - \frac{2}{\sqrt{n}}\right) \left(1 + \frac{4e\tau L_{\text{res}}}{\sqrt{n}L_{\max}}\right) = \left(1 + \frac{4e\tau L_{\text{res}}}{\sqrt{n}L_{\max}}\right) - \left(1 + \frac{4e\tau L_{\text{res}}}{\sqrt{n}L_{\max}}\right) \frac{2}{\sqrt{n}} \\ & \geq \left(1 + \frac{4e\tau}{\sqrt{n}}\right) - \left(1 + \frac{1}{\tau + 1}\right) \frac{2}{\sqrt{n}} = 1 + \left(2e\tau - 1 - \frac{1}{\tau + 1}\right) \frac{2}{\sqrt{n}} > 1, \end{aligned}$$

where the last inequality uses  $\tau \geq 1$ . Note that for  $\rho$  defined by (24), and using (23), we have

$$\rho^\tau \leq \rho^{\tau+1} = \left( \left(1 + \frac{4e\tau L_{\text{res}}}{\sqrt{n}L_{\max}}\right)^{\frac{\sqrt{n}L_{\max}}{4e\tau L_{\text{res}}}} \right)^{\frac{4e\tau L_{\text{res}}(\tau+1)}{\sqrt{n}L_{\max}}} \leq e^{\frac{4e\tau L_{\text{res}}(\tau+1)}{\sqrt{n}L_{\max}}} \leq e.$$

Thus from the definition of  $\psi$  (18), we have that

$$\begin{aligned} \psi & = 1 + \frac{L_{\text{res}}\tau\rho^\tau}{\sqrt{n}L_{\max}} \left( 2 + \frac{L_{\max}}{\sqrt{n}L_{\text{res}}} + \frac{2\tau}{n} \right) \leq 1 + \frac{L_{\text{res}}\tau\rho^\tau}{4eL_{\text{res}}\tau(\tau+1)} \left( 2 + \frac{1}{\sqrt{n}} + \frac{2\tau}{n} \right) \\ & \leq 1 + \frac{1}{4(\tau+1)} \left( 2 + \frac{1}{\sqrt{n}} + \frac{2\tau}{n} \right) \leq 1 + \left( \frac{1}{4} + \frac{1}{16} + \frac{1}{10} \right) \leq 2. \end{aligned} \quad (58)$$

(The second last inequality uses  $n \geq 5$  and  $\tau \geq 1$ .) Thus, the steplength parameter choice  $\gamma = 1/2$  satisfies the first bound in (19). To show that the second bound in (19) holds also,

we have

$$\begin{aligned} \left(1 - \frac{1}{\rho} - \frac{2}{\sqrt{n}}\right) \frac{\sqrt{n}L_{\max}}{4L_{\text{res}}\tau\rho^{\tau}} &= \left(\frac{\rho-1}{\rho} - \frac{2}{\sqrt{n}}\right) \frac{\sqrt{n}L_{\max}}{4L_{\text{res}}\tau\rho^{\tau}} \\ &= \frac{4e\tau L_{\text{res}}}{4L_{\text{res}}\tau\rho^{\tau+1}} - \frac{L_{\max}}{2L_{\text{res}}\tau\rho^{\tau}} \geq 1 - \frac{1}{2} = \frac{1}{2}. \end{aligned}$$

We can thus set  $\gamma = 1/2$ , and by substituting this choice into (21), we obtain (25). We obtain (26) by making the same substitution into (22). ■

## References

- A. Agarwal and J. C. Duchi. Distributed delayed stochastic optimization. In *Advances in Neural Information Processing Systems 24*, pages 873–881. 2011. URL <http://papers.nips.cc/paper/4247-distributed-delayed-stochastic-optimization.pdf>.
- H. Avron, A. Druinsky, and A. Gupta. Revisiting asynchronous linear solvers: Provable convergence rate through randomization. *IPDPS*, 2014.
- A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sciences*, 2(1):183–202, 2009.
- A. Beck and L. Tetruashvili. On the convergence of block coordinate descent type methods. *SIAM Journal on Optimization*, 23(4):2037–2060, 2013.
- D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Pentice Hall, 1989.
- S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2011.
- C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines, 2011. URL <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.
- C. Cortes and V. Vapnik. Support vector networks. *Machine Learning*, pages 273–297, 1995.
- A. Cotter, O. Shamir, N. Srebro, and K. Sridharan. Better mini-batch algorithms via accelerated gradient methods. In *Advances in Neural Information Processing Systems 24*, pages 1647–1655. 2011. URL <http://papers.nips.cc/paper/4432-better-mini-batch-algorithms-via-accelerated-gradient-methods.pdf>.
- J. C. Duchi, A. Agarwal, and M. J. Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic Control*, 57(3):592–606, 2012.
- M. C. Ferris and O. L. Mangasarian. Parallel variable distribution. *SIAM Journal on Optimization*, 4(4):815–832, 1994.

- D. Goldfarb and S. Ma. Fast multiple-splitting algorithms for convex optimization. *SIAM Journal on Optimization*, 22(2):533–556, 2012.
- Z. Lu and L. Xiao. On the complexity analysis of randomized block-coordinate descent methods. Technical Report arXiv:1305.4723, Simon Fraser University, 2013.
- Z. Q. Luo and P. Tseng. On the convergence of the coordinate descent method for convex differentiable minimization. *Journal of Optimization Theory and Applications*, 72:7–35, 1992.
- O. L. Mangasarian. Parallel gradient distribution in unconstrained optimization. *SIAM Journal on Optimization*, 33(1):916–1925, 1995.
- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19:1574–1609, 2009.
- Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Kluwer Academic Publishers, 2004.
- Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- F. Niu, B. Recht, C. Ré, and S. J. Wright. HOGWILD!: A lock-free approach to parallelizing stochastic gradient descent. *Advances in Neural Information Processing Systems 24*, pages 693–701, 2011.
- Z. Peng, M. Yan, and W. Yin. Parallel and distributed sparse optimization. Preprint, 2013.
- P. Richtárik and M. Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144:1–38, 2012a.
- P. Richtárik and M. Takáč. Parallel coordinate descent methods for big data optimization. Technical Report arXiv:1212.0873, 2012b.
- C. Scherrer, A. Tewari, M. Halappanavar, and D. Haglin. Feature clustering for accelerating parallel coordinate descent. *Advances in Neural Information Processing Systems 25*, pages 28–36, 2012.
- S. Shalev-Shwartz and T. Zhang. Accelerated mini-batch stochastic dual coordinate ascent. *Advances in Neural Information Processing Systems 26*, pages 378–385, 2013.
- O. Shamir and T. Zhang. Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. In *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1996.

- P. Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications*, 109:475–494, 2001.
- P. Tseng and S. Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming, Series B*, 117:387–423, June 2009.
- P. Tseng and S. Yun. A coordinate gradient descent method for linearly constrained smooth optimization and support vector machines training. *Computational Optimization and Applications*, 47(2):179–206, 2010.
- P.-W. Wang and C.-J. Lin. Iteration complexity of feasible descent methods for convex optimization. *Journal of Machine Learning Research*, 15:1523–1548, 2014.
- S. J. Wright. Accelerated block-coordinate relaxation for regularized optimization. *SIAM Journal on Optimization*, 22(1):159–186, 2012.
- T. Yang. Trading computation for communication: Distributed stochastic dual coordinate ascent. *Advances in Neural Information Processing Systems 26*, pages 629–637, 2013.