
Accelerated Cyclic Coordinate Dual Averaging with Extrapolation for Composite Convex Optimization

Cheuk Yin Lin¹ Chaobing Song¹ Jelena Diakonikolas¹

Abstract

Exploiting partial first-order information in a cyclic way is arguably the most natural strategy to obtain scalable first-order methods. However, despite their wide use in practice, cyclic schemes are far less understood from a theoretical perspective than their randomized counterparts. Motivated by a recent success in analyzing an extrapolated cyclic scheme for generalized variational inequalities, we propose an *Accelerated Cyclic Coordinate Dual Averaging with Extrapolation* (A-CODER) method for composite convex optimization, where the objective function can be expressed as the sum of a smooth convex function accessible via a gradient oracle and a convex, possibly nonsmooth, function accessible via a proximal oracle. We show that A-CODER attains the optimal convergence rate with improved dependence on the number of blocks compared to prior work. Furthermore, for the setting where the smooth component of the objective function is expressible in a finite sum form, we introduce a variance-reduced variant of A-CODER, VR-A-CODER, with state-of-the-art complexity guarantees. Finally, we demonstrate the effectiveness of our algorithms through numerical experiments.

1. Introduction

Block coordinate descent methods are broadly used in machine learning due to their effectiveness on large datasets brought by cheap iterations requiring only partial access to problem information (Wright, 2015; Nesterov, 2012). They are frequently applied to problems such as feature selection (Wu et al., 2008; Friedman et al., 2010; Mazumder et al., 2011), empirical risk minimization (Nesterov, 2012; Zhang

& Lin, 2015; Lin et al., 2015; Allen-Zhu et al., 2016; Alacaoglu et al., 2017; Gürbüzbalaban et al., 2017; Diakonikolas & Orecchia, 2018), and in distributed computing (Liu et al., 2014; Fercoq & Richtárik, 2015; Richtárik & Takáč, 2016). In the more recent literature, coordinate updates on either the primal or the dual side in primal-dual settings have been used to attain variance-reduced guarantees in finite sum settings (Chambolle et al., 2018; Alacaoglu et al., 2017; 2020; Song et al., 2020; 2021b).

Most of the existing theoretical results for (block) coordinate-type methods have been established for algorithms that select coordinate blocks to be updated by random sampling without replacement (Nesterov, 2012; Wright, 2015; Chambolle et al., 2018; Alacaoglu et al., 2017; 2020; Song et al., 2020; 2021b; Zhang & Lin, 2015; Lin et al., 2015; Allen-Zhu et al., 2016; Diakonikolas & Orecchia, 2018). Such methods are commonly referred to as the randomized block coordinate methods (RBCMs). What makes these methods particularly appealing from the aspect of convergence analysis is that the gradient evaluated on the sampled coordinate block can be related to the full gradient, by taking the expectation over the random choice of a coordinate block.

An alternative class of block coordinate methods is the class of cyclic block coordinate methods (CBCMs), which update blocks of coordinates in a cyclic order. CBCMs are frequently used in practice due to often superior empirical performance compared to RBCMs (Beck & Tetrushvili, 2013; Chow et al., 2017; Sun & Ye, 2019) and are also part of standard software packages for high-dimensional computational statistics such as GLMNet (Friedman et al., 2010) and SparseNet (Mazumder et al., 2011). However, CBCMs have traditionally been considered much more challenging to analyze than RBCMs.

The first convergence rate analysis of CBCMs for smooth convex optimization problems, obtained by Beck & Tetrushvili (2013), relied on relating the partial coordinate blocks of the gradient to the full gradient. For this reason, the dependence of iteration complexity on the number of coordinate blocks in Beck & Tetrushvili (2013) scaled linearly and as a square root for vanilla CBCM and its accelerated variant, respectively. Such a high dependence on

¹Department of Computer Sciences, University of Wisconsin-Madison. Correspondence to: Cheuk Yin Lin <cylin@cs.wisc.edu>.

the number of blocks (equal to the dimension in the coordinate case) makes the complexity guarantee of CBCMs seem worse than not only RBCMs but even full gradient methods such as gradient descent and the fast gradient method of Nesterov (1983), bringing into question their usefulness. This is further exacerbated by a result that shows that such a high gap in complexity does happen in the worst case (Sun & Ye, 2019), prompting research that would explain the gap between the theory and practice of CBCMs. However, most of the results that improved the dependence on the number of blocks only did so for structured classes of convex quadratic problems (Wright & Lee, 2020; Lee & Wright, 2019; Gürbüzbalaban et al., 2017).

On the other hand, a very recent work in Song & Diakonikolas (2021) introduced an extrapolated CBCM for variational inequalities whose complexity guarantee does not involve explicit dependence on the number of blocks. This result is enabled by a novel Lipschitz condition introduced in the same work. While the result from Song & Diakonikolas (2021) applies to convex minimization settings as a special case, the obtained convergence rates are not accelerated. Our main motivation in this work is to close this convergence gap by providing accelerated extrapolated CBCMs for convex composite minimization.

1.1. Contributions

We study the following composite convex problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x}) \right\}, \quad (\text{P})$$

where f is smooth and convex and g is proper, (possibly strongly) convex, and lower semicontinuous. This is a standard and broadly studied setting of structured nonsmooth optimization; see, e.g. Beck & Teboulle (2009); Nesterov (2007) and the follow-up work. To further make the problem amenable to optimization via block coordinate methods, we assume that g is block separable, with each component function admitting an efficiently computable prox operator (see Section 2 for a precise statement of the assumptions).

Similar to Song & Diakonikolas (2021), we define a summary Lipschitz constant L of f obtained from Lipschitz conditions of individual blocks. Our summary Lipschitz condition is similar to that of Song & Diakonikolas (2021) (although not exactly the same) and enjoys the same favorable properties as the condition introduced in that paper; see Section 2 for more details.

We introduce a new accelerated cyclic algorithm for (P) whose full gradient oracle complexity (number of full gradient passes or, equivalently, number of full cycles) is of the order $O\left(\min\left\{\sqrt{L}k\mathbf{x}_0 \quad \mathbf{x} \quad k_2, \sqrt{L} \log\left(\frac{Lk\mathbf{x}_0 \quad \mathbf{x} \quad k_2}{\gamma}\right)\right\}\right)$, where γ is the strong convexity parameter of g (equal to zero if g is only convex, by convention), \mathbf{x}^* is an optimal

solution to (P), and $\mathbf{x}_0 \in \text{dom}(g)$ is an arbitrary initial point. This complexity result matches the gradient oracle complexity of the fast gradient method (Nesterov, 1983), but with the traditional Lipschitz constant being replaced by the Lipschitz constant introduced in our work. In the very worst case, this constant is no higher than $\frac{1}{m}$ times the traditional one, where m is the number of blocks, giving an $m^{1/4}$ improvement in the resulting complexity over the accelerated cyclic method from Beck & Tretuashvili (2013). Even in this worst case, the obtained improvement in the dependence on the number of blocks is the first such improvement for accelerated methods since the work of Beck & Tretuashvili (2013). We note, however, that for both synthetic data and real data sets and on an example problem where both Lipschitz constants are explicitly computable, our Lipschitz constant is within a small constant factor (smaller than 1.5) of the traditional one (see Figure 1, Table 1, and the related discussion in Section 2).

Some key ingredients in our analysis are the following. First, we construct an estimate of the optimality gap we want to bound, where we replace the gradient terms with a vector composed of partial, or block, extrapolated gradient terms evaluated at intermediate points within a cycle. Crucially, we show that the error introduced by doing so can be controlled and bounded via our Lipschitz condition. An auxiliary result allowing us to carry out the analysis and appropriately bound the error terms resulting from our approach is Lemma 1, which shows that our Lipschitz condition translates into inequalities of the form

$$\begin{aligned} f(\mathbf{y}) &\leq f(\mathbf{x}) + \langle \mathbf{h}_r f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2, \\ k_r f(\mathbf{y}) &\leq r f(\mathbf{x}) + k^2 - 2L \langle f(\mathbf{y}) - f(\mathbf{x}), \mathbf{h}_r f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle, \end{aligned}$$

similar to the standard inequalities that hold for the traditional, full-gradient, Lipschitz constant. Finally, we note that the accelerated algorithm that we introduce is novel even in the single block (i.e., full-gradient) setting, due to the employed gradient extrapolation.

We further consider the finite sum setting, where f is expressible as $f(\mathbf{x}) = \frac{1}{n} \sum_{t=1}^n f_t(\mathbf{x})$, and where n is typically very large. We then propose a variance-reduced variant of our accelerated method, which further reduces the full gradient oracle complexity to $O\left(\min\left\{\sqrt{\frac{L}{n}}k\mathbf{x}_0 \quad \mathbf{x} \quad k_2, \sqrt{\frac{L}{n}} \log\left(\frac{Lk\mathbf{x}_0 \quad \mathbf{x} \quad k_2}{\gamma}\right)\right\}\right)$. The variance reduction that we employ is of the SVRG type (Johnson & Zhang, 2013). While following a similar approach as the basic accelerated algorithm described above, the analysis in this case turns out to be much more technical, due to the need to simultaneously handle error terms arising from variance reduction as well as the error terms arising from the cyclic updates. Through utilizing the novel smoothness properties obtained in Lemma 1 specific to convex minimization,

we are able to obtain the desired error bounds without using the additional point extrapolation step in the gradient estimator as Song & Diakonikolas (2021), but rather only with an SVRG estimator. This important change paves a path to achieving accelerated convergence rates while also simplifying the implementation of our algorithms.

Last but not least, we demonstrate the practical efficacy of our novel accelerated algorithms A-CODER and VR-A-CODER through numerical experiments, comparing against other relevant block coordinate descent methods. The use of A-CODER and VR-A-CODER achieves faster convergence in primal gap with respect to both the number of full-gradient evaluations and wall-clock time.

1.2. Further Discussion of Related Work

As discussed at the beginning of this section, cyclic block coordinate methods constitute a fundamental class of optimization methods whose convergence is not yet well understood. In the worst case, the full gradient oracle complexity of vanilla cyclic block coordinate gradient update is worse than that of vanilla gradient descent, by a factor scaling with the number of blocks m (equal to the dimension in the coordinate case) (Sun & Ye, 2019; Beck & Tetrushvili, 2013). Since the initial results providing such an upper bound (Sun & Ye, 2019), there were no improvements on the dependence on the number of blocks in the convergence guarantees of cyclic methods until the very recent work of Song & Diakonikolas (2021), which in the worst case improves the dependence on m by a factor $\rho_{\overline{m}}$. Our work further contributes to this line of work by improving the dependence on m in accelerated methods from $\rho_{\overline{m}}$ to m^{1-d} in the worst case.

In the finite-sum settings, variance reduction has been widely explored; e.g., in Johnson & Zhang (2013); Defazio et al. (2014); Allen-Zhu (2017); Reddi et al. (2016); Lei et al. (2017); Song et al. (2020); Schmidt et al. (2017) for the case of full-gradient methods and in Chen & Gu (2016); Lei & Shanbhag (2018) for randomized block coordinate methods. However, variance reduced schemes for cyclic methods are much more rare, with nonasymptotic guarantees being obtained very recently for the case of variational inequalities (Song & Diakonikolas, 2021) and nonconvex optimization (Cai et al., 2022; Xu & Yin, 2014). We are not aware of any existing variance reduced results for accelerated cyclic block coordinate methods.

1.3. Outline of the Paper

Section 2 introduces the necessary notation and background and outlines our main problem assumptions. Section 3 introduces the A-CODER algorithm and outlines the analysis. For space constraints, the full convergence analysis of A-CODER is provided in Appendix A. Section 4 presents

VR-A-CODER and outlines its convergence analysis, while the full technical details are deferred to Appendix B. Finally, Section 5 provides numerical experiments for our results and concludes the paper with a discussion.

2. Notation and Preliminaries

For a positive integer K , we use $[K]$ to denote the set $\{1, 2, \dots, K\}$. We consider the d -dimensional Euclidean space $(\mathbb{R}^d, \|\cdot\|)$, where $\|\cdot\|$ denotes the Euclidean norm, $\langle \cdot, \cdot \rangle$ denotes the (standard) inner product, and d is assumed to be finite. Throughout the paper, we assume that there is a given partition of the set $[d]$ into sets S^j , $j \in [m]$, where $|S^j| = d^j > 0$. For convenience of notation, we assume that sets S^j are comprised of consecutive elements from $[d]$, that is, $S^1 = \{1, 2, \dots, d^1\}$, $S^2 = \{d^1 + 1, d^1 + 2, \dots, d^1 + d^2\}$, \dots , $S^m = \{d^1 + d^2 + \dots + d^{m-1} + 1, d^1 + d^2 + \dots + d^{m-1} + 2, \dots, d^1 + d^2 + \dots + d^m\}$. This assumption is without loss of generality, as all our results are invariant to permutations of the coordinates (though the value of the Lipschitz constant of the gradients defined in our work depends on the ordering of the coordinates; see Assumption 2). For a vector $\mathbf{x} \in \mathbb{R}^d$, we use $\mathbf{x}^{(j)}$ to denote its coordinate components indexed by S^j . Similarly for a gradient $\mathbf{r}f$ of a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we use $\mathbf{r}^{(j)}f$ to denote its coordinate components indexed by S^j . We use $(\cdot)_j$ to denote an operator for vectors and square matrices that replaces the first $j-1$ elements of rows and columns with zeros, i.e., keeping elements with indices j the same, otherwise zeros.

Given a proper, convex, lower semicontinuous function $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$, we use $\partial g(\mathbf{x})$ to denote the subdifferential set (the set of all subgradients) of g . Of particular interests to us are functions g whose proximal operator (or resolvent), defined by

$$\text{prox}_g(\mathbf{u}) := \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \tau g(\mathbf{x}) + \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|^2 \right\} \quad (1)$$

is efficiently computable for all $\tau > 0$ and $\mathbf{u} \in \mathbb{R}^d$. To unify the cases in which g are convex and strongly convex respectively, we say that g is γ -strongly convex with modulus $\gamma \geq 0$, if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\mathbf{g}^0(\mathbf{x}) \in \partial g(\mathbf{x})$,

$$g(\mathbf{y}) \leq g(\mathbf{x}) + \langle \mathbf{g}^0(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\gamma}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Problem definition. We consider Problem (P), under the following assumptions.

Assumption 1. $g(\mathbf{x})$ is γ -strongly convex, where $\gamma \geq 0$, and block-separable over coordinate sets $f^{S^j} g_{j=1}^m : g(\mathbf{x}) = \sum_{j=1}^m g^j(\mathbf{x}^{(j)})$. Each $g^j(\mathbf{x}^{(j)})$ for $j \in [m]$ admits an efficiently computable proximal operator.

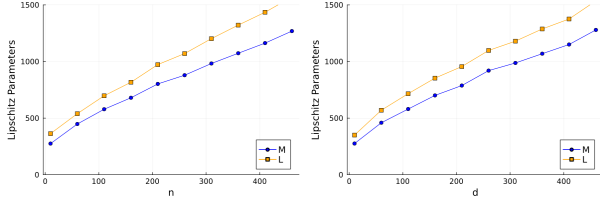


Figure 1. Comparisons of Lipschitz constants for elastic-net problems on synthetic datasets, where M denotes the commonly known Lipschitz constant and L is our new Lipschitz constant as defined in Assumption 2.

Table 1. Comparisons of Lipschitz constants for elastic-net problems on LibSVM datasets. M is the classical gradient Lipschitz constant and L is our novel smoothness constant. We use each coordinate as a block, i.e., $m = d$.

DATASET	#FEATURES	M	L
SONAR	60	12.5	15.8
COLON	2000	310.6	394.7
A9A	123	6.1	7.7
PHISHING	68	0.60	0.76
MADELON	500	1.2	1.5

Assumption 2. There exist positive semidefinite matrices $\tilde{f}Q^1, Q^2, \dots, Q^m g$ such that $r^{(j)}f(\cdot)$ is 1-Lipschitz continuous w.r.t. the seminorm $k k_{Q^j}$, i.e., $\delta \mathbf{x}, \mathbf{y} \geq \mathbb{R}^d$,

$$k r^{(j)}f(\mathbf{x}) - r^{(j)}f(\mathbf{y}) k^2 \leq k \mathbf{x} - \mathbf{y} k_{Q^j}^2, \quad (2)$$

where $k \mathbf{x} - \mathbf{y} k_{Q^j}^2 := (\mathbf{x} - \mathbf{y})^T Q^j (\mathbf{x} - \mathbf{y})$ is the Mahalanobis (semi)norm. Moreover, we define a new Lipschitz constant L such that $L^2 = 2k Q k < 1$ where $Q = \sum_{j=1}^m [(Q^j)_{j,j} + (Q^j)_{j,j+1}]$.

Observe that when f is M -smooth in a traditional sense (i.e., when f has M -Lipschitz gradients w.r.t. the Euclidean norm), Assumption 2 can be trivially satisfied using $Q^j = M I$ for all $j \in [m]$, where I is the identity matrix. Consequently, it can be argued that $L = 2\sqrt{m}M$ (Song & Diakonikolas, 2021); however, we show that this bound is much tighter in practice as illustrated in Figure 1 and in Table 1. In particular, we follow the experiments in Song & Diakonikolas (2021) and show empirically that the standard Lipschitz constant M and our new Lipschitz constant L scale within the same factor for both synthetic and real data.

3. Accelerated Cyclic Algorithm

In this section, we introduce and analyze A-CODER, whose pseudocode is provided in Algorithm 1. A-CODER can be seen as a Nesterov-style accelerated variant of CODER, previously introduced for solving variational inequalities by Song & Diakonikolas (2021). A-CODER is related to

other accelerated algorithms in the following sense. In the case of a single block ($m = 1$) and when gradient extrapolation is not used (i.e., when $\mathbf{q}_k = \mathbf{p}_k$), A-CODER reduces to a generalized variant of AGD+ (Cohen et al., 2018; Diakonikolas & Guzmán, 2021) or the method of similar triangles (Gasnikov & Nesterov, 2018). The analysis of A-CODER follows the general gap bounding argument (Diakonikolas & Orecchia, 2019; Song et al., 2021a) and it is based on three key ingredients: (i) gradient extrapolation, which enables the use of partial information about the gradients within a full epoch of cyclic updates, (ii) Lipschitz condition for the gradients based on the Mahalanobis norm as defined in Assumption 2, and (iii) upper and lower bounds on the difference between the function and its linear approximation that are compatible with the gradient Lipschitz condition that we use, as stated in Lemma 1.

Lemma 1. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex and smooth function whose gradients satisfy Assumption 2. Then, $\delta \mathbf{x}, \mathbf{y} \geq \mathbb{R}^d$:

$$f(\mathbf{y}) - f(\mathbf{x}) \leq \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} k \mathbf{y} - \mathbf{x} k^2,$$

$$k r f(\mathbf{y}) - r f(\mathbf{x}) k^2 \leq 2L(f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle).$$

We now derive the A-CODER algorithm. We define $\tilde{f}a_k g_{k-1}$ and $\tilde{f}A_k g_{k-1}$ to be sequences of positive numbers with $A_k = \sum_{i=1}^k a_i, a_0 = A_0 = 0$. Let $\tilde{f}x_k g_{k-1}$ be an arbitrary sequence of points in $\text{dom}(g)$. Our goal here is to bound the function value gap $f(\mathbf{y}_k) - f(\mathbf{u})$ above for all $\mathbf{u} \in \text{dom}(g)$. Towards this goal, we define an estimation sequence ψ_k recursively by $\psi_0(\mathbf{u}) = \frac{1}{2} k \mathbf{u} - \mathbf{x}_0 k^2$ and

$$\psi_k(\mathbf{u}) := \psi_{k-1}(\mathbf{u}) + a_k(f(\mathbf{x}_k) + \langle \mathbf{h} \mathbf{q}_k, \mathbf{u} - \mathbf{x}_k \rangle + g(\mathbf{u}))$$

for $k \geq 1$. Meanwhile, \mathbf{v}_k and \mathbf{y}_k are defined as $\mathbf{v}_k := \arg \min_{\mathbf{u} \in \mathbb{R}^d} \psi_k(\mathbf{u})$ and $\mathbf{y}_k := \frac{1}{A_k} \sum_{i=1}^k a_i \mathbf{v}_i$ respectively. We start our analysis by characterizing the gap function in the following lemma.

Lemma 2. For any $\mathbf{u} \in \mathbb{R}^d$ and any sequence of vectors $\tilde{f}q_i g_{i-1}$, we have

$$A_k(f(\mathbf{y}_k) - f(\mathbf{u})) \quad (3)$$

$$\sum_{i=1}^k E_i(\mathbf{u}) + \frac{1}{2} k \mathbf{u} - \mathbf{x}_0 k^2 \leq \frac{1 + A_k \gamma}{2} k \mathbf{u} - \mathbf{v}_k k^2, \quad (4)$$

where

$$E_i(\mathbf{u}) = A_i(f(\mathbf{y}_i) - f(\mathbf{x}_i)) - A_{i-1}(f(\mathbf{y}_{i-1}) - f(\mathbf{x}_i))$$

$$+ a_i \langle \mathbf{h} \mathbf{q}_i, \mathbf{v}_i - \mathbf{x}_i \rangle + a_i \langle \nabla f(\mathbf{x}_i), \mathbf{q}_i - \mathbf{x}_i \rangle - \langle \mathbf{u}, \mathbf{v}_i - \mathbf{x}_i \rangle$$

$$+ \frac{1 + A_{i-1} \gamma}{2} k \mathbf{v}_i - \mathbf{v}_{i-1} k^2. \quad (5)$$

Lemma 2 applies to an arbitrary algorithm that satisfies its assumptions. From now on, we make the analysis specific

Algorithm 1 Accelerated Cyclic cOordinate Dual avEraging with extRapolation (A-CODER)

```

1: Input:  $\mathbf{x}_0 \in \text{dom}(g)$ ,  $\gamma \in [0, L]$ ,  $L > 0$ ,  $m$ ,
    $f \in S^1, \dots, S^m$ 
2: Initialization:  $\mathbf{x}_1 = \mathbf{x}_0 = \mathbf{v}_1 = \mathbf{v}_0 = \mathbf{y}_0$ ;  $\mathbf{p}_0 =$ 
    $r f(\mathbf{x}_0)$ ;  $\mathbf{z}_0 = \mathbf{0}$ ;  $a_0 = A_0 = 0$ 
3: for  $k = 1$  to  $K$  do
4:   Set  $a_k > 0$  be largest value s.t.  $\frac{a_k^2}{A_k} \leq \frac{2(1+A_{k-1})}{5L}$ 
   where  $A_k = A_{k-1} + a_k$ 
5:    $\mathbf{x}_k = \frac{A_{k-1}}{A_k} \mathbf{y}_{k-1} + \frac{a_k}{A_k} \mathbf{v}_k$ 
6:   for  $j = m$  to  $1$  do
7:      $\mathbf{p}_k^{(j)} = r^{(j)} f(\mathbf{x}_k^{(1)}, \dots, \mathbf{x}_k^{(j)}, \mathbf{y}_k^{(j+1)}, \dots, \mathbf{y}_k^{(m)})$ 
8:      $\mathbf{q}_k^{(j)} = \mathbf{p}_k^{(j)} + \frac{a_{k-1}}{a_k} (r^{(j)} f(\mathbf{x}_{k-1}) - \mathbf{p}_k^{(j-1)})$ 
9:      $\mathbf{z}_k^{(j)} = \mathbf{z}_{k-1}^{(j)} + a_k \mathbf{q}_k^{(j)}$ 
10:     $\mathbf{v}_k^{(j)} = \text{prox}_{A_k g_j}(\mathbf{x}_0, \mathbf{z}_k^{(j)})$ 
11:     $\mathbf{y}_k^{(j)} = \frac{A_{k-1}}{A_k} \mathbf{y}_{k-1}^{(j)} + \frac{a_k}{A_k} \mathbf{v}_k^{(j)}$ 
12:   end for
13: end for
14: return  $\mathbf{v}_K, \mathbf{y}_K$ 
    
```

to A-CODER (Algorithm 1). In Lemma 2, $fE_i(\mathbf{u})g$ are the error terms that we need to bound above. If $\sum_{i=1}^k E_i(\mathbf{u}) \leq \frac{1+A_k}{2} k\mathbf{u} - \mathbf{v}_k k^2$, then we get the desired $1/A_k$ rate. To this end, we bound each term $E_k(\mathbf{u})$ in Lemma 3 by using the extrapolation direction \mathbf{q}_k , the definition of $\mathbf{y}_k, \mathbf{x}_k$ and the parameter setting of a_k .

Lemma 3. *Let $\mathbf{x}_0 \in \text{dom}(g)$ be an arbitrary initial point and consider the updates in Algorithm 1. If, for $k \geq 1$, $\frac{a_k^2}{A_k} \leq \frac{2(1+A_{k-1})}{5L}$, then $\delta \mathbf{u}$,*

$$\begin{aligned}
 E_k(\mathbf{u}) &= a_k \langle r f(\mathbf{x}_k) - \mathbf{p}_k, \mathbf{v}_k - \mathbf{u} \rangle \\
 &\quad + a_{k-1} \langle r f(\mathbf{x}_{k-1}) - \mathbf{p}_{k-1}, \mathbf{v}_{k-1} - \mathbf{u} \rangle \\
 &\quad + \frac{1+A_{k-1}\gamma}{10} k \mathbf{v}_k - \mathbf{v}_{k-1} k^2 \\
 &\quad + \frac{1+A_{k-1}2\gamma}{10} k \mathbf{v}_{k-1} - \mathbf{v}_{k-2} k^2.
 \end{aligned}$$

We are now ready to state the main convergence result of this section.

Theorem 1. *Let $\mathbf{x}_0 \in \text{dom}(g)$ be an arbitrary initial point and consider the updates in Algorithm 1. Then, $\delta k \geq 1$ and any $\mathbf{u} \in \text{dom}(g)$:*

$$f(\mathbf{y}_k) - f(\mathbf{u}) + \frac{3(1+A_{k-1}\gamma)}{10A_k} k\mathbf{u} - \mathbf{v}_k k^2 \leq \frac{k\mathbf{u} - \mathbf{x}_0 k^2}{2A_k}.$$

In particular, if $\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$ exists, then

$$f(\mathbf{y}_k) - f(\mathbf{x}^*) \leq \frac{k\mathbf{x}^* - \mathbf{x}_0 k^2}{2A_k}.$$

Further, in this case we also have:

$$k\mathbf{v}_k - \mathbf{x} k^2 \leq \frac{5}{3(1+A_{k-1}\gamma)} k\mathbf{x}_0 - \mathbf{x} k^2,$$

$$k\mathbf{y}_k - \mathbf{x} k^2 \leq \left(\frac{5}{3A_k} \sum_{i=1}^k \frac{a_i}{1+A_{i-1}\gamma} \right) k\mathbf{x}_0 - \mathbf{x} k^2.$$

Finally, in all the bounds we have

$$A_k \leq \max \left\{ \frac{2}{5L} \left(1 + \sqrt{\frac{2\gamma}{5L}} \right)^k, \frac{k^2}{10L} \right\}.$$

Adaptive A-CODER. The Lipschitz parameter L used in the statement of A-CODER (Algorithm 1) is usually not readily available for typical instances of convex composite minimization problems; however, as we argue in Appendix A, this parameter can be adaptively estimated using the standard backtracking line search. A variant of A-CODER implementing this adaptive estimation of L is provided in Algorithm 3. This is enabled by our analysis, which only requires the stated Lipschitz condition to hold between the successive iterates of the algorithm. Notably, unlike randomized algorithms which estimate Lipschitz constants for each of the coordinate blocks (see, e.g., Nesterov (2012)), we only need to estimate one summary Lipschitz parameter L .

4. Variance Reduced A-CODER

In this section, we assume that the problem (P) has a finite sum structure, i.e., $f(\mathbf{x}) = \frac{1}{n} \sum_{t=1}^n f_t(\mathbf{x})$, where n may be very large. For this case, we can further reduce the per-iteration cost and improve the complexity results by combining the well-known SVRG-style variance reduction strategy (Johnson & Zhang, 2013) with the results from the previous section to obtain our variance reduced A-CODER (VR-A-CODER). From another perspective, VR-A-CODER can be seen as a cyclic gradient-extrapolated version of the recent VRADA algorithm for finite-sum composite convex minimization (Song et al., 2020).

For this finite-sum setting, we need to make the following stronger assumption for each $f_t(\mathbf{x})$.

Assumption 3. *For all $t \in [n]$, $f_t(\mathbf{x})$ is convex. Moreover for all $t \in [n]$, there exist positive semidefinite matrices $\mathbf{Q}^1, \mathbf{Q}^2, \dots, \mathbf{Q}^m \in \mathbb{R}^m$ such that $r^{(j)} f_t(\cdot)$ is 1-Lipschitz continuous w.r.t. the norm $\|\cdot\|_{\mathbf{Q}^j}$ i.e., $\delta \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, t \in [n]$,*

$$\|r^{(j)} f_t(\mathbf{x}) - r^{(j)} f_t(\mathbf{y})\|_{\mathbf{Q}^j} \leq \|\mathbf{x} - \mathbf{y}\|_{\mathbf{Q}^j}.$$

Lemma 4. *If $f(\mathbf{x}) = \frac{1}{n} \sum_{t=1}^n f_t(\mathbf{x})$ satisfies Assumption 3, then it satisfies Assumption 2 and thus Lemma 1 holds.*

Algorithm 2 Variance Reduced A-CODER (Implementable Version)

- 1: **Input:** $\mathbf{x}_0 \in \text{dom}(g)$, $\gamma \geq 0$, $L > 0$, m , $fS^1, \dots, S^m g$
- 2: **Initialization:** $\mathbf{y}_0 = \mathbf{v}_{1,0} = \mathbf{y}_{1,0} = \mathbf{x}_{1,1} = \mathbf{x}_0$; $\mathbf{z}_{1,0} = \mathbf{0}$
- 3: $a_0 = A_0 = 0$; $A_1 = a_1 = \frac{1}{4L}$
- 4: $\mathbf{z}_{1,1} = \gamma f(\mathbf{x}_0)$; $\mathbf{v}_{1,1} = \text{prox}_{a_1 g}(\mathbf{x}_0 - \mathbf{z}_{1,1})$
- 5: $\mathbf{y}_1 = \mathbf{y}_{1,1} = \mathbf{v}_{1,1}$
- 6: $\mathbf{w}_{1,1:j} = (\mathbf{x}_{1,1}^{(1)}, \dots, \mathbf{x}_{1,1}^{(j)}, \mathbf{y}_{1,1}^{(j+1)}, \dots, \mathbf{y}_{1,1}^{(m)})$
- 7: $\mathbf{v}_{2,0} = \mathbf{v}_{1,1}$; $\mathbf{w}_{2,0:j} = \mathbf{w}_{1,1:j}$; $\mathbf{x}_{2,0} = \mathbf{x}_{1,1}$; $\mathbf{y}_{2,0} = \mathbf{y}_{1,1}$; $\mathbf{z}_{2,0} = \mathbf{z}_{1,1}$
- 8: **for** $s = 2$ **to** S **do**
- 9: $a_s = \sqrt{\frac{KA_{s-1}(1+A_{s-1})}{8L}}$; $A_s = A_{s-1} + a_s$
- 10: $a_{s,0} = a_{s-1}$; $a_{s,1} = a_{s,2} = \dots = a_{s,K} = a_s$
- 11: $\mathbf{v}_{s,0} = \mathbf{v}_{s-1,K}$; $\mathbf{w}_{s,0:j} = \mathbf{w}_{s-1,K;j}$; $\mathbf{x}_{s,0} = \mathbf{x}_{s-1,K}$; $\mathbf{y}_{s,0} = \mathbf{y}_{s-1,K}$; $\mathbf{z}_{s,0} = \mathbf{z}_{s-1,K}$
- 12: $\mathbf{y}_s = \gamma f(\mathbf{y}_{s-1})$
- 13: **for** $k = 1$ **to** K **do**
- 14: $\mathbf{x}_{s;k} = \frac{A_{s-1}}{A_s} \mathbf{y}_{s-1} + \frac{a_s}{A_s} \mathbf{v}_{s;k-1}$
- 15: **for** $j = m$ **to** 1 **do**
- 16: $\mathbf{w}_{s;k;j} = (\mathbf{x}_{s;k}^{(1)}, \dots, \mathbf{x}_{s;k}^{(j)}, \mathbf{y}_{s;k}^{(j+1)}, \dots, \mathbf{y}_{s;k}^{(m)})$
- 17: Choose t in $[m]$ uniformly at random
- 18: $\tilde{r}_{s;k}^{(j)} = \gamma f_t(\mathbf{w}_{s;k;j}) - \gamma f_t(\mathbf{y}_{s-1}) + \frac{\gamma}{s}$
- 19: $\mathbf{q}_{s;k}^{(j)} = \tilde{r}_{s;k}^{(j)} + \frac{a_{s;k-1}}{a_s} (\gamma f_t(\mathbf{x}_{s;k-1}) - \gamma f_t(\mathbf{w}_{s;k-1;j}))$
- 20: $\mathbf{z}_{s;k}^{(j)} = \mathbf{z}_{s;k-1}^{(j)} + a_s \mathbf{q}_{s;k}^{(j)}$
- 21: $\mathbf{v}_{s;k}^{(j)} = \text{prox}_{(A_{s-1} + \frac{a_{s;k}}{K})g}(\mathbf{x}_0 - \mathbf{z}_{s;k}^{(j)}/K)$
- 22: $\mathbf{y}_{s;k}^{(j)} = \frac{A_{s-1}}{A_s} \mathbf{y}_{s-1} + \frac{a_s}{A_s} \mathbf{v}_{s;k}^{(j)}$
- 23: **end for**
- 24: **end for**
- 25: $\mathbf{y}_s = \frac{1}{K} \sum_{k=1}^K \mathbf{y}_{s;k}$
- 26: **end for**
- 27: **return** $\mathbf{v}_{S;K}, \mathbf{y}_S$

With this assumption, we can now derive the VR-A-CODER algorithm. Similar to Section 3, we define $\{a_s g_{s-1}\}$ and $\{A_s g_s\}$ to be sequences of positive numbers with $A_s = \sum_{i=1}^s a_i$, $a_0 = A_0 = 0$. Let $\{\mathbf{y}_s g_{s-1}\}$ be a sequence of points in $\text{dom}(g)$ which will be determined by the VR-A-CODER algorithm. Our goal here is to bound the function value gap $f(\mathbf{y}_s) - f(\mathbf{u})$ above $(\mathbf{u} \in \text{dom}(g))$. To attain this, we define the estimate sequence $\{\psi_{s;k} g_{s-1;k \in [K]}\}$ recursively by $\psi_{1,0}(\mathbf{u}) = \frac{K}{2} k \mathbf{u} - \mathbf{x}_0 k^2$,

$$\begin{aligned} \psi_{1,1}(\mathbf{u}) &= \psi_{1,0}(\mathbf{u}) + K a_1 (f(\mathbf{x}_0) \\ &\quad + h \gamma f(\mathbf{x}_0), \mathbf{u} - \mathbf{x}_0) + g(\mathbf{u}), \end{aligned} \quad (6)$$

and $\psi_{2,0} = \psi_{1,1}$; for $s = 2, 1 \leq k \leq K$,

$$\begin{aligned} \psi_{s;k}(\mathbf{u}) &= \psi_{s;k-1}(\mathbf{u}) + a_s (f(\mathbf{x}_{s;k}) \\ &\quad + h \mathbf{q}_{s;k}, \mathbf{u} - \mathbf{x}_{s;k}) + g(\mathbf{u}), \end{aligned} \quad (7)$$

and $\psi_{s+1,0} = \psi_{s;K}$. In Eqs. (6) and (7), \mathbf{x}_0 is the initial point, $\mathbf{x}_{s;k}$ and $\mathbf{y}_{s;k}$ are computed as convex combinations of two points, which is commonly used in Nesterov-style acceleration, and $\mathbf{q}_{s;k} = (\mathbf{q}_{s;k}^{(1)}, \mathbf{q}_{s;k}^{(2)}, \dots, \mathbf{q}_{s;k}^{(m)})$ is a variance reduced stochastic gradient with extrapolation, which is the main novelty in our algorithm design. Meanwhile, we define $\mathbf{v}_{s;k}$ by $\mathbf{v}_{s;k} := \arg \min_{\mathbf{u} \in \mathbb{R}^d} \psi_{s;k}(\mathbf{u})$ and note that due to the specific choice of $\mathbf{q}_{s;k}$, $\mathbf{v}_{s;k}$ is updated in a cyclic (block) coordinate way. Furthermore, in VR-A-CODER, for $s = 2$, we define $\mathbf{y}_s = \frac{1}{K} \sum_{k=1}^K \mathbf{y}_{s;k}$.

We start our analysis by characterizing the gap function, in the following lemma, similar to Lemma 2 in Section 3, although the proof is much more technical in this case. Due to space constraints, we include the full analysis of Variance Reduced A-CODER in Appendix B.

Lemma 5. For any $\mathbf{u} \in \mathbb{R}^d$ and any sequence of vectors $\{\mathbf{q}_{s;k} g_{s-2;k \in [K]}\}$, for all $S \geq 2$, we have

$$\begin{aligned} &KA_S (f(\mathbf{y}_S) - f(\mathbf{u})) \\ &\leq \frac{K}{2} k \mathbf{x}_0 - \mathbf{u} k^2 - \frac{K(1 + A_S \gamma)}{2} k \mathbf{v}_{S;K} - \mathbf{u} k^2 \\ &\quad - \frac{K}{4} k \mathbf{v}_{1,1} - \mathbf{v}_{1,0} k^2 + \sum_{s=2}^S \sum_{k=1}^K E_{s;k}(\mathbf{u}), \end{aligned}$$

where

$$\begin{aligned} E_{s;k}(\mathbf{u}) &= A_s (f(\mathbf{y}_{s;k}) - f(\mathbf{x}_{s;k})) \\ &\quad - A_{s-1} (f(\mathbf{y}_{s-1}) - f(\mathbf{x}_{s;k})) \\ &\quad + a_s h \gamma f(\mathbf{x}_{s;k}) - \mathbf{q}_{s;k}, \mathbf{x}_{s;k} - \mathbf{u} \\ &\quad + a_s h \mathbf{q}_{s;k}, \mathbf{x}_{s;k} - \mathbf{v}_{s;k} \\ &\quad - \frac{K(1 + A_{s-1} \gamma)}{2} k \mathbf{v}_{s;k} - \mathbf{v}_{s;k-1} k^2. \end{aligned} \quad (8)$$

In the following lemma, we bound the expected error terms $\sum_{s=2}^S \sum_{k=1}^K \mathbb{E}[E_{s;k}(\mathbf{u})]$ arising from the gap bound stated in the previous lemma. This bound is then finally used in Theorem 2 to obtain the claimed convergence results.

Lemma 6. With $a_s^2 = \frac{KA_{s-1}(1+A_{s-1})}{8L}$, $a_{s;k} = a_s$ and $A_{s;k} = A_s$ for $k \in [K]$, $a_{s,0} = a_{s-1}$ and $A_{s,0} = A_{s-1}$,

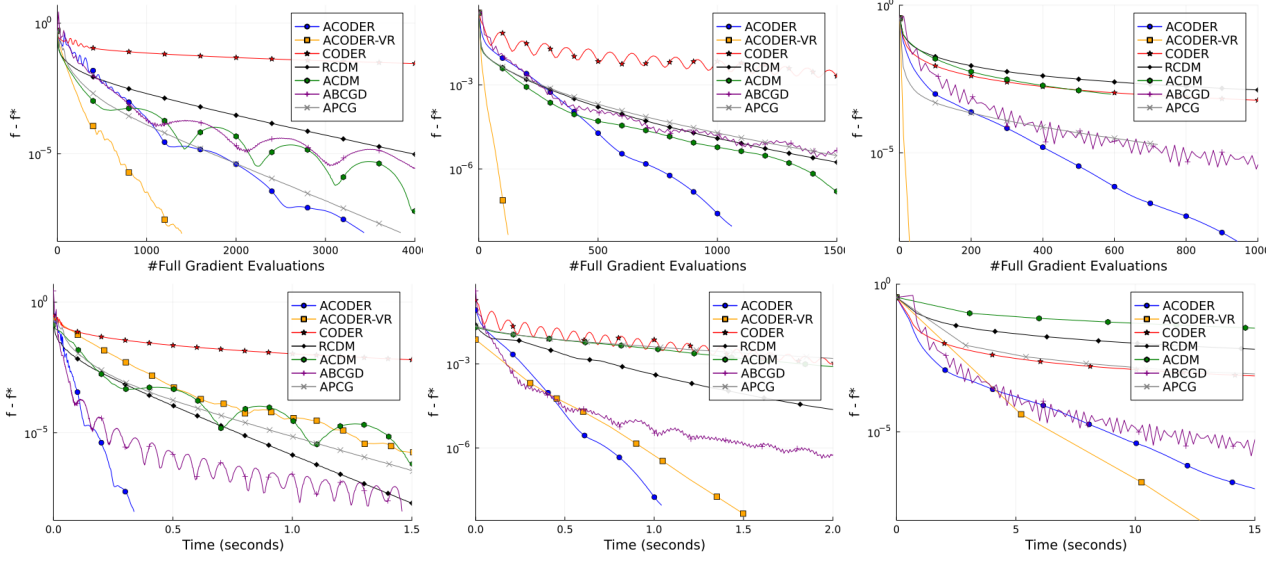


Figure 2. Performance comparisons between implemented algorithms in terms of the number of full-gradient evaluations and wall-clock time for logistic regression with ridge regularized problems. The top row contains plots against the number of full-gradient evaluations, and the bottom row contains plots against the wall-clock time. The left column is for the sonar dataset, the middle column is for the a1a dataset and the rightmost column is for the a9a dataset, all obtained from LIBSVM (Chang & Lin, 2011).

then for any fixed $\mathbf{u} \in \text{dom}(g)$ we have

$$\begin{aligned} & \sum_{s=2}^S \sum_{k=1}^K \mathbb{E} [E_{S;k}(\mathbf{u})] \\ & \sum_{j=1}^m a_1 \left\langle r^{(j)} f(\mathbf{x}_{1,1}) - r^{(j)} f(\mathbf{w}_{1,1;j}), \mathbf{v}_{1,1}^{(j)} - \mathbf{u}^{(j)} \right\rangle \\ & + \sum_{j=1}^m a_S \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{x}_{S;K}) - r^{(j)} f(\mathbf{w}_{S;K;j}), \mathbf{v}_{S;K}^{(j)} - \mathbf{u}^{(j)} \right\rangle + \frac{K}{64} \|\mathbf{v}_{1,1} - \mathbf{v}_{1,0}\|^2 \right] \\ & \frac{5K(1 + A_S \gamma)}{32} \mathbb{E} \left[\|\mathbf{v}_{S;K} - \mathbf{v}_{S;K}^*\|^2 \right], \end{aligned}$$

where $\mathbf{x}_{1,1}, \mathbf{v}_{1,0} \in \text{dom}(g)$ can be chosen arbitrarily and $\mathbf{w}_{1,1;j}$ is defined in Algorithm 4.

Our main result for this section is summarized in the following theorem.

Theorem 2. Let $\mathbf{x}_0 \in \text{dom}(g)$ be an arbitrary initial point. Fix $K \geq 1$ and consider the updates in Algorithm 4. Then for $S \geq 2$ and $\mathbf{u} \in \text{dom}(g)$, we have

$$\begin{aligned} \mathbb{E} [f(\mathbf{y}_S) - f(\mathbf{u})] & + \frac{9(1 + A_S \gamma)}{64A_S} \mathbb{E} \left[\|\mathbf{v}_{S;K} - \mathbf{u}\|^2 \right] \\ & \leq \frac{5}{8A_S} \|\mathbf{x}_0 - \mathbf{u}\|^2. \end{aligned}$$

In particular if $\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$ exists, then we have

$$\mathbb{E} [f(\mathbf{y}_S) - f(\mathbf{x}^*)] \leq \frac{5}{8A_S} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

and

$$\mathbb{E} \left[\|\mathbf{v}_{S;K} - \mathbf{x}^*\|^2 \right] \leq \frac{40}{9(1 + A_S \gamma)} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Finally in all the bounds above we have

$$A_S \leq \max \left\{ \frac{S^2 K}{64L}, \frac{1}{4L} \left(1 + \sqrt{\frac{K\gamma}{8L}} \right)^S \right\}.$$

Note that in Theorem 2, we can set the number of inner iterations K to be any positive integer. However, in order to balance the computational cost between the outer loop of each epoch and the inner loops, it is optimal to set $K = \sqrt{n}$ and for simplicity we can set $K = n$. Therefore, the total number of arithmetic operations required to obtain an ϵ -accurate solution \mathbf{y}_S by applying Algorithm 2 such that $\mathbb{E}[f(\mathbf{y}_S) - f(\mathbf{x}^*)] \leq \epsilon$ is at most $O\left(nd\sqrt{\frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|}{n}}\right)$ for the general convex case when $\gamma = 0$, and $O\left(\frac{nd \log\left(\frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|}{\epsilon}\right)}{\log\left(1 + \frac{\gamma}{n}\right)}\right)$ for the strongly convex case when $\gamma > 0$.

Adaptive VR-A-CODER. Similar to A-CODER, VR-A-CODER can adaptively estimate the Lipschitz parameter. For completeness, we have included the adaptive version of VR-A-CODER in Algorithm 5 (Appendix B).

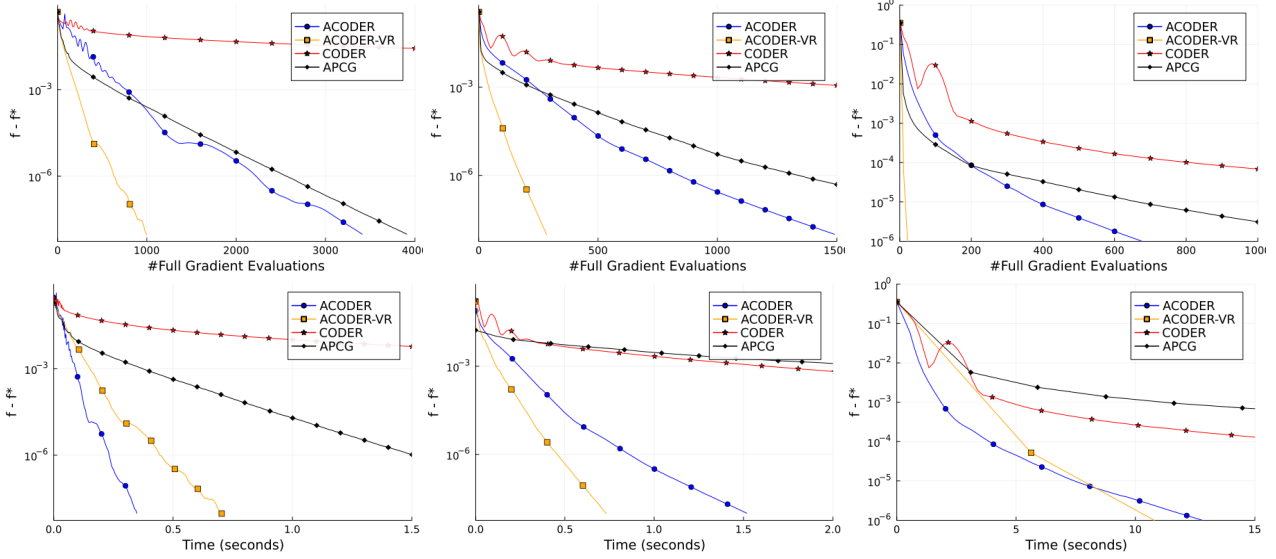


Figure 3. Performance comparisons between implemented algorithms in terms of the number of full-gradient evaluations and wall-clock time for logistic regression with elastic net regularized problems. The top row contains plots against the number of full-gradient evaluations, and the bottom row contains plots against the wall-clock time. The left column is for the sonar dataset, the middle column is for the a1a dataset and the rightmost column is for the a9a dataset, all obtained from LIBSVM (Chang & Lin, 2011).

5. Numerical Experiments and Discussion

To verify the effectiveness of our proposed algorithms, we conducted a set of numerical experiments to demonstrate that both A-CODER (Algorithm 1) and VR-A-CODER (Algorithm 2) almost completely outperform other comparable block-coordinate descent methods in terms of both iteration count and wall-clock time. In particular, we compare against a number of representative methods: CODER (Song & Diakonikolas, 2021), RCDM, ACDM (Nesterov, 2012), ABCGD (Beck & Tretushvili, 2013) and APCG (Lin et al., 2015). For all the methods, we use the function value gap $f(\mathbf{x}) - f(\mathbf{x}^*)$ as the performance measure and we plot our results against the total number of full-gradient evaluations and against wall-clock time in seconds. We implement our experiments in Julia, a high performance programming language designed for numerical analysis and computational science, while optimizing all implementations to the best of our ability. Our code can be found at <https://github.com/ericlincc/Accelerated-CODER>. We set the block size to one in all the experiments, i.e., each block corresponds to one coordinate. We discussed in Section 4 that in theory it is optimal to choose $K = \lfloor n \rfloor$ in order to balance the computational costs of outer loop and inner loop in VR-A-CODER. We observed in our experiments that it is beneficial to choose K to be slightly smaller than n ($K \approx n/10$) to balance the computational time and the number of full-gradient evaluations.

We consider instances of ℓ_2 -norm (Ridge), ℓ_1 -norm (LASSO) ($\gamma = 0$) and elastic net ($\gamma > 0$) regularized lo-

gistic regression problems using three LIBSVM datasets: sonar, a1a and a9a. In the ridge regularized logistic regression problem (Figure 2), we use $\lambda_2 = 10^{-5}$ for sonar dataset and $\lambda_2 = 10^{-4}$ for a1a and a9a datasets. In the elastic net regularized logistic regression problem (Figure 3), we use $\lambda_1 = \lambda_2 = 10^{-5}$ for sonar dataset and $\lambda_1 = \lambda_2 = 10^{-4}$ for a1a and a9a datasets. In the ℓ_1 -norm regularized logistic regression problem (Figure 4), we use $\lambda_1 = 10^{-5}$ for sonar dataset and $\lambda_1 = 10^{-4}$ for a1a and a9a datasets. Figures 3 and 4 provide performance comparisons between algorithms considered in terms of the number of full-gradient evaluations and wall-clock time for the elastic net regularized logistic regression problems. We search for the best L or M for each algorithm individually at intervals of 2^i for $i \geq \mathbb{Z}$, and display the best performing runs in the plots. As predicted by our theoretical results, A-CODER and VR-A-CODER exhibit accelerated convergence rates and improved dependence on the number of blocks m even in the worst case, outperforming all other algorithms. In terms of wall-clock time, due to different per-iteration cost of each algorithm in practice, we see a mildly different set of convergence behaviors. However, A-CODER and VR-A-CODER still both perform significantly better than comparable methods.

Combined with the best known theoretical convergence rates guarantee, we believe that this work provides strong supporting arguments for cyclic methods in modern machine learning applications.

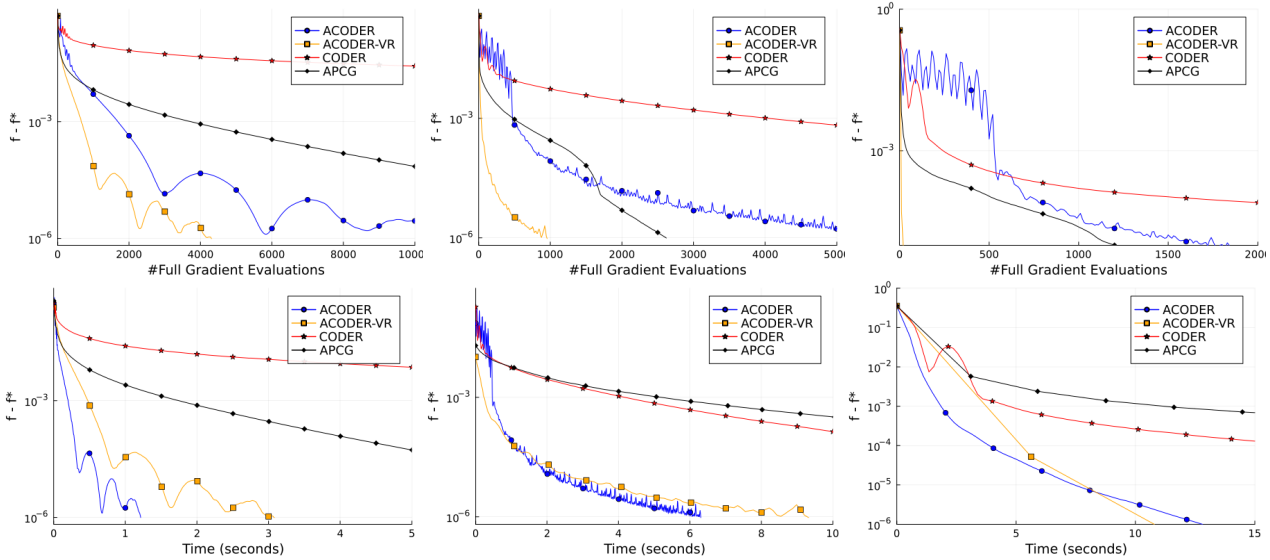


Figure 4. Performance comparisons between various algorithms in terms of number of full-gradient evaluations and wall-clock time for logistic regression with LASSO regularized problems. The top row contains plots against the number of full-gradient evaluations, and the bottom two contains plots against wall-clock time. The left column is on sonar dataset, the middle column is on a1a dataset and the rightmost column is on a9a dataset.

Acknowledgments

This research was supported in part by the NSF awards 2007757 and 2023239, and by the Office of Naval Research under contract number N00014-22-1-2348.

References

- Alacaoglu, A., Dinh, Q. T., Fercoq, O., and Cevher, V. Smooth primal-dual coordinate descent algorithms for nonsmooth convex optimization. In *Proc. NIPS'17*, 2017.
- Alacaoglu, A., Fercoq, O., and Cevher, V. Random extrapolation for primal-dual coordinate descent. In *Proc. ICML'20*, 2020.
- Allen-Zhu, Z. Katyusha: The first direct acceleration of stochastic gradient methods. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 1200–1205. ACM, 2017.
- Allen-Zhu, Z., Qu, Z., Richtárik, P., and Yuan, Y. Even faster accelerated coordinate descent using non-uniform sampling. In *Proc. ICML'16*, 2016.
- Beck, A. and Teboulle, M. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- Beck, A. and Tetrushvili, L. On the convergence of block coordinate descent type methods. *SIAM Journal on Optimization*, 23(4):2037–2060, 2013.
- Cai, X., Song, C., Wright, S. J., and Diakonikolas, J. Cyclic block coordinate descent with variance reduction for composite nonconvex optimization. *ArXiv*, abs/2212.05088, 2022.
- Chambolle, A., Ehrhardt, M. J., Richtárik, P., and Schonlieb, C.-B. Stochastic primal-dual hybrid gradient algorithm

- with arbitrary sampling and imaging applications. *SIAM Journal on Optimization*, 28(4):2783–2808, 2018.
- Chang, C.-C. and Lin, C.-J. Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.
- Chen, J. and Gu, Q. Accelerated stochastic block coordinate gradient descent for sparsity constrained nonconvex optimization. In *Conference on Uncertainty in Artificial Intelligence*, 2016.
- Chow, Y. T., Wu, T., and Yin, W. Cyclic coordinate-update algorithms for fixed-point problems: Analysis and applications. *SIAM Journal on Scientific Computing*, 39(4): A1280–A1300, 2017.
- Cohen, M., Diakonikolas, J., and Orecchia, L. On acceleration with noise-corrupted gradients. In *Proc. ICML’18*, 2018.
- Defazio, A., Bach, F., and Lacoste-Julien, S. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in neural information processing systems*, pp. 1646–1654, 2014.
- Diakonikolas, J. and Guzmán, C. Complementary composite minimization, small gradients in general norms, and applications to regression problems. *arXiv preprint arXiv:2101.11041*, 2021.
- Diakonikolas, J. and Orecchia, L. Alternating randomized block coordinate descent. In *Proc. ICML’18*, 2018.
- Diakonikolas, J. and Orecchia, L. The approximate duality gap technique: A unified theory of first-order methods. *SIAM Journal on Optimization*, 29(1):660–689, 2019.
- Fercoq, O. and Richtárik, P. Accelerated, parallel, and proximal coordinate descent. *SIAM Journal on Optimization*, 25(4):1997–2023, 2015.
- Friedman, J., Hastie, T., and Tibshirani, R. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010.
- Gasnikov, A. V. and Nesterov, Y. E. Universal method for stochastic composite optimization problems. *Computational Mathematics and Mathematical Physics*, 58(1): 48–64, 2018.
- Gürbüzbalaban, M., Ozdaglar, A., Parrilo, P. A., and Vanli, N. D. When cyclic coordinate descent outperforms randomized coordinate descent. In *Proc. NIPS’17*, 2017.
- Johnson, R. and Zhang, T. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in neural information processing systems*, pp. 315–323, 2013.
- Lee, C.-P. and Wright, S. J. Random permutations fix a worst case for cyclic coordinate descent. *IMA Journal of Numerical Analysis*, 39(3):1246–1275, 2019.
- Lei, J. and Shanbhag, U. V. Asynchronous variance-reduced block schemes for composite non-convex stochastic optimization: block-specific steplengths and adapted batch-sizes. *Optimization Methods and Software*, 37:264 – 294, 2018.
- Lei, L., Ju, C., Chen, J., and Jordan, M. I. Non-convex finite-sum optimization via SCSG methods. In *Proc. NIPS’17*, 2017.
- Lin, Q., Lu, Z., and Xiao, L. An accelerated randomized proximal coordinate gradient method and its application to regularized empirical risk minimization. *SIAM Journal on Optimization*, 25(4):2244–2273, 2015.
- Liu, J., Wright, S., Ré, C., Bittorf, V., and Sridhar, S. An asynchronous parallel stochastic coordinate descent algorithm. In *Proc. ICML’14*, 2014.
- Mazumder, R., Friedman, J. H., and Hastie, T. Sparsenet: Coordinate descent with nonconvex penalties. *Journal of the American Statistical Association*, 106(495):1125–1138, 2011.
- Nesterov, Y. A method for unconstrained convex minimization problem with the rate of convergence $\mathcal{O}(1/k^2)$. In *Doklady AN USSR*, volume 269, pp. 543–547, 1983.
- Nesterov, Y. Gradient methods for minimizing composite objective function, 2007.
- Nesterov, Y. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- Reddi, S. J., Hefny, A., Sra, S., Póczos, B., and Smola, A. Stochastic variance reduction for nonconvex optimization. In *Proc. ICML’16*, 2016.
- Richtárik, P. and Takáč, M. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, 156(1-2):433–484, 2016.
- Schmidt, M., Le Roux, N., and Bach, F. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1):83–112, 2017.
- Song, C. and Diakonikolas, J. Cyclic coordinate dual averaging with extrapolation for generalized variational inequalities. *arXiv preprint arXiv:2102.13244*, 2021.
- Song, C., Jiang, Y., and Ma, Y. Variance reduction via accelerated dual averaging for finite-sum optimization. *Advances in Neural Information Processing Systems*, 33, 2020.

- Song, C., Jiang, Y., and Ma, Y. Unified acceleration of high-order algorithms under general hölder continuity. *SIAM Journal on Optimization*, 31(3):1797–1826, 2021a.
- Song, C., Lin, C. Y., Wright, S. J., and Diakonikolas, J. Coordinate linear variance reduction for generalized linear programming. *arXiv preprint arXiv:2111.01842*, 2021b.
- Sun, R. and Ye, Y. Worst-case complexity of cyclic coordinate descent: $O(n^2)$ gap with randomized version. *Mathematical Programming*, pp. 1–34, 2019.
- Wright, S. and Lee, C.-p. Analyzing random permutations for cyclic coordinate descent. *Mathematics of Computation*, 89(325):2217–2248, 2020.
- Wright, S. J. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.
- Wu, T. T., Lange, K., et al. Coordinate descent algorithms for lasso penalized regression. *Annals of Applied Statistics*, 2(1):224–244, 2008.
- Xu, Y. and Yin, W. Block stochastic gradient iteration for convex and nonconvex optimization. *ArXiv*, abs/1408.2597, 2014.
- Zhang, Y. and Lin, X. Stochastic primal-dual coordinate method for regularized empirical risk minimization. In *Proc. ICML’15*, 2015.

Supplementary Material for Accelerated Cyclic Coordinate Dual Averaging with Extrapolation for Composite Convex Optimization

Outline. The appendix of the paper is organized as follows:

- Section A presents the proofs related to the A-CODER algorithm in the main body of the paper, as well as the implementable and adaptive versions of A-CODER.
- Section B presents the proofs related to the A-CODER-VR algorithm in the main body of the paper. We also include the implementable and adaptive versions of A-CODER-VR in this section.

A. Omitted Proofs and Pseudocode for A-CODER

Lemma 1. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex and smooth function whose gradients satisfy Assumption 2. Then, $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$:

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &\leq h r f(\mathbf{x}, \mathbf{y}) - \mathbf{x}^T \left(\mathbf{y} - \mathbf{x} \right) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2, \\ r f(\mathbf{y}) - r f(\mathbf{x}) &\leq k^2 \|\mathbf{y} - \mathbf{x}\|^2 + 2L(f(\mathbf{y}) - f(\mathbf{x}) - h r f(\mathbf{x}, \mathbf{y}) - \mathbf{x}^T (\mathbf{y} - \mathbf{x})). \end{aligned}$$

Proof. Let $\mathbf{z}_j = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(j)}, \mathbf{y}^{(j+1)}, \dots, \mathbf{y}^{(m)})$ and observe that $\mathbf{z}_m = \mathbf{x}$ and $\mathbf{z}_0 = \mathbf{y}$. Then we have

$$f(\mathbf{y}) - f(\mathbf{x}) = \sum_{j=1}^m (f(\mathbf{z}_{j-1}) - f(\mathbf{z}_j)). \quad (9)$$

As f is continuously differentiable and \mathbf{z}_j and \mathbf{z}_{j-1} only differ over the j^{th} block, we further have, by Taylor's theorem,

$$\begin{aligned} f(\mathbf{z}_{j-1}) - f(\mathbf{z}_j) &= \int_0^1 h r f(\mathbf{z}_j + t(\mathbf{z}_{j-1} - \mathbf{z}_j)), \mathbf{z}_{j-1} - \mathbf{z}_j dt \\ &= \int_0^1 \left\langle r^{(j)} f(\mathbf{z}_j + t(\mathbf{z}_{j-1} - \mathbf{z}_j)), \mathbf{y}^{(j)} - \mathbf{x}^{(j)} \right\rangle dt \\ &= \left\langle r^{(j)} f(\mathbf{x}), \mathbf{y}^{(j)} - \mathbf{x}^{(j)} \right\rangle + \int_0^1 \left\langle r^{(j)} f(\mathbf{z}_j + t(\mathbf{z}_{j-1} - \mathbf{z}_j)) - r^{(j)} f(\mathbf{x}), \mathbf{y}^{(j)} - \mathbf{x}^{(j)} \right\rangle dt. \end{aligned} \quad (10)$$

Using Young's inequality, we have, for any $\alpha > 0$,

$$\begin{aligned} &\left\langle r^{(j)} f(\mathbf{z}_j + t(\mathbf{z}_{j-1} - \mathbf{z}_j)) - r^{(j)} f(\mathbf{x}), \mathbf{y}^{(j)} - \mathbf{x}^{(j)} \right\rangle \\ &\leq \frac{\alpha}{2} k r^{(j)} f(\mathbf{z}_j + t(\mathbf{z}_{j-1} - \mathbf{z}_j)) - r^{(j)} f(\mathbf{x}) k^2 + \frac{1}{2\alpha} k \mathbf{y}^{(j)} - \mathbf{x}^{(j)} k^2 \\ &\leq \frac{\alpha}{2} k \mathbf{z}_j + t(\mathbf{z}_{j-1} - \mathbf{z}_j) - \mathbf{x} k_{\mathcal{O}_j}^2 + \frac{1}{2\alpha} k \mathbf{y}^{(j)} - \mathbf{x}^{(j)} k^2 \\ &\leq \frac{\alpha}{2} \left[(1-t) k \mathbf{z}_j - \mathbf{x} k_{\mathcal{O}_j}^2 + t k \mathbf{z}_{j-1} - \mathbf{x} k_{\mathcal{O}_j}^2 \right] + \frac{1}{2\alpha} k \mathbf{y}^{(j)} - \mathbf{x}^{(j)} k^2, \end{aligned}$$

where the second inequality is by our block Lipschitz assumption from Assumption 2 and the last line is by Jensen's inequality. Now observe that \mathbf{z}_j and \mathbf{x} agree on the first j blocks. Thus, we can write $\mathbf{z}_j - \mathbf{x} = (\mathbf{y} - \mathbf{x})_{j+1}$ and $\mathbf{z}_{j-1} - \mathbf{x} = (\mathbf{y} - \mathbf{x})_j$, while noting that we have $k(\mathbf{y} - \mathbf{x})_j k_{\mathcal{O}_j}^2 = k \mathbf{y} - \mathbf{x} k_{(\mathcal{O}_j)_j}^2$ and $k(\mathbf{y} - \mathbf{x})_{j+1} k_{\mathcal{O}_j}^2 = k \mathbf{y} - \mathbf{x} k_{(\mathcal{O}_j)_{j+1}}^2$. So by combining with Eq. (10) and integrating over t , we have, $\forall \alpha > 0$,

$$\begin{aligned} f(\mathbf{z}_{j-1}) - f(\mathbf{z}_j) &\leq \left\langle r^{(j)} f(\mathbf{x}), \mathbf{y}^{(j)} - \mathbf{x}^{(j)} \right\rangle + \frac{1}{2\alpha} k \mathbf{y}^{(j)} - \mathbf{x}^{(j)} k^2 \\ &\quad + \frac{\alpha}{4} \left(k \mathbf{y} - \mathbf{x} k_{(\mathcal{O}_j)_j}^2 + k \mathbf{y} - \mathbf{x} k_{(\mathcal{O}_j)_{j+1}}^2 \right). \end{aligned} \quad (11)$$

Summing Eq. (11) over $j \in [m]$ and using the definition of Mahalanobis norm, we finally get

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &= \sum_{j=1}^m (f(\mathbf{z}_{j-1}) - f(\mathbf{z}_j)) \\ &\leq \text{hr } f(\mathbf{x}), \mathbf{y} - \mathbf{x} + \frac{\alpha}{4} (\mathbf{y} - \mathbf{x})^T \mathcal{Q} (\mathbf{y} - \mathbf{x}) + \frac{1}{2\alpha} k \mathbf{y} - \mathbf{x} k^2 \\ &\leq \text{hr } f(\mathbf{x}), \mathbf{y} - \mathbf{x} + \left(\frac{1}{2\alpha} + \frac{\alpha L^2}{8} \right) k \mathbf{y} - \mathbf{x} k^2, \end{aligned}$$

where we used Holder's inequality and the definition of L in Assumption 2. Letting $\alpha = \frac{2}{L}$ now completes the proof of the first part.

The second part of the proof is standard and is provided for completeness. Let \mathbf{x}, \mathbf{y} be any two points from \mathbb{R}^d . Define $h_{\mathbf{x}}(\mathbf{y}) = f(\mathbf{y}) - \text{hr } f(\mathbf{x}), \mathbf{y}$. Observe that $h_{\mathbf{x}}(\mathbf{y})$ is convex (as the sum of a convex function $f(\mathbf{y})$ and a linear function $-\text{hr } f(\mathbf{x}), \mathbf{y}$) and is minimized at $\mathbf{y} = \mathbf{x}$ (as for any $\mathbf{y} \in \mathbb{R}^d$, $h_{\mathbf{x}}(\mathbf{y}) - h_{\mathbf{x}}(\mathbf{x}) = f(\mathbf{y}) - f(\mathbf{x}) - \text{hr } f(\mathbf{x}), \mathbf{y} - \mathbf{x} \geq 0$, by convexity of f). Observe further that for any $\mathbf{y}, \mathbf{z} \in \mathbb{R}^d$, we have

$$\begin{aligned} h_{\mathbf{x}}(\mathbf{y}) - h_{\mathbf{x}}(\mathbf{z}) &\leq \text{hr } h_{\mathbf{x}}(\mathbf{z}), \mathbf{y} - \mathbf{z} = f(\mathbf{y}) - f(\mathbf{z}) - \text{hr } f(\mathbf{z}), \mathbf{y} - \mathbf{z} \\ &\leq \frac{L}{2} k \mathbf{y} - \mathbf{z} k^2, \end{aligned}$$

where the last inequality is by the first part of the proof. The last inequality and the fact that \mathbf{x} minimizes $h_{\mathbf{x}}$ now allow us to conclude that

$$\begin{aligned} h_{\mathbf{x}}(\mathbf{x}) - h_{\mathbf{x}}(\mathbf{y}) &\leq \frac{1}{L} \text{hr } h_{\mathbf{x}}(\mathbf{y}) \\ h_{\mathbf{x}}(\mathbf{y}) &\leq \frac{1}{2L} k \text{hr } h_{\mathbf{x}}(\mathbf{y}) k^2. \end{aligned}$$

To complete the proof, it remains to plug the definition of $h_{\mathbf{x}}(\cdot)$ into the last inequality, and rearrange. \square

Lemma 2. For any $\mathbf{u} \in \mathbb{R}^d$ and any sequence of vectors $\{ \mathbf{q}_i, \mathbf{v}_i \}_{i=1}^k$, we have

$$A_k (f(\mathbf{y}_k) - f(\mathbf{u})) \tag{3}$$

$$\leq \sum_{i=1}^k E_i(\mathbf{u}) + \frac{1}{2} k \mathbf{u} - \mathbf{x}_0 k^2 + \frac{1 + A_k \gamma}{2} k \mathbf{u} - \mathbf{v}_k k^2, \tag{4}$$

where

$$\begin{aligned} E_i(\mathbf{u}) &= A_i (f(\mathbf{y}_i) - f(\mathbf{x}_i)) - A_{i-1} (f(\mathbf{y}_{i-1}) - f(\mathbf{x}_{i-1})) \\ &\quad + a_i \text{hr } \mathbf{q}_i, \mathbf{v}_i - \mathbf{x}_i + a_i \text{hr } f(\mathbf{x}_i) - \mathbf{q}_i, \mathbf{x}_i - \mathbf{u} \\ &\leq \frac{1 + A_{i-1} \gamma}{2} k \mathbf{v}_i - \mathbf{v}_{i-1} k^2. \end{aligned} \tag{5}$$

Proof. As $\mathbf{y}_k = \frac{1}{A_k} \sum_{i=1}^k a_i \mathbf{v}_i$ and g is convex, we have $g(\mathbf{y}_k) \leq \frac{1}{A_k} \sum_{i=1}^k a_i g(\mathbf{v}_i)$ and thus,

$$\begin{aligned} A_k f(\mathbf{y}_k) &\leq A_k f(\mathbf{y}_k) + \sum_{i=1}^k a_i g(\mathbf{v}_i) \\ &= \sum_{i=1}^k (A_i f(\mathbf{y}_i) - A_{i-1} f(\mathbf{y}_{i-1})) + \sum_{i=1}^k a_i g(\mathbf{v}_i), \end{aligned} \tag{12}$$

where the equality is by $A_0 = 0$. Then, as f is convex and $f = f + g$, we have, $\forall \mathbf{u}$,

$$\begin{aligned}
 A_k f(\mathbf{u}) &= \sum_{i=1}^k a_i f(\mathbf{u}) - \sum_{i=1}^k a_i (f(\mathbf{x}_i) + h r f(\mathbf{x}_i), \mathbf{u} - \mathbf{x}_i) + g(\mathbf{u}) \\
 &= \sum_{i=1}^k a_i (f(\mathbf{x}_i) + h \mathbf{q}_i, \mathbf{u} - \mathbf{x}_i) + g(\mathbf{u}) + \sum_{i=1}^k a_i h r f(\mathbf{x}_i) - \mathbf{q}_i, \mathbf{u} - \mathbf{x}_i \\
 &= \psi_k(\mathbf{u}) - \psi_0(\mathbf{u}) + \sum_{i=1}^k a_i h r f(\mathbf{x}_i) - \mathbf{q}_i, \mathbf{u} - \mathbf{x}_i \\
 &\leq \psi_k(\mathbf{v}_k) + \frac{1 + A_k \gamma}{2} k \mathbf{u} - \mathbf{v}_k k^2 + \frac{1}{2} k \mathbf{u} - \mathbf{x}_0 k^2 \\
 &\quad + \sum_{i=1}^k a_i h r f(\mathbf{x}_i) - \mathbf{q}_i, \mathbf{u} - \mathbf{x}_i,
 \end{aligned}$$

where the first inequality is by the convexity of f , the third equality is by the recursive definition of $\psi_k(\mathbf{u})$, and the last inequality is by the $(1 + A_k \gamma)$ -strong convexity of $\psi_k(\mathbf{u})$, $\mathbf{v}_k = \arg \min_{\mathbf{u}} \psi_k(\mathbf{u})$ which implies $\psi_k(\mathbf{u}) \leq \psi_k(\mathbf{v}_k) + \frac{1 + A_k \gamma}{2} k \mathbf{u} - \mathbf{v}_k k^2$, and the definition of $\psi_0(\mathbf{u})$.

Then as $\psi_0(\mathbf{v}_0) = 0$, using the recursive definition of ψ_k , we have

$$\begin{aligned}
 \psi_k(\mathbf{v}_k) &= \sum_{i=1}^k (\psi_i(\mathbf{v}_i) - \psi_{i-1}(\mathbf{v}_{i-1})) \\
 &= \sum_{i=1}^k \left((\psi_{i-1}(\mathbf{v}_i) - \psi_{i-1}(\mathbf{v}_{i-1})) + a_i (f(\mathbf{x}_i) + h \mathbf{q}_i, \mathbf{v}_i - \mathbf{x}_i) + g(\mathbf{v}_i) \right) \\
 &\leq \sum_{i=1}^k \left(\frac{1 + A_i \gamma}{2} k \mathbf{v}_i - \mathbf{v}_{i-1} k^2 + a_i (f(\mathbf{x}_i) + h \mathbf{q}_i, \mathbf{v}_i - \mathbf{x}_i) + g(\mathbf{v}_i) \right), \tag{13}
 \end{aligned}$$

where the last inequality is by the $(1 + A_i \gamma)$ -strong convexity of ψ_{i-1} and the optimality of \mathbf{v}_{i-1} . Combining Eqs. (12) and (13), we have

$$A_k (f(\mathbf{y}_k) - f(\mathbf{u})) \leq \sum_{i=1}^k E_i(\mathbf{u}) - \frac{1 + A_k \gamma}{2} k \mathbf{u} - \mathbf{v}_k k^2 + \frac{1}{2} k \mathbf{u} - \mathbf{x}_0 k^2, \tag{14}$$

where $E_i(\mathbf{u})$ is defined in (5). □

Lemma 3. Let $\mathbf{x}_0 \in \text{dom}(g)$ be an arbitrary initial point and consider the updates in Algorithm 1. If, for $k \geq 1$, $\frac{a_k^2}{A_k} \leq \frac{2(1 + A_{k-1})}{5L}$, then $\forall \mathbf{u}$,

$$\begin{aligned}
 E_k(\mathbf{u}) &\leq a_k h r f(\mathbf{x}_k) - \mathbf{p}_k, \mathbf{v}_k - \mathbf{u} \\
 &\quad + a_{k-1} h r f(\mathbf{x}_{k-1}) - \mathbf{p}_{k-1}, \mathbf{v}_{k-1} - \mathbf{u} \\
 &\quad + \frac{1 + A_{k-1} \gamma}{10} k \mathbf{v}_k - \mathbf{v}_{k-1} k^2 \\
 &\quad + \frac{1 + A_{k-2} \gamma}{10} k \mathbf{v}_{k-1} - \mathbf{v}_{k-2} k^2.
 \end{aligned}$$

Proof. By the convexity of f , we have

$$f(\mathbf{y}_{k-1}) \leq f(\mathbf{x}_k) + h r f(\mathbf{x}_k), \mathbf{y}_{k-1} - \mathbf{x}_k.$$

Then by applying Lemma 1, we have

$$\begin{aligned} A_k(f(\mathbf{y}_k) - f(\mathbf{x}_k)) - A_{k-1}(f(\mathbf{y}_{k-1}) - f(\mathbf{x}_k)) &= h\Gamma f(\mathbf{x}_k), A_k \mathbf{y}_k - A_{k-1} \mathbf{y}_{k-1} - a_k \mathbf{x}_{k-1} + \frac{A_k L}{2} k \mathbf{y}_k - \mathbf{x}_k k^2 \\ &= a_k h\Gamma f(\mathbf{x}_k), \mathbf{v}_k - \mathbf{x}_{k-1} + \frac{a_k^2 L}{2 A_k} k \mathbf{v}_k - \mathbf{v}_{k-1} k^2, \end{aligned} \quad (15)$$

where we used the definitions of \mathbf{y}_k and \mathbf{x}_k from Algorithm 1 in the last equality. Combining Eq. (5) (with $i = k$) in Lemma 2 and Eq. (15), we have

$$E_k(\mathbf{u}) = \left(\frac{a_k^2 L}{2 A_k} - \frac{1 + A_k \Gamma}{2} \right) k \mathbf{v}_k - \mathbf{v}_{k-1} k^2 + a_k h\Gamma f(\mathbf{x}_k) - \mathbf{q}_k, \mathbf{v}_k - \mathbf{u}. \quad (16)$$

Thus by rewriting the second term as the sum of inner products over the m blocks and by using the definition of $\mathbf{q}_k^{(j)}$ in Algorithm 1, we have

$$\begin{aligned} a_k h\Gamma f(\mathbf{x}_k) - \mathbf{q}_k, \mathbf{v}_k - \mathbf{u} &= a_k \sum_{j=1}^m \left\langle r^{(j)} f(\mathbf{x}_k) - \mathbf{q}_k^{(j)}, \mathbf{v}_k^{(j)} - \mathbf{u}^{(j)} \right\rangle \\ &= \sum_{j=1}^m \left[a_k \left\langle r^{(j)} f(\mathbf{x}_k) - \mathbf{p}_k^{(j)}, \mathbf{v}_k^{(j)} - \mathbf{u}^{(j)} \right\rangle - a_{k-1} \left\langle r^{(j)} f(\mathbf{x}_{k-1}) - \mathbf{p}_{k-1}^{(j)}, \mathbf{v}_{k-1}^{(j)} - \mathbf{u}^{(j)} \right\rangle \right] \\ &\quad + a_{k-1} \sum_{j=1}^m \left\langle r^{(j)} f(\mathbf{x}_{k-1}) - \mathbf{p}_{k-1}^{(j)}, \mathbf{v}_{k-1}^{(j)} - \mathbf{v}_k^{(j)} \right\rangle. \end{aligned} \quad (17)$$

Notice that the first two inner product terms in the first line of Eq. (17) telescope when summed over k , therefore it remains to bound $a_{k-1} \sum_{j=1}^m \left\langle r^{(j)} f(\mathbf{x}_{k-1}) - \mathbf{p}_{k-1}^{(j)}, \mathbf{v}_{k-1}^{(j)} - \mathbf{v}_k^{(j)} \right\rangle$. In particular we let $\mathbf{w}_{k,j} = (\mathbf{x}_{k-1}^1, \dots, \mathbf{x}_{k-1}^j, \mathbf{y}_{k-1}^{j+1}, \dots, \mathbf{y}_{k-1}^m)$ so that $\mathbf{p}_{k-1}^{(j)} = r^{(j)} f(\mathbf{w}_{k,j})$, then we have

$$\begin{aligned} \left\langle r^{(j)} f(\mathbf{x}_{k-1}) - \mathbf{p}_{k-1}^{(j)}, \mathbf{v}_{k-1}^{(j)} - \mathbf{v}_k^{(j)} \right\rangle &= \left\langle r^{(j)} f(\mathbf{x}_{k-1}) - r^{(j)} f(\mathbf{w}_{k-1,j}), \mathbf{v}_{k-1}^{(j)} - \mathbf{v}_k^{(j)} \right\rangle \\ &\leq \frac{\alpha}{2} \left\| r^{(j)} f(\mathbf{x}_{k-1}) - r^{(j)} f(\mathbf{w}_{k-1,j}) \right\|^2 + \frac{1}{2\alpha} \left\| \mathbf{v}_{k-1}^{(j)} - \mathbf{v}_k^{(j)} \right\|^2 \\ &\leq \frac{\alpha}{2} k \mathbf{x}_{k-1} - \mathbf{w}_{k-1,j} k_{\mathcal{Q}}^2 + \frac{1}{2\alpha} \left\| \mathbf{v}_{k-1}^{(j)} - \mathbf{v}_k^{(j)} \right\|^2 \end{aligned} \quad (18)$$

where the first inequality holds for any $\alpha > 0$ by Young's inequality and the second inequality is by Assumption 2. Notice that \mathbf{x}_{k-1} and $\mathbf{w}_{k-1,j}$ agree on the first j blocks, so similar to the proof of Lemma 1 we can write $\mathbf{x}_{k-1} - \mathbf{w}_{k-1,j} = (\mathbf{y}_{k-1} - \mathbf{x}_{k-1})_{j+1}$. Therefore by applying similar arguments as Lemma 1, we get

$$\begin{aligned} \sum_{j=1}^m k \mathbf{x}_{k-1} - \mathbf{w}_{k-1,j} k_{\mathcal{Q}}^2 &= \sum_{j=1}^m k \mathbf{y}_{k-1} - \mathbf{x}_{k-1} k_{(\mathcal{Q})_{j+1}}^2 \\ &= \sum_{j=1}^m k \mathbf{y}_{k-1} - \mathbf{x}_{k-1} k_{(\mathcal{Q})_{j+1}}^2 + \sum_{j=1}^m k \mathbf{y}_{k-1} - \mathbf{x}_{k-1} k_{(\mathcal{Q})_j}^2 \\ &= k \mathbf{y}_{k-1} - \mathbf{x}_{k-1} k_{\mathcal{Q}}^2 \\ &\leq \frac{a_{k-1}^2 L^2}{2 A_{k-1}^2} k \mathbf{v}_{k-1} - \mathbf{v}_{k-2} k^2 \end{aligned} \quad (19)$$

where we used the non-negativity of Mahalanobis norm w.r.t. semi-positive definite matrix in the first inequality and the definition of $\mathbf{x}_k, \mathbf{y}_k$ and L in the last inequality. Lastly, by combining Eqs. (16), (17), (18) and (19), we have

$$E_k(\mathbf{u}) = a_k \text{hr} f(\mathbf{x}_k) - \langle \mathbf{p}_k, \mathbf{v}_k \rangle \langle \mathbf{u}, \mathbf{i} \rangle - a_{k-1} \text{hr} f(\mathbf{x}_{k-1}) - \langle \mathbf{p}_{k-1}, \mathbf{v}_{k-1} \rangle \langle \mathbf{u}, \mathbf{i} \rangle \\ + \left(\frac{a_k^2 L}{2A_k} - \frac{1 + A_{k-1}\gamma}{2} + \frac{a_{k-1}}{2\alpha} \right) k \|\mathbf{v}_k - \mathbf{v}_{k-1}\|^2 + \left(\frac{\alpha a_{k-1}^3 L^2}{4A_{k-1}^2} \right) k \|\mathbf{v}_{k-1} - \mathbf{v}_{k-2}\|^2.$$

It remains to choose $\alpha = \frac{A_{k-1}}{a_{k-1}L}$ and some sequence $\{a_i\}_i^k$ such that $\frac{a_k^2}{A_k} = \frac{2(1+A_{k-1})}{5L}$. \square

Theorem 1. Let $\mathbf{x}_0 \in \text{dom}(g)$ be an arbitrary initial point and consider the updates in Algorithm 1. Then, $\forall k \geq 1$ and any $\mathbf{u} \in \text{dom}(g)$:

$$f(\mathbf{y}_k) - f(\mathbf{u}) + \frac{3(1 + A_{k-1}\gamma)}{10A_k} k \|\mathbf{u} - \mathbf{v}_k\|^2 \leq \frac{k \|\mathbf{u} - \mathbf{x}_0\|^2}{2A_k}.$$

In particular, if $\mathbf{x} = \arg \min_{\mathbf{x}} f(\mathbf{x})$ exists, then

$$f(\mathbf{y}_k) - f(\mathbf{x}) \leq \frac{k \|\mathbf{x} - \mathbf{x}_0\|^2}{2A_k}.$$

Further, in this case we also have:

$$k \|\mathbf{v}_k - \mathbf{x}\|^2 \leq \frac{5}{3(1 + A_{k-1}\gamma)} k \|\mathbf{x}_0 - \mathbf{x}\|^2, \\ k \|\mathbf{y}_k - \mathbf{x}\|^2 \leq \left(\frac{5}{3A_k} \sum_{i=1}^k \frac{a_i}{1 + A_{i-1}\gamma} \right) k \|\mathbf{x}_0 - \mathbf{x}\|^2.$$

Finally, in all the bounds we have

$$A_k \leq \max \left\{ \frac{2}{5L} \left(1 + \sqrt{\frac{2\gamma}{5L}} \right)^k, \frac{k^2}{10L} \right\}.$$

Proof. By Lemma 3, and using the fact $A_0 = a_0 = 0$ and $\mathbf{v}_0 = \mathbf{v}_{-1}$, we have

$$\sum_{i=1}^k E_i(\mathbf{u}) \leq \frac{1 + A_{k-1}\gamma}{10} k \|\mathbf{v}_k - \mathbf{v}_{k-1}\|^2 + a_k \text{hr} f(\mathbf{x}_k) - \langle \mathbf{p}_k, \mathbf{v}_k \rangle \langle \mathbf{u}, \mathbf{i} \rangle. \quad (20)$$

Same as in the proof of Lemma 3, we can bound $a_k \text{hr} f(\mathbf{x}_k) - \langle \mathbf{p}_k, \mathbf{v}_k \rangle \langle \mathbf{u}, \mathbf{i} \rangle$ using Young's inequality and the definition of smoothness for f . In particular, for any $\alpha > 0$,

$$a_k \text{hr} f(\mathbf{x}_k) - \langle \mathbf{p}_k, \mathbf{v}_k \rangle \langle \mathbf{u}, \mathbf{i} \rangle = a_k \left(\frac{\alpha L^2}{4} k \|\mathbf{y}_k - \mathbf{x}_k\|^2 + \frac{1}{2\alpha} k \|\mathbf{u} - \mathbf{v}_k\|^2 \right) \\ = a_k \left(\frac{\alpha L^2 a_k^2}{4A_k^2} k \|\mathbf{v}_k - \mathbf{v}_{k-1}\|^2 + \frac{1}{2\alpha} k \|\mathbf{u} - \mathbf{v}_k\|^2 \right).$$

Choosing $\alpha = \frac{A_k}{a_k L}$ and using $\frac{a_k^2}{A_k} = \frac{2(1+A_{k-1})}{5L}$, we get

$$a_k \text{hr} f(\mathbf{x}_k) - \langle \mathbf{p}_k, \mathbf{v}_k \rangle \langle \mathbf{u}, \mathbf{i} \rangle \leq \frac{(1 + A_{k-1}\gamma)}{10} k \|\mathbf{v}_k - \mathbf{v}_{k-1}\|^2 + \frac{(1 + A_{k-1}\gamma)}{5} k \|\mathbf{u} - \mathbf{v}_k\|^2. \quad (21)$$

Then combining Lemma 2, Eq. (21) and Eq. (20) with the fact $A_{k-1} \leq A_k$, we have

$$f(\mathbf{y}_k) - f(\mathbf{u}) + \frac{3(1 + A_{k-1}\gamma)}{10A_k} k \|\mathbf{u} - \mathbf{v}_k\|^2 \leq \frac{1}{2A_k} k \|\mathbf{u} - \mathbf{x}_0\|^2. \quad (22)$$

Assume now that $\mathbf{x} = \arg \min_{\mathbf{x}} f(\mathbf{x})$ exists. As $f(\mathbf{y}_k) - f(\mathbf{x}) \geq 0$, Eq. (22) implies

$$k \|\mathbf{v}_k - \mathbf{x}\|^2 \leq \frac{5}{3(1 + A_{k-1}\gamma)} k \|\mathbf{x}_0 - \mathbf{x}\|^2. \quad (23)$$

Algorithm 3 Adaptive Accelerated Cyclic cOordinate Dual avEraging with extRapolation (Ada-A-CODER)

```

1: Input:  $\mathbf{x}_0 \in \text{dom}(g), \gamma \geq 0, L_0 > 0, m, f \in \mathcal{S}^1, \dots, \mathcal{S}^m g$ 
2: Initialization:  $\mathbf{x}_1 = \mathbf{x}_0 = \mathbf{v}_1 = \mathbf{v}_0 = \mathbf{y}_0, \mathbf{p}_0 = \nabla f(\mathbf{x}_0), \mathbf{z}_0 = \mathbf{0}, a_0 = A_0 = 0$ 
3: for  $k = 1$  to  $K$  do
4:    $L_k = L_{k-1}/2$ 
5:   repeat
6:      $L_k = 2L_k$ 
7:     Set  $a_k > 0$  be largest value s.t.  $\frac{a_k^2}{A_k} \leq \frac{2(1+A_{k-1})}{5L_k}$  where  $A_k = A_{k-1} + a_k$ 
8:      $\mathbf{x}_k = \frac{A_{k-1}}{A_k} \mathbf{y}_{k-1} + \frac{a_k}{A_k} \mathbf{v}_{k-1}$ 
9:     for  $j = m$  to  $1$  do
10:       $\mathbf{p}_k^{(j)} = \nabla^{(j)} f(\mathbf{x}_k^{(1)}, \dots, \mathbf{x}_k^{(j)}, \mathbf{y}_k^{(j+1)}, \dots, \mathbf{y}_k^{(m)})$ 
11:       $\mathbf{q}_k^{(j)} = \mathbf{p}_k^{(j)} + \frac{a_{k-1}}{a_k} (\nabla^{(j)} f(\mathbf{x}_{k-1}) - \mathbf{p}_k^{(j-1)})$ 
12:       $\mathbf{z}_k^{(j)} = \mathbf{z}_{k-1}^{(j)} + a_k \mathbf{q}_k^{(j)}$ 
13:       $\mathbf{v}_k^{(j)} = \text{prox}_{A_k g_j}(\mathbf{x}_0^{(j)} - \mathbf{z}_k^{(j)})$ 
14:       $\mathbf{y}_k^{(j)} = \frac{A_{k-1}}{A_k} \mathbf{y}_{k-1}^{(j)} + \frac{a_k}{A_k} \mathbf{v}_k^{(j)}$ 
15:     end for
16:   until  $f(\mathbf{y}_k) - f(\mathbf{x}_k) + h r f(\mathbf{x}_k), \mathbf{y}_k - \mathbf{x}_k \leq \frac{L_k}{2} k \mathbf{y}_k - \mathbf{x}_k k^2$ 
17:   end for
18: return  $\mathbf{v}_K, \mathbf{y}_K$ 

```

Using Jensen's inequality, as $\mathbf{y}_k = \frac{1}{A_k} \sum_{i=1}^k a_i \mathbf{v}_i$, we also have from Eq. (23)

$$k \mathbf{y}_k - \mathbf{x}_k \leq k^2 \left(\frac{5}{3A_k} \sum_{i=1}^k \frac{a_i}{1 + A_{i-1} \gamma} \right) k \mathbf{x}_0 - \mathbf{x}_k k^2.$$

Finally, recall once again that $a_k g_{k-1}$ is chosen so that $\frac{a_k^2}{A_k} = \frac{2(1+A_{k-1})}{5L}$. When $\gamma = 0$, this leads to the standard $A_k \sim \frac{k^2}{10L}$ growth of accelerated algorithms by choosing $a_k = \frac{k}{5L}$ for $k \geq 1$. When $\gamma > 0$, we have $\frac{a_k}{A_{k-1}} > \sqrt{\frac{2}{5L}}$, and it remains to use that $A_k = \frac{A_k}{A_{k-1}} \dots \frac{A_2}{A_1} A_1 = A_1 \left(1 + \sqrt{\frac{2}{5L}}\right)^{k-1}$ where $a_1 = A_1 = \frac{2}{5L}$ using the choice of a_k in Algorithm 1 and $A_0 = a_0 = 0$, completing the proof. \square

A.1. (Lipschitz) Parameter-Free A-CODER

Similar to CODER, it is possible to adaptively estimate the Lipschitz parameter L_k for A-CODER. Note that in the case of A-CODER, all that is needed for the analysis from Section 3 to apply is that the quadratic bound from Lemma 1 holds between \mathbf{x}_k and \mathbf{y}_k . A variant of A-CODER that implements this adaptive estimation is provided in Algorithm 3.

B. Omitted Proofs and Pseudocode for ACODER-VR

Lemma 4. *If $f(\mathbf{x}) = \frac{1}{n} \sum_{t=1}^n f_t(\mathbf{x})$ satisfies Assumption 3, then it satisfies Assumption 2 and thus Lemma 1 holds.*

Proof. By using Jensen's inequality and Assumption 3, we have

$$\left\| \nabla^{(j)} f(\mathbf{x}) - \nabla^{(j)} f(\mathbf{y}) \right\|^2 \leq \frac{1}{n} \sum_{t=1}^n \left\| \nabla^{(j)} f_t(\mathbf{x}) - \nabla^{(j)} f_t(\mathbf{y}) \right\|^2 \leq k \mathbf{x} - \mathbf{y} k_Q^2.$$

\square

Lemma 5. *For any $\mathbf{u} \in \mathbb{R}^d$ and any sequence of vectors $f \mathbf{q}_{s,k} g_s$, $s=2:k \geq 2[K]$, for all $S \geq 2$, we have*

$$K A_S (f(\mathbf{y}_S) - f(\mathbf{u}))$$

Algorithm 4 Variance Reduced A-CODER (Analysis Version)

```

1: Input:  $\mathbf{x}_0 \in \text{dom}(g), \gamma \in [0, 1], L > 0, m, f \in \mathcal{S}^1, \dots, S^m g$ 
2: Initialization:  $\mathbf{y}_0 = \mathbf{v}_{1,0} = \mathbf{y}_{1,0} = \mathbf{x}_{1,1} = \mathbf{x}_0$ 
3:  $a_0 = A_0 = 0; A_1 = a_1 = \frac{1}{4L}$ 
4:  $\psi_{1,0}(\cdot) = \frac{K}{2} k \|\mathbf{x}_0\|^2$ 
5:  $\mathbf{v}_{1,1} = \arg \min_{\mathbf{v}} \tilde{f} \psi_{1,1}(\mathbf{v}) := \psi_{1,0}(\mathbf{v}) + K a_1 (f(\mathbf{x}_0) + h r f(\mathbf{x}_0), \mathbf{v} - \mathbf{x}_0)_i + g(\mathbf{v})_g$ 
6:  $\mathbf{w}_{1,1:j} = (\mathbf{x}_{1,1}^{(1)}, \dots, \mathbf{x}_{1,1}^{(j)}, \mathbf{y}_{1,1}^{(j+1)}, \dots, \mathbf{y}_{1,1}^{(m)})$ 
7:  $\mathbf{y}_1 = \mathbf{v}_{2,0} = \mathbf{y}_{1,1} = \mathbf{v}_{1,1}; \mathbf{w}_{2,0:j} = \mathbf{w}_{1,1:j}; \psi_{2,0} = \psi_{1,1}$ 
8: for  $s = 2$  to  $S$  do
9:   Set  $a_s > 0$  s.t.  $a_s^2 = \frac{K A_{s-1} (1 + A_{s-1})}{8L}; A_s = A_{s-1} + a_s$ 
10:   $a_{s,0} = a_{s-1}; a_{s,1} = a_{s,2} = \dots = a_{s,K} = a_s$ 
11:   $\mathbf{x}_{s,0} = \mathbf{x}_{s-1,K}; \mathbf{y}_{s,0} = \mathbf{x}_{s-1,K}; \mathbf{w}_{s,0:j} = \mathbf{w}_{s-1,K:j}; \mathbf{v}_{s,0} = \mathbf{v}_{s-1,K}; \psi_{s,0} = \psi_{s-1,K}$ 
12:   $\mathbf{y}_s = r f(\mathbf{y}_{s-1})$ 
13:  for  $k = 1$  to  $K$  do
14:     $\mathbf{x}_{s;k} = \frac{A_{s-1}}{A_s} \mathbf{y}_{s-1} + \frac{a_s}{A_s} \mathbf{v}_{s;k-1}$ 
15:    for  $j = m$  to  $1$  do
16:       $\mathbf{w}_{s;k;j} = (\mathbf{x}_{s;k}^{(1)}, \dots, \mathbf{x}_{s;k}^{(j)}, \mathbf{y}_{s;k}^{(j+1)}, \dots, \mathbf{y}_{s;k}^{(m)})$ 
17:      Choose  $t$  in  $[n]$  uniformly at random
18:       $\tilde{r}_{s;k}^{(j)} = r^{(j)} f_t(\mathbf{w}_{s;k;j}) - r^{(j)} f_t(\mathbf{y}_{s-1}) + \frac{a_s}{A_s} \tilde{r}_{s;k}^{(j)}$ 
19:       $\mathbf{q}_{s;k}^{(j)} = r_{s;k}^{(j)} + \frac{a_{s,k-1}}{a_s} (r^{(j)} f_t(\mathbf{x}_{s;k-1}) - r^{(j)} f_t(\mathbf{w}_{s;k-1;j}))$ 
20:       $\mathbf{v}_{s;k}^{(j)} = \arg \min_{\mathbf{v}^{(j)} \in \mathbb{R}^d} \tilde{f} \psi_{s;k}^{(j)}(\mathbf{v}^{(j)}) := \psi_{s;k-1}^{(j)}(\mathbf{v}^{(j)}) + a_s (\frac{1}{m} f(\mathbf{x}_{s;k}) + h \mathbf{q}_{s;k}^{(j)}, \mathbf{v}^{(j)} - \mathbf{y}_{s;k-1}^{(j)})_i + g^j(\mathbf{v}^{(j)})_g$ 
21:       $\mathbf{y}_{s;k}^{(j)} = \frac{A_{s-1}}{A_s} \mathbf{y}_{s-1} + \frac{a_s}{A_s} \mathbf{v}_{s;k}^{(j)}$ 
22:    end for
23:  end for
24:   $\mathbf{y}_s = \frac{1}{K} \sum_{k=1}^K \mathbf{y}_{s;k}$ 
25: end for
26: return  $\mathbf{v}_{S,L}, \mathbf{y}_S$ 

```

$$\frac{K}{2} k \|\mathbf{x}_0\|^2 + \frac{K(1 + A_S \gamma)}{2} k \|\mathbf{v}_{S,K}\|^2 + \frac{K}{4} k \|\mathbf{v}_{1,1}\|^2 + \sum_{s=2}^S \sum_{k=1}^K E_{S;k}(\mathbf{u}),$$

where

$$E_{S;k}(\mathbf{u}) = A_S (f(\mathbf{y}_{S;k}) - f(\mathbf{x}_{S;k})) - A_{S-1} (f(\mathbf{y}_{S-1}) - f(\mathbf{x}_{S;k})) + a_S h r f(\mathbf{x}_{S;k}) - \mathbf{q}_{S;k}, \mathbf{x}_{S;k} - \mathbf{u}_i + a_S h \mathbf{q}_{S;k}, \mathbf{x}_{S;k} - \mathbf{v}_{S;k} - \frac{K(1 + A_S \gamma)}{2} k \|\mathbf{v}_{S;k-1}\|^2. \quad (8)$$

Proof. As f is convex and $f = f + g$, we have: $\partial \mathbf{u}$,

$$\begin{aligned} K A_S f(\mathbf{u}) &= \sum_{s=1}^S \sum_{k=1}^K a_s f(\mathbf{u}) \\ &= K a_1 (f(\mathbf{x}_{1,1}) + h r f(\mathbf{x}_{1,1}), \mathbf{u} - \mathbf{x}_{1,1})_i + g(\mathbf{u}) \\ &\quad + \sum_{s=2}^S \sum_{k=1}^K a_s (f(\mathbf{x}_{S;k}) + h r f(\mathbf{x}_{S;k}), \mathbf{u} - \mathbf{x}_{S;k})_i + g(\mathbf{u}) \\ &= \psi_{S,K}(\mathbf{u}) - \psi_{1,0}(\mathbf{u}) + \sum_{s=2}^S \sum_{k=1}^K a_s h r f(\mathbf{x}_{S;k}) - \mathbf{q}_{S;k}, \mathbf{u} - \mathbf{x}_{S;k} \end{aligned}$$

$$\begin{aligned} & \psi_{S;K}(\mathbf{v}_{S;K}) + \frac{K(1+A_S\gamma)}{2}k\mathbf{u} \quad \mathbf{v}_{S;K}k^2 \quad \frac{K}{2}k\mathbf{x}_0 \quad \mathbf{u}k^2 \\ & + \sum_{s=2}^S \sum_{k=1}^K a_s h r f(\mathbf{x}_{S;k}) \quad \mathbf{q}_{S;k}, \mathbf{u} \quad \mathbf{x}_{S;k}i, \end{aligned} \quad (24)$$

where the first inequality is by the convexity of f , the second equality is by the recursive definition of $\psi_{S;K}(\mathbf{u})$, the second inequality is by the $K(1+A_S\gamma)$ -strong convexity of $\psi_{S;K}(\mathbf{u})$ and $\mathbf{v}_{S;K} = \arg \min_{\mathbf{u}} \psi_{S;K}(\mathbf{u})$ leading to $\psi_{S;K}(\mathbf{u}) \leq \psi_{S;K}(\mathbf{v}_{S;K}) + \frac{K(1+A_S\gamma)}{2}k\mathbf{u} - \mathbf{v}_{S;K}k^2$.

Then using our recursive definition of the estimate sequences again, we have

$$\begin{aligned} & \psi_{S;K}(\mathbf{v}_{S;K}) \\ & = \psi_{1,1}(\mathbf{v}_{1,1}) + \sum_{s=2}^S \sum_{k=1}^K (\psi_{s;k}(\mathbf{v}_{s;k}) - \psi_{s;k-1}(\mathbf{v}_{s;k-1})) \\ & = \psi_{1,1}(\mathbf{v}_{1,1}) + \sum_{s=2}^S \sum_{k=1}^K (\psi_{s;k-1}(\mathbf{v}_{s;k}) - \psi_{s;k-1}(\mathbf{v}_{s;k-1})) \\ & \quad + \sum_{s=2}^S \sum_{k=1}^K a_s (f(\mathbf{x}_{S;k}) + h\mathbf{q}_{S;k}, \mathbf{v}_{S;k} - \mathbf{x}_{S;k}i + g(\mathbf{v}_{S;k})) \\ & \quad - \frac{K}{2}k\mathbf{v}_{1,1} - \mathbf{v}_{1,0}k^2 + K a_1 (f(\mathbf{x}_{1,1}) + h r f(\mathbf{x}_{1,1}), \mathbf{v}_{1,1} - \mathbf{x}_{1,1}i + g(\mathbf{v}_{1,1})) \\ & \quad + \sum_{s=2}^S \sum_{k=1}^K \frac{K(1+A_S\gamma)}{2}k\mathbf{v}_{S;k} - \mathbf{v}_{S;k-1}k^2 \\ & \quad + \sum_{s=2}^S \sum_{k=1}^K a_s (f(\mathbf{x}_{S;k}) + h\mathbf{q}_{S;k}, \mathbf{v}_{S;k} - \mathbf{x}_{S;k}i + g(\mathbf{v}_{S;k})), \end{aligned} \quad (25)$$

where the first equality is by $\psi_{S+1,0} = \psi_{S;K}$ and $\mathbf{v}_{S+1,0} = \mathbf{v}_{S;K}$, the second equality is by the definition of $\psi_{s;k}$, the last inequality is by the definition of $\psi_{1,1}(\mathbf{v}_{1,1})$ and the $K(1+A_S\gamma)$ -strong convexity of $\psi_{s;k-1}(s=2, k \geq [K])$. Then by Lemmas 4 and 1, we have

$$\begin{aligned} & f(\mathbf{v}_{1,1}) - f(\mathbf{x}_{1,1}) + h r f(\mathbf{x}_{1,1}), \mathbf{v}_{1,1} - \mathbf{x}_{1,1}i + \frac{L}{2}k\mathbf{v}_{1,1} - \mathbf{x}_{1,1}k^2 \\ & \leq f(\mathbf{x}_{1,1}) + h r f(\mathbf{x}_{1,1}), \mathbf{v}_{1,1} - \mathbf{x}_{1,1}i + \frac{1}{4a_1}k\mathbf{v}_{1,1} - \mathbf{v}_{1,0}k^2, \end{aligned} \quad (26)$$

where the last inequality is by $a_1 = \frac{1}{4L}$ and $\mathbf{v}_{1,0} = \mathbf{x}_{1,1}$.

Using $\mathbf{y}_s = \frac{1}{K} \sum_{k=1}^K \mathbf{y}_{S;k} = \frac{A_s-1}{A_s} \mathbf{y}_{s-1} + \frac{a_s}{KA_s} \sum_{k=1}^K \mathbf{v}_{S;k}$, the convexity of g , and $A_0 = 0$, we have

$$\begin{aligned} & \sum_{s=2}^S \sum_{k=1}^K a_s g(\mathbf{v}_{S;k}) - \sum_{s=2}^S K a_s g\left(\frac{1}{K} \sum_{k=1}^K \mathbf{v}_{S;k}\right) - \sum_{s=2}^S (K A_s g(\mathbf{y}_s) - K A_{s-1} g(\mathbf{y}_{s-1})) \\ & = K A_S g(\mathbf{y}_S) - K A_1 g(\mathbf{y}_1) \\ & = K A_S g(\mathbf{y}_S) - K A_1 g(\mathbf{v}_{1,1}). \end{aligned} \quad (27)$$

Thus, combining Eqs. (24)–(27), we have

$$\begin{aligned}
 KA_S f(\mathbf{u}) &= \frac{K(1+A_S\gamma)}{2} k\mathbf{u} \quad \mathbf{v}_{S;K} k^2 \quad \frac{K}{2} k\mathbf{x}_0 \quad \mathbf{u} k^2 + \frac{K}{4} k\mathbf{v}_{1;1} \quad \mathbf{v}_{1;0} k^2 + Ka_1 f(\mathbf{v}_{1;1}) \\
 &+ \sum_{s=2}^S \sum_{k=1}^K \left(a_s h r f(\mathbf{x}_{S;k}) \quad \mathbf{q}_{S;k}, \mathbf{u} \quad \mathbf{x}_{S;k} i + \frac{K(1+A_S \gamma)}{2} k\mathbf{v}_{S;k} \quad \mathbf{v}_{S;k-1} k^2 \right) \\
 &+ \sum_{s=2}^S \sum_{k=1}^K a_s (f(\mathbf{x}_{S;k}) + h\mathbf{q}_{S;k}, \mathbf{v}_{S;k} \quad \mathbf{x}_{S;k} i) + KA_S g(\mathbf{y}_S).
 \end{aligned} \tag{28}$$

Then with $A_0 = 0$ and $\mathbf{y}_s = \frac{1}{K} \sum_{k=1}^K \mathbf{y}_{S;k}$, we also have

$$\begin{aligned}
 KA_S f(\mathbf{y}_S) &= KA_1 f(\mathbf{y}_1) + K \sum_{s=2}^S (A_s f(\mathbf{y}_s) \quad A_{s-1} f(\mathbf{y}_{s-1})) \\
 &+ Ka_1 f(\mathbf{v}_{1;1}) + \sum_{s=2}^S \sum_{k=1}^K A_s f(\mathbf{y}_{S;k}) \quad K \sum_{s=2}^S A_{s-1} f(\mathbf{y}_{s-1}),
 \end{aligned} \tag{29}$$

where the last equality is by $A_1 = a_1$, $\mathbf{y}_1 = \mathbf{v}_{1;1}$. Subtracting Eq. (28) from (29) and noting that $f(\mathbf{y}_S) = f(\mathbf{y}_S) + g(\mathbf{y}_S)$ now completes the proof. \square

Lemma 7. *The error sequence $fE_{S;k}(\mathbf{u})g_{s=2;k \in [K]}$ in Lemma 5 satisfies*

$$\begin{aligned}
 E_{S;k}(\mathbf{u}) &= A_{s-1} (f(\mathbf{y}_{s-1}) \quad f(\mathbf{x}_{S;k}) \quad h r f(\mathbf{x}_{S;k}), \mathbf{y}_{s-1} \quad \mathbf{x}_{S;k} i) \\
 &+ a_s h r f(\mathbf{x}_{S;k}) \quad \mathbf{q}_{S;k}, \mathbf{v}_{S;k} \quad \mathbf{u} i + \left(\frac{La_s^2}{2A_s} \quad \frac{K(1+A_S \gamma)}{2} \right) k\mathbf{v}_{S;k} \quad \mathbf{v}_{S;k-1} k^2.
 \end{aligned}$$

Proof. Using Assumption 3, Lemma 4, and Lemma 1, and applying the definition of $\mathbf{y}_{S;k}$, we have

$$\begin{aligned}
 f(\mathbf{y}_{S;k}) &= f(\mathbf{x}_{S;k}) \\
 &+ h r f(\mathbf{x}_{S;k}), \mathbf{y}_{S;k} \quad \mathbf{x}_{S;k} i + \frac{L}{2} k\mathbf{y}_{S;k} \quad \mathbf{x}_{S;k} k^2 \\
 &= h r f(\mathbf{x}_{S;k}), \frac{A_{s-1}}{A_s} \mathbf{y}_{s-1} + \frac{a_s}{A_s} \mathbf{v}_{S;k} \quad \mathbf{x}_{S;k} i + \frac{La_s^2}{2A_s^2} k\mathbf{v}_{S;k} \quad \mathbf{v}_{S;k-1} k^2 \\
 &= \frac{A_{s-1}}{A_s} h r f(\mathbf{x}_{S;k}), \mathbf{y}_{s-1} \quad \mathbf{x}_{S;k} i + \frac{a_s}{A_s} h r f(\mathbf{x}_{S;k}), \mathbf{v}_{S;k} \quad \mathbf{x}_{S;k} i + \frac{La_s^2}{2A_s^2} k\mathbf{v}_{S;k} \quad \mathbf{v}_{S;k-1} k^2.
 \end{aligned} \tag{30}$$

It remains to plug Eq. (30) into the definition of $E_{S;k}(\mathbf{u})$, and rearrange. \square

The definition of the variance reduced extrapolation point $\mathbf{q}_{S;k}$ is crucial for bounding the error terms $fE_{S;k}(\mathbf{u})g$ from Lemma 5. The next three auxiliary lemmas apply the definition of $\mathbf{q}_{S;k}^{(j)}$ to bound the inner product term $h r f(\mathbf{x}_{S;k}) \quad \mathbf{q}_{S;k}, \mathbf{v}_{S;k} \quad \mathbf{u} i$ in $E_{S;k}(\mathbf{u})$ when we take the expectation over all randomness in the algorithm. We will use $\mathcal{F}_{S;k;i}$ to denote the natural filtration, containing all randomness up to and including epoch s , outer iteration k , and inner iteration i . Note that in Algorithm 4, the index of the inner iteration goes from $j = m$ to 1, therefore inner iteration i corresponds to when index of the inner iteration is $j = m - i + 1$. This detail however does not play a important role in our analysis.

Lemma 8. For all $s \geq 2$, $k \in [K]$ and $\mathbf{u} \in \text{dom}(g)$, we have

$$\begin{aligned}
 & a_s \mathbb{E}[hr f(\mathbf{x}_{s;k}) - \mathbf{q}_{s;k}, \mathbf{v}_{s;k} - \mathbf{u}] \\
 &= \sum_{j=1}^m a_s \mathbb{E}[hr^{(j)} f(\mathbf{x}_{s;k}) - r^{(j)} f(\mathbf{w}_{s;k;j}), \mathbf{v}_{s;k}^{(j)} - \mathbf{u}^{(j)}] + \sum_{j=1}^m a_{s;k-1} \mathbb{E}[hr^{(j)} f(\mathbf{x}_{s;k-1}) - r^{(j)} f(\mathbf{w}_{s;k-1;j}), \mathbf{v}_{s;k-1}^{(j)} - \mathbf{u}^{(j)}] \\
 &\quad + \sum_{j=1}^m a_{s;k-1} \mathbb{E}[hr^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) - r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j}), \mathbf{v}_{s;k-1}^{(j)} - \mathbf{v}_{s;k-1}^{(j)}] \\
 &\quad + \sum_{j=1}^m a_s \mathbb{E}[hr^{(j)} f(\mathbf{w}_{s;k;j}) - (r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) + \frac{\langle \mathbf{y}_{s-1}, \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \rangle}{s}), \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)}],
 \end{aligned}$$

Proof. Using the definition of $\mathbf{q}_{s;k}^{(j)}$, we have

$$\begin{aligned}
 & a_s (r^{(j)} f(\mathbf{x}_{s;k}) - \mathbf{q}_{s;k}^{(j)}) \\
 &= a_s (r^{(j)} f(\mathbf{x}_{s;k}) - r^{(j)} f(\mathbf{w}_{s;k;j})) + a_s (r^{(j)} f(\mathbf{w}_{s;k;j}) - \mathbf{q}_{s;k}^{(j)}) \\
 &= a_s (r^{(j)} f(\mathbf{x}_{s;k}) - r^{(j)} f(\mathbf{w}_{s;k;j})) + a_s (r^{(j)} f(\mathbf{w}_{s;k;j}) - (r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) + \frac{\langle \mathbf{y}_{s-1}, \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \rangle}{s})) \\
 &\quad + a_{s;k-1} (r^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) - r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j})).
 \end{aligned} \tag{31}$$

First, for $j \in [m]$ and any fixed $\mathbf{u}^{(j)}$, we have

$$\begin{aligned}
 & \mathbb{E}[a_s hr^{(j)} f(\mathbf{w}_{s;k;j}) - (r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) + \frac{\langle \mathbf{y}_{s-1}, \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \rangle}{s}), \mathbf{v}_{s;k}^{(j)} - \mathbf{u}^{(j)}] \\
 &= \mathbb{E}[a_s hr^{(j)} f(\mathbf{w}_{s;k;j}) - (r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) + \frac{\langle \mathbf{y}_{s-1}, \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \rangle}{s}), \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)}] \\
 &\quad + a_s \mathbb{E}[(r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) + \frac{\langle \mathbf{y}_{s-1}, \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \rangle}{s}) | F_{s;k;j-1}, \mathbf{v}_{s;k-1}^{(j)} - \mathbf{u}^{(j)}] \\
 &= \mathbb{E}[a_s hr^{(j)} f(\mathbf{w}_{s;k;j}) - (r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) + \frac{\langle \mathbf{y}_{s-1}, \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \rangle}{s}), \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)}],
 \end{aligned} \tag{32}$$

where the first equality follows from $\mathbf{v}_{s;k-1}^{(j)} \in F_{s;k;j-1}$ and the second equality follows from $\mathbb{E}[r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) | F_{s;k;j-1}] = r^{(j)} f(\mathbf{w}_{s;k;j})$ and $\mathbb{E}[r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) | F_{s;k;j-1}] = r^{(j)} f(\mathbf{y}_{s-1}) = \frac{\langle \mathbf{y}_{s-1}, \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \rangle}{s}$. Meanwhile, for $j \in [m]$ and any fixed $\mathbf{u}^{(j)}$, we have

$$\begin{aligned}
 & \mathbb{E}[a_{s;k-1} hr^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) - r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j}), \mathbf{v}_{s;k}^{(j)} - \mathbf{u}^{(j)}] \\
 &= \mathbb{E}[a_{s;k-1} hr^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) - r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j}), \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)}] \\
 &\quad + \mathbb{E}[\mathbb{E}[a_{s;k-1} hr^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) - r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j}), \mathbf{v}_{s;k-1}^{(j)} - \mathbf{u}^{(j)} | F_{s;k;j-1}]] \\
 &= \mathbb{E}[a_{s;k-1} hr^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) - r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j}), \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)}] \\
 &\quad + \mathbb{E}[a_{s;k-1} hr^{(j)} f(\mathbf{x}_{s;k-1}) - r^{(j)} f(\mathbf{w}_{s;k-1;j}), \mathbf{v}_{s;k-1}^{(j)} - \mathbf{u}^{(j)}],
 \end{aligned} \tag{33}$$

where the last equality is by $\mathbf{v}_{s;k-1}^{(j)} \in F_{s;k;j-1}$, $\mathbb{E}[r^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) | F_{s;k;j-1}] = r^{(j)} f(\mathbf{x}_{s;k-1})$ and $\mathbb{E}[r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j}) | F_{s;k;j-1}] = r^{(j)} f(\mathbf{w}_{s;k-1;j})$. Combining Eqs. (31)–(33) completes the proof. \square

In the following two lemmas, we will bound the third and the fourth terms of the R.H.S. in Lemma 8 by above using our novel Lipschitz Assumption 2 and Assumption 3.

Lemma 9. For $s \geq 2$ and $k \geq [K]$, we have

$$\sum_{j=1}^m a_{s;k} \mathbb{E} \left[\left\langle r^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) - r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j}), \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\rangle \right] \\ \mathbb{E} \left[\frac{K(1+A_s-1\gamma)}{8} k \|\mathbf{v}_{s;k} - \mathbf{v}_{s;k-1}\|^2 + \frac{a_{s;k-1}^4 L^2}{K A_{s;k-1}^2 (1+A_s-1\gamma)} k \|\mathbf{v}_{s;k-1}\|^2 \right],$$

where $a_{s;0} = a_{s-1}$, $A_{s;0} = A_{s-1}$ and $a_{s;k} = a_s$, $A_{s;k} = A_s$ for $k \geq [K]$.

Proof. Using Cauchy–Schwarz and Young’s inequalities, we have

$$a_{s;k} \mathbb{E} \left[\left\langle r^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) - r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j}), \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\rangle \right] \\ \mathbb{E} \left[\frac{2a_{s;k-1}^2}{K(1+A_s-1\gamma)} \left\| r^{(j)} f_{t_j}(\mathbf{x}_{s;k-1}) - r^{(j)} f_{t_j}(\mathbf{w}_{s;k-1;j}) \right\|^2 + \frac{K(1+A_s-1\gamma)}{8} \left\| \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\|^2 \right] \\ \mathbb{E} \left[\frac{2a_{s;k-1}^2}{K(1+A_s-1\gamma)} k \|\mathbf{x}_{s;k-1} - \mathbf{w}_{s;k-1;j}\|_{\mathcal{Q}}^2 + \frac{K(1+A_s-1\gamma)}{8} \left\| \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\|^2 \right] \\ = \mathbb{E} \left[\frac{2a_{s;k-1}^2}{K(1+A_s-1\gamma)} k \|\mathbf{x}_{s;k-1} - \mathbf{y}_{s;k-1}\|_{\mathcal{Q}}^2 + \frac{K(1+A_s-1\gamma)}{8} \left\| \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\|^2 \right], \quad (34)$$

where we used Assumption 3 in the first inequality and the definitions of $\mathbf{x}_{s;k-1}$ and $\mathbf{w}_{s;k-1;j}$ in the last equality. Finally by including the summation and using the definition of L , $\mathbf{x}_{s;k-1}$ and $\mathbf{y}_{s;k-1}$, the first term of the above expression becomes

$$\sum_{j=1}^m \mathbb{E} \left[\frac{2a_{s;k-1}^2}{K(1+A_s-1\gamma)} k \|\mathbf{x}_{s;k-1} - \mathbf{y}_{s;k-1}\|_{\mathcal{Q}}^2 \right] = \mathbb{E} \left[\frac{2a_{s;k-1}^2}{K(1+A_s-1\gamma)} k \|\mathbf{x}_{s;k-1} - \mathbf{y}_{s;k-1}\|_{\sum_{j=1}^m \mathcal{Q}}^2 \right] \\ \mathbb{E} \left[\frac{a_{s;k-1}^4 L^2}{K A_{s;k-1}^2 (1+A_s-1\gamma)} k \|\mathbf{v}_{s;k-1}\|^2 \right], \quad (35)$$

where $a_{s;0} = a_{s-1}$, $A_{s;0} = A_{s-1}$ and $a_{s;k} = a_s$, $A_{s;k} = A_s$ for $k \geq [K]$. Taking summation over j and combining Eqs. (34) and (35) give the lemma statement. \square

Lemma 10. For $s \geq 2$ and $k \geq [K]$, we have

$$\sum_{j=1}^m a_s \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{w}_{s;k;j}) - \left(r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) + r^{(j)} f(\mathbf{y}_{s-1}) \right), \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\rangle \right] \\ \mathbb{E} \left[\left(\frac{2L^2 a_s^4}{K A_s^2 (1+A_s-1\gamma)} + \frac{K(1+A_s-1\gamma)}{8} \right) k \|\mathbf{v}_{s;k} - \mathbf{v}_{s;k-1}\|^2 \right] \\ + \frac{8a_s^2 L}{K(1+A_s-1\gamma)} \mathbb{E} [f(\mathbf{y}_{s-1}) - f(\mathbf{x}_{s;k}) - h r f(\mathbf{x}_{s;k}), \mathbf{y}_{s-1} - \mathbf{x}_{s;k}]$$

Proof. Using similar arguments in the proof of Lemma 9, we have

$$a_s \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{w}_{s;k;j}) - \left(r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) + r^{(j)} f(\mathbf{y}_{s-1}) \right), \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\rangle \right] \\ \mathbb{E} \left[\frac{2a_s^2}{K(1+A_s-1\gamma)} \left\| r^{(j)} f(\mathbf{w}_{s;k;j}) - \left(r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) + r^{(j)} f(\mathbf{y}_{s-1}) \right) \right\|^2 \right] \\ + \frac{K(1+A_s-1\gamma)}{8} \left\| \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\|^2 \right]$$

$$\begin{aligned}
 &= \mathbb{E} \left[\frac{2a_s^2}{K(1+A_s-1\gamma)} \mathbb{E} \left[\left\| \left(r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) + r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) \right) - \left(r^{(j)} f(\mathbf{w}_{s;k;j}) - r^{(j)} f(\mathbf{y}_{s-1}) \right) \right\|^2 \middle| \mathcal{F}_{s;k;j-1} \right] \right. \\
 &\quad \left. + \mathbb{E} \left[\frac{K(1+A_s-1\gamma)}{8} \left\| \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\|^2 \right] \right] \\
 &= \mathbb{E} \left[\frac{2a_s^2}{K(1+A_s-1\gamma)} \mathbb{E} \left[\left\| r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) \right\|^2 \middle| \mathcal{F}_{s;k;j-1} \right] \right] + \mathbb{E} \left[\frac{K(1+A_s-1\gamma)}{8} \left\| \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\|^2 \right] \\
 &\mathbb{E} \left[\frac{4a_s^2}{K(1+A_s-1\gamma)} \left(\left\| r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{x}_{s;k}) \right\|^2 + \left\| r^{(j)} f_{t_j}(\mathbf{x}_{s;k}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) \right\|^2 \right) \right. \\
 &\quad \left. + \frac{K(1+A_s-1\gamma)}{8} \left\| \mathbf{v}_{s;k}^{(j)} - \mathbf{v}_{s;k-1}^{(j)} \right\|^2 \right], \tag{36}
 \end{aligned}$$

where the first equality comes from $\mathbb{E} [r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) | \mathcal{F}_{s;k;j-1}] = r^{(j)} f(\mathbf{w}_{s;k;j}) - r^{(j)} f(\mathbf{y}_{s-1})$ since the only randomness is in t_j when conditioned at $\mathcal{F}_{s;k;j-1}$, and the last inequality comes from $(a+b)^2 \leq 2(a^2+b^2)$. In order to bound the second term in Eq. (36), we will include the outer summation with respect to j and apply the results from Lemma 1 to get

$$\begin{aligned}
 \sum_{j=1}^m \mathbb{E} \left[\left\| r^{(j)} f_{t_j}(\mathbf{x}_{s;k}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) \right\|^2 \right] &= \sum_{j=1}^m \mathbb{E} \left[\mathbb{E} \left[\left\| r^{(j)} f_{t_j}(\mathbf{x}_{s;k}) - r^{(j)} f_{t_j}(\mathbf{y}_{s-1}) \right\|^2 \middle| \mathcal{F}_{s;k;0} \right] \right] \\
 &= \mathbb{E} \left[\sum_{j=1}^m \sum_{l=1}^n \frac{1}{n} \left\| r^{(j)} f_l(\mathbf{x}_{s;k}) - r^{(j)} f_l(\mathbf{y}_{s-1}) \right\|^2 \right] \\
 &= \mathbb{E} \left[\sum_{l=1}^n \frac{1}{n} k r f_l(\mathbf{x}_{s;k}) - r f_l(\mathbf{y}_{s-1}) k^2 \right] \\
 &\mathbb{E} [2L(f(\mathbf{y}_{s-1}) - f(\mathbf{x}_{s;k})) - hr f(\mathbf{x}_{s;k}), \mathbf{y}_{s-1} - \mathbf{x}_{s;k} l)], \tag{37}
 \end{aligned}$$

where the second equality comes from $\mathbf{x}_{s;k}, \mathbf{y}_{s-1} \in \mathcal{F}_{s;k;0}$ and the last inequality is by applying Lemma 1 and the definition of $f(\mathbf{x}) = \frac{1}{n} \sum_{l=1}^n f_l(\mathbf{x})$. To bound the first term of Eq. (36), we apply similar arguments as in Lemma 9 and get

$$\sum_{j=1}^m \mathbb{E} \left[\left\| r^{(j)} f_{t_j}(\mathbf{w}_{s;k;j}) - r^{(j)} f_{t_j}(\mathbf{x}_{s;k}) \right\|^2 \right] \leq \mathbb{E} \left[\frac{a_s^2 L^2}{2A_s^2} k \mathbf{v}_{s;k} - \mathbf{v}_{s;k-1} k^2 \right]. \tag{38}$$

Combining Eqs. (36) – (38) gives the lemma statement. \square

Lemma 6. With $a_s^2 = \frac{KA_{s-1}(1+A_{s-1})}{8L}$, $a_{s;k} = a_s$ and $A_{s;k} = A_s$ for $k \geq [K]$, $a_{s;0} = a_{s-1}$ and $A_{s;0} = A_{s-1}$, then for any fixed $\mathbf{u} \in \text{dom}(g)$ we have

$$\begin{aligned}
 &\sum_{s=2}^S \sum_{k=1}^K \mathbb{E} [E_{s;k}(\mathbf{u})] \\
 &\sum_{j=1}^m a_1 \left\langle r^{(j)} f(\mathbf{x}_{1;1}) - r^{(j)} f(\mathbf{w}_{1;1;j}), \mathbf{v}_{1;1}^{(j)} - \mathbf{u}^{(j)} \right\rangle + \sum_{j=1}^m a_S \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{x}_{S;K}) - r^{(j)} f(\mathbf{w}_{S;K;j}), \mathbf{v}_{S;K}^{(j)} - \mathbf{u}^{(j)} \right\rangle \right] \\
 &+ \frac{K}{64} k \mathbf{v}_{1;1} - \mathbf{v}_{1;0} k^2 - \frac{5K(1+A_S-1\gamma)}{32} \mathbb{E} \left[k \mathbf{v}_{S;K} - \mathbf{v}_{S;K-1} k^2 \right],
 \end{aligned}$$

where $\mathbf{x}_{1;1}, \mathbf{v}_{1;0} \in \text{dom}(g)$ can be chosen arbitrarily and $\mathbf{w}_{1;1;j}$ is defined in Algorithm 4.

Proof. Combining Lemma 7, 8, 9, 10, setting a_s such that $a_s^2 = \frac{KA_{s-1}(1+A_{s-1})}{8L}$ and using $\frac{A_{s-1}}{A_s} \leq 1$, we have

$$\begin{aligned} \mathbb{E}[E_{S;k}(\mathbf{u})] &= \sum_{j=1}^m a_s \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{x}_{S;k}) - r^{(j)} f(\mathbf{w}_{S;k;j}), \mathbf{v}_{S;k}^{(j)} - \mathbf{u}^{(j)} \right\rangle \right] \\ &= \sum_{j=1}^m a_{S;k-1} \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{x}_{S;k-1}) - r^{(j)} f(\mathbf{w}_{S;k-1;j}), \mathbf{v}_{S;k-1}^{(j)} - \mathbf{u}^{(j)} \right\rangle \right] \\ &= \left(\frac{5K(1+A_S\gamma)}{32} \right) \mathbb{E} \left[k\mathbf{v}_{S;k} - \mathbf{v}_{S;k-1} k^2 \right] + \left(\frac{a_{S;k-1}^4 L^2}{KA_{S;k-1}^2(1+A_S\gamma)} \right) \mathbb{E} \left[k\mathbf{v}_{S;k-1} - \mathbf{v}_{S;k-2} k^2 \right]. \end{aligned}$$

Next, by setting $a_{s,0} = a_{s-1}$ and $a_{s;k} = a_s$ for $k = [K]$, we can telescope the error terms and get

$$\begin{aligned} \sum_{s=2}^S \sum_{k=1}^K \mathbb{E}[E_{S;k}(\mathbf{u})] &= \sum_{j=1}^m \sum_{s=2}^S \sum_{k=1}^K a_s \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{x}_{S;k}) - r^{(j)} f(\mathbf{w}_{S;k;j}), \mathbf{v}_{S;k}^{(j)} - \mathbf{u}^{(j)} \right\rangle \right] \\ &= \sum_{j=1}^m \sum_{s=2}^S \sum_{k=1}^K a_{s;k-1} \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{x}_{S;k-1}) - r^{(j)} f(\mathbf{w}_{S;k-1;j}), \mathbf{v}_{S;k-1}^{(j)} - \mathbf{u}^{(j)} \right\rangle \right] \\ &= \sum_{s=2}^S \sum_{k=1}^K \frac{5K(1+A_S\gamma)}{32} \mathbb{E} \left[k\mathbf{v}_{S;k} - \mathbf{v}_{S;k-1} k^2 \right] \\ &\quad + \sum_{s=2}^S \left[\frac{K(1+A_S2\gamma)}{64} k\mathbf{v}_{S,0} - \mathbf{v}_{S,1} k^2 + \sum_{k=2}^K \frac{K(1+A_S\gamma)}{64} k\mathbf{v}_{S;k-1} - \mathbf{v}_{S;k-2} k^2 \right] \\ &= \sum_{j=1}^m a_S \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{x}_{S;K}) - r^{(j)} f(\mathbf{w}_{S;K;j}), \mathbf{v}_{S;K}^{(j)} - \mathbf{u}^{(j)} \right\rangle \right] \\ &\quad + \sum_{j=1}^m a_1 \mathbb{E} \left[\left\langle r^{(j)} f(\mathbf{x}_{1,1}) - r^{(j)} f(\mathbf{w}_{1,1;j}), \mathbf{v}_{1,1}^{(j)} - \mathbf{u}^{(j)} \right\rangle \right] \\ &\quad + \frac{K(1+A_0\gamma)}{64} \mathbb{E} \left[k\mathbf{v}_{2,0} - \mathbf{v}_{2,1} k^2 \right] - \frac{5K(1+A_S\gamma)}{32} \mathbb{E} \left[k\mathbf{v}_{S;K} - \mathbf{v}_{S;K-1} k^2 \right]. \end{aligned}$$

The lemma follows by setting $A_0 = 0$, $\mathbf{v}_{2,1} = \mathbf{v}_{1,0}$, $\mathbf{x}_{2,0} = \mathbf{x}_{1,1}$ and $\mathbf{w}_{2,0;j} = \mathbf{w}_{1,1;j}$. \square

Theorem 2. Let $\mathbf{x}_0 \in \text{dom}(g)$ be an arbitrary initial point. Fix $K \geq 1$ and consider the updates in Algorithm 4. Then for $S \geq 2$ and $8\mathbf{u} \in \text{dom}(g)$, we have

$$\mathbb{E} [f(\mathbf{y}_S) - f(\mathbf{u})] + \frac{9(1+A_S\gamma)}{64A_S} \mathbb{E} \left[k\mathbf{v}_{S;K} - \mathbf{u} k^2 \right] \leq \frac{5}{8A_S} k\mathbf{x}_0 - \mathbf{u} k^2.$$

In particular if $\mathbf{x} = \arg \min_{\mathbf{x}} f(\mathbf{x})$ exists, then we have

$$\mathbb{E} [f(\mathbf{y}_S) - f(\mathbf{x})] \leq \frac{5}{8A_S} k\mathbf{x}_0 - \mathbf{x} k^2$$

and

$$\mathbb{E} \left[k\mathbf{v}_{S;K} - \mathbf{x} k^2 \right] \leq \frac{40}{9(1+A_S\gamma)} k\mathbf{x}_0 - \mathbf{x} k^2.$$

Finally in all the bounds above we have

$$A_S \leq \max \left\{ \frac{S^2 K}{64L}, \frac{1}{4L} \left(1 + \sqrt{\frac{K\gamma}{8L}} \right)^{S-1} \right\}.$$

Proof. Combining Lemma 5 and Lemma 6, and by setting $\mathbf{y}_{1,0} = \mathbf{x}_{1,1} = \mathbf{x}_0$ and $\mathbf{y}_{1,1} = \mathbf{v}_{1,1}$, we have

$$\begin{aligned}
 KA_S E [f(\mathbf{y}_S) - f(\mathbf{u})] & \leq \frac{K}{2} k \mathbf{x}_0 \quad \mathbf{u} k^2 + \frac{K(1+A_S)}{2} E [k \mathbf{v}_{S,K} \quad \mathbf{u} k^2] \\
 & \quad + \frac{15K}{64} k \mathbf{v}_{1,1} \quad \mathbf{x}_0 k^2 + \frac{5K(1+A_S+1\gamma)}{32} E [k \mathbf{v}_{S,K} \quad \mathbf{v}_{S,K-1} k^2] \\
 & \quad + \sum_{j=1}^m a_1 \langle r^{(j)} f(\mathbf{x}_{1,1}) - r^{(j)} f(\mathbf{w}_{1,1;j}), \mathbf{v}_{1,1}^{(j)} \quad \mathbf{u}^{(j)} \rangle \\
 & \quad + \sum_{j=1}^m a_S E \left[\langle r^{(j)} f(\mathbf{x}_{S,K}) - r^{(j)} f(\mathbf{w}_{S,K;j}), \mathbf{v}_{S,K}^{(j)} \quad \mathbf{v}_{S,K-1}^{(j)} \rangle \right].
 \end{aligned} \tag{39}$$

Using the same approach as Lemma 9 and Lemma 10, we can upper bound the first inner product term by

$$\begin{aligned}
 \sum_{j=1}^m a_1 \langle r^{(j)} f(\mathbf{x}_{1,1}) - r^{(j)} f(\mathbf{w}_{1,1;j}), \mathbf{v}_{1,1}^{(j)} \quad \mathbf{u}^{(j)} \rangle & \leq \frac{1}{8K} k \mathbf{v}_{1,1} \quad \mathbf{x}_0 k^2 + \frac{K}{16} k \mathbf{v}_{1,1} \quad \mathbf{u} k^2 \\
 & \quad + \frac{15K}{64} k \mathbf{v}_{1,1} \quad \mathbf{x}_0 k^2 + \frac{K}{8} k \mathbf{x}_0 \quad \mathbf{u} k^2,
 \end{aligned} \tag{40}$$

where we used $(a+b)^2 \leq 2(a^2+b^2)$, $a_1 = \frac{1}{4L}$ and $K \leq 2$ in the last inequality. Similarly, we have

$$\begin{aligned}
 \sum_{j=1}^m a_S E \left[\langle r^{(j)} f(\mathbf{x}_{S,K}) - r^{(j)} f(\mathbf{w}_{S,K;j}), \mathbf{v}_{S,K}^{(j)} \quad \mathbf{v}_{S,K-1}^{(j)} \rangle \right] & \tag{41} \\
 & \leq \frac{K(1+A_S+1\gamma)}{8} E [k \mathbf{v}_{S,K} \quad \mathbf{v}_{S,K-1} k^2] + \frac{K(1+A_S+1\gamma)}{64} E [k \mathbf{v}_{S,K} \quad \mathbf{u} k^2],
 \end{aligned}$$

where we also used $a_S^2 = \frac{KA_S+1(1+A_S-1)}{8L}$ here. Combining Eqs. (39)–(41) gives us our main bounds in the theorem. Lastly, recall that $\bar{a}_S g_{S-1}$ is chosen so that $a_S^2 = \frac{KA_S+1(1+A_S-1)}{8L}$. When $\gamma = 0$, this leads to the standard $A_S = \frac{k^2 K}{64L}$ growth of accelerated algorithms by choosing $a_S = \frac{SK}{32L}$ for $k \geq 1$. When $\gamma > 0$, we have $\frac{\bar{a}_S}{A_S+1} > \sqrt{\frac{K}{8L}}$, and it remains to use that $A_k = \frac{A_k}{A_{k-1}} \dots \frac{A_2}{A_1} A_1 = A_1 \left(1 + \sqrt{\frac{K}{8L}}\right)^{k-1}$ where $a_1 = A_1 = \frac{1}{4L}$ using the choice of a_k in Algorithm 1 and $A_0 = a_0 = 0$, completing the proof. \square

B.1. Adaptive Variance Reduced A-CODER

Similar to A-CODER, VR-A-CODER can adaptively estimate the Lipschitz parameter by checking the quadratic bounds between \mathbf{y}_S and $\mathbf{x}_{S,k}$ as well as between \mathbf{y}_S and $\mathbf{x}_{S,k}$. For completeness, we have included the adaptive version of VR-A-CODER in Algorithm 5 below.

Algorithm 5 Variance Reduced A-CODER (Adaptive Version)

```

1: Input:  $\mathbf{x}_0 \in \text{dom}(g), \gamma = 0, L_0 > 0, m, f \in S^1, \dots, S^m g$ 
2: Initialization:  $\mathbf{y}_0 = \mathbf{v}_{1,0} = \mathbf{y}_{1,0} = \mathbf{x}_{1,1} = \mathbf{x}_0; \mathbf{z}_{1,0} = \mathbf{0}$ 
3:  $L_1 = L_0/2$ 
4: repeat
5:    $L_1 = 2L_1$ 
6:    $a_0 = A_0 = 0; A_1 = a_1 = \frac{1}{4L_0}$ 
7:    $\mathbf{z}_{1,1} = r f(\mathbf{x}_0); \mathbf{v}_{1,1} = \text{prox}_{a_1 g}(\mathbf{x}_0 - \mathbf{z}_{1,1})$ 
8:   until  $f(\mathbf{v}_{1,1}) \leq f(\mathbf{x}_0) + hr f(\mathbf{x}_0), \mathbf{v}_{1,1} = \mathbf{x}_0 + \frac{L_1}{2} k \mathbf{v}_{1,1} - \mathbf{x}_0 k^2$ 
9:    $\mathbf{y}_1 = \mathbf{y}_{1,1} = \mathbf{v}_{1,1}$ 
10:   $\mathbf{w}_{1,1;j} = (\mathbf{x}_{1,1}^{(1)}, \dots, \mathbf{x}_{1,1}^{(j)}, \mathbf{y}_{1,1}^{(j+1)}, \dots, \mathbf{y}_{1,1}^{(m)})$ 
11:   $\mathbf{v}_{2,0} = \mathbf{v}_{1,1}; \mathbf{w}_{2,0;j} = \mathbf{w}_{1,1;j}; \mathbf{x}_{2,0} = \mathbf{x}_{1,1}; \mathbf{y}_{2,0} = \mathbf{y}_{1,1}; \mathbf{z}_{2,0} = \mathbf{z}_{1,1}$ 
12:  for  $s = 2$  to  $S$  do
13:     $L_s = L_{s-1}/2$ 
14:    repeat
15:       $L_s = 2L_s$ 
16:      Set  $a_s > 0$  s.t.  $a_s^2 = \frac{KA_{s-1}(1+A_{s-1})}{8L_s}; A_s = A_{s-1} + a_s$ 
17:       $a_{s,0} = a_{s-1}; a_{s,1} = a_{s,2} = \dots = a_{s,K} = a_s$ 
18:       $\mathbf{v}_{s,0} = \mathbf{v}_{s-1,K}; \mathbf{w}_{s,0;j} = \mathbf{w}_{s-1,K;j}; \mathbf{x}_{s,0} = \mathbf{x}_{s-1,K}; \mathbf{y}_{s,0} = \mathbf{y}_{s-1,K}; \mathbf{z}_{s,0} = \mathbf{z}_{s-1,K}$ 
19:       $\mathbf{y}_s = r f(\mathbf{y}_{s-1})$ 
20:      for  $k = 1$  to  $K$  do
21:         $\mathbf{x}_{s;k} = \frac{A_{s-1}}{A_s} \mathbf{y}_{s-1} + \frac{a_s}{A_s} \mathbf{v}_{s;k-1}$ 
22:        for  $j = m$  to  $1$  do
23:           $\mathbf{w}_{s;k;j} = (\mathbf{x}_{s;k}^{(1)}, \dots, \mathbf{x}_{s;k}^{(j)}, \mathbf{y}_{s;k}^{(j+1)}, \dots, \mathbf{y}_{s;k}^{(m)})$ 
24:          Choose  $t$  in  $[n]$  uniformly at random
25:           $r_{s;k}^{(j)} = r^{(j)} f_t(\mathbf{w}_{s;k;j}) - r^{(j)} f_t(\mathbf{y}_{s-1}) + \frac{a_s}{A_s} r_{s;k-1}^{(j)}$ 
26:           $\mathbf{q}_{s;k}^{(j)} = r_{s;k}^{(j)} + \frac{a_{s,k-1}}{a_s} (r^{(j)} f_t(\mathbf{x}_{s;k-1}) - r^{(j)} f_t(\mathbf{w}_{s;k-1;j}))$ 
27:           $\mathbf{z}_{s;k}^{(j)} = \mathbf{z}_{s;k-1}^{(j)} + a_s \mathbf{q}_{s;k}^{(j)}$ 
28:           $\mathbf{v}_{s;k}^{(j)} = \text{prox}_{(A_{s-1} + \frac{a_s k}{K})g^j}(\mathbf{x}_0 - \mathbf{z}_{s;k}^{(j)}/K)$ 
29:           $\mathbf{y}_{s;k}^{(j)} = \frac{A_{s-1}}{A_s} \mathbf{y}_{s-1} + \frac{a_s}{A_s} \mathbf{v}_{s;k}^{(j)}$ 
30:        end for
31:      end for
32:       $\mathbf{y}_s = \frac{1}{K} \sum_{k=1}^K \mathbf{y}_{s;k}$ 
33:      until  $f(\mathbf{y}_{s;k}) \leq f(\mathbf{x}_{s;k}) + hr f(\mathbf{x}_{s;k}), \mathbf{y}_{s;k} = \mathbf{x}_{s;k} + \frac{L_s}{2} k \mathbf{y}_{s;k} - \mathbf{x}_{s;k} k^2$ 
      and  $\frac{1}{n} \sum_{t=1}^n k r f_t(\mathbf{x}_{s;k}) - r f_t(\mathbf{y}_{s-1}) k^2 \leq 2L_s (f(\mathbf{y}_{s-1}) - f(\mathbf{x}_{s;k}) + hr f(\mathbf{x}_{s;k}), \mathbf{y}_{s-1} - \mathbf{x}_{s;k})$ 
34:    end for
35:  return  $\mathbf{v}_{S,K}, \mathbf{y}_S$ 

```
