



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

A stochastic proximal alternating method for non-smooth non-convex optimization

Citation for published version:

Driggs, D, Tang, J, Liang, J, Davies, M & Schönlieb, C-B 2021, 'A stochastic proximal alternating method for non-smooth non-convex optimization', *Siam journal on imaging sciences*, vol. 14, no. 4, pp. 1932–1970.
<https://doi.org/10.1137/20M1387213>

Digital Object Identifier (DOI):

[10.1137/20M1387213](https://doi.org/10.1137/20M1387213)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Siam journal on imaging sciences

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



A Stochastic Proximal Alternating Minimization Algorithm for Non-smooth and Non-convex Optimization*

Derek Driggs^{†‡}, Junqi Tang^{†§}, Jingwei Liang[¶], Mike Davies[§], and Carola-Bibiane Schönlieb[‡]

Abstract. In this work, we introduce a novel stochastic *proximal alternating linearized minimization* (PALM) algorithm [6] for solving a class of non-smooth and non-convex optimization problems. Large-scale imaging problems are becoming increasingly prevalent due to the advances in data acquisition and computational capabilities. Motivated by the success of stochastic optimization methods, we propose a stochastic variant of proximal alternating linearized minimization. We provide global convergence guarantees, demonstrating that our proposed method with variance-reduced stochastic gradient estimators, such as SAGA [16] and SARAH [27], achieves state-of-the-art oracle complexities. We also demonstrate the efficacy of our algorithm via several numerical examples including sparse non-negative matrix factorization, sparse principal component analysis and blind image deconvolution.

Key words. Non-convex and non-smooth optimization, Stochastic optimization, Variance reduction, Alternating minimization, Stochastic PALM, Kurdyka-Łojasiewicz inequality, Sparse principle component analysis

AMS subject classifications. 90C26, 90C15, 90C30, 49M27

1. Introduction. With the advent of large-scale machine learning, developing efficient and reliable algorithms for (empirical) risk minimization has become an intense focus of the optimization community. These tasks involve minimizing a loss function measuring the fit between observed data, x , and a model's predicted result, b : $\min_{x \in \mathbb{R}^{m_1}} \frac{1}{n} \sum_{i=1}^n F(x_i, b_i)$ where n denotes the number of samples and F is the loss function. The two defining qualities of these problems are their large scale (in many applications, n is on the order of billions), and finite-sum structure.

When the value of n is very large, computing the gradient of the loss function is often prohibitively expensive, rendering most traditional deterministic first-order optimization algorithms ineffective. Over the years, randomized optimization algorithms [7, 32] have become increasingly popular due to their efficiency and simplicity. For these algorithms, the full gradient is replaced by a stochastic approximation that is cheap to compute, so that their per-iteration complexity grows slowly with n . For objectives with a finite-sum structure, many works have shown that certain randomized algorithms achieve convergence rates similar to those of full-gradient methods, even though their per-iteration complexity is often a factor of n smaller [16, 21, 38].

Outside machine learning, objectives with a finite-sum structure also arise in problems from image processing and computer vision. Recently, randomized optimization algorithms have been explored for image processing tasks including PET reconstruction, deblurring and tomography [12, 36]. As stochastic methods expand into new applications, they move further from smooth, strongly convex finite-sum objectives where they are well-understood theoretically. In this work, we aim to provide a

[†]Contributed equally

*Submitted to the editors DATE.

[‡]Department of Applied Mathematics and Theoretical Physics, University of Cambridge, Cambridge, UK (d.driggs@maths.cam.ac.uk, cbs31@cam.ac.uk).

[§]School of Engineering, University of Edinburgh, Edinburgh, UK (j.tang@ed.ac.uk, mike.davies@ed.ac.uk).

[¶]Institute of Natural Sciences and School of Mathematical Sciences, Shanghai Jiao Tong University, Shanghai, China (jingwei.liang@sjtu.edu.cn).

better understanding of stochastic algorithms for problems that are neither smooth nor convex.

1.1. Non-smooth, non-convex optimization. Our goal is to minimize composite objectives of the following form:

$$(1.1) \quad \min_{x \in \mathbb{R}^{m_1}, y \in \mathbb{R}^{m_2}} \{ \Phi(x, y) \stackrel{\text{def}}{=} J(x) + F(x, y) + R(y) \},$$

where $F(x, y) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n F_i(x, y)$ has a finite-sum structure. In general, functions J and R are non-smooth regularizations that promote structures in the solutions, *e.g.* sparsity or non-negativity. The blocks x and y represent differently structured elements of the solution that are coupled through the loss term, $F(x, y)$. Throughout this work, we impose the following assumptions:

(A.1) $J : \mathbb{R}^{m_1} \rightarrow \mathbb{R} \cup \{+\infty\}$ and $R : \mathbb{R}^{m_2} \rightarrow \mathbb{R} \cup \{+\infty\}$ are proper lower semi-continuous (lsc) functions that are bounded from below;

(A.2) $F_i : \mathbb{R}^{m_1} \times \mathbb{R}^{m_2} \rightarrow \mathbb{R}$ are finite-valued, differentiable, and their gradients ∇F_i are $M(\mathcal{X}, \mathcal{Y})$ -Lipschitz continuous on bounded sets $\mathcal{X} \times \mathcal{Y} \subset \mathbb{R}^{m_1} \times \mathbb{R}^{m_2}$ for all $i \in \{1, \dots, n\}$;¹

(A.3) The partial gradients $\nabla_x F_i$ are Lipschitz continuous with modulus $L_1(y)$, and $\nabla_y F_i$ are Lipschitz continuous with modulus $L_2(x)$ for all $i \in \{1, \dots, n\}$;

(A.4) The function Φ is bounded from below.

No convexity is imposed on any of the functions involved. Problem (1.1) departs from the sum-of-convex-objectives models that populate the majority of the optimization literature. Many models in machine learning, statistics and image processing require the full generality of (1.1). Archetypal examples include non-negative or sparse matrix factorization [20], Sparse PCA [13, 42], Robust PCA [11], trimmed least-squares [1] and blind image deconvolution [10]. Despite the prevalence of these problems, few numerical methods can solve the general problem (1.1), and none that realize match the efficiency that randomized algorithms provide. We outline some existing options below.

Proximal alternating minimization. One approach to solve (1.1) is the Proximal Alternating Minimization (PAM) method [3], whose iterations take the following form:

$$(1.2) \quad \begin{aligned} x_{k+1} &\in \operatorname{Argmin}_{x \in \mathbb{R}^{m_1}} \left\{ \Phi(x, y_k) + \frac{1}{2\gamma_{x,k}} \|x - x_k\|^2 \right\}, \\ y_{k+1} &\in \operatorname{Argmin}_{y \in \mathbb{R}^{m_2}} \left\{ \Phi(x_{k+1}, y) + \frac{1}{2\gamma_{y,k}} \|y - y_k\|^2 \right\}, \end{aligned}$$

where $\gamma_{x,k}, \gamma_{y,k} > 0$ are step-sizes. A significant limitation of PAM is that the subproblems in (1.2) do not have closed-form solutions in general. As a consequence, each subproblem requires its own set of inner iterations, which makes PAM inefficient in practice.

Proximal alternating linearized minimization [6]. To circumvent this limitation of PAM, Proximal Alternating Linearized Minimization (PALM) [6] replaces PAM's two subproblems with their proximal linearizations. PALM's iterations take the form

$$(1.3) \quad \begin{aligned} x_{k+1} &\in \operatorname{prox}_{\gamma_{x,k} J} \left(x_k - \gamma_{x,k} \nabla_x F(x_k, y_k) \right), \\ y_{k+1} &\in \operatorname{prox}_{\gamma_{y,k} R} \left(y_k - \gamma_{y,k} \nabla_y F(x_{k+1}, y_k) \right), \end{aligned}$$

¹Because we consider a particular bounded set in our analysis, we drop the dependence on \mathcal{X} and \mathcal{Y} for the remainder of the paper, writing the Lipschitz constant as M .

Algorithm 1.1 SPRING: Stochastic Proximal Alternating Linearized Minimization

Initialize: $x_0 \in \mathbb{R}^{m_1}, y_0 \in \mathbb{R}^{m_2}$.
for $k = 1, 2, \dots, T - 1$ **do**
 $x_{k+1} \in \text{prox}_{\gamma_{x,k}J}(x_k - \gamma_{x,k}\tilde{\nabla}_x(x_k, y_k))$
 $y_{k+1} \in \text{prox}_{\gamma_{y,k}R}(y_k - \gamma_{y,k}\tilde{\nabla}_y(x_{k+1}, y_k))$
end for
return (x_T, y_T)

where $\nabla_x F$ and $\nabla_y F$ are partial derivatives, and $\text{prox}_{\gamma_{x,k}J}$ is called “proximal operator” of J and defined by

$$\text{prox}_{\gamma J}(\cdot) \stackrel{\text{def}}{=} \text{Argmin}_x \gamma J(x) + \frac{1}{2} \|x - \cdot\|^2.$$

67 The proximal mapping is set-valued in general, and becomes single-valued if J is convex.

68 In contrast to PAM, each subproblem of PALM can be efficiently computed if the proximal maps of
69 J and R are easy to calculate, which is true in many applications. PALM also has the same convergence
70 guarantees as PAM, so linearizing F in each proximal step is a clear improvement over PAM. PALM
71 with momentum is considered in [29], where the authors show that inertia allows PALM to converge to
72 critical points with lower objective values, although accelerated rates might not be obtained.

73 **1.2. Stochastic PALM.** In this work, we introduce SPRING, a randomized version of PALM
74 where the partial gradients $\nabla_x F(x_k, y_k)$ and $\nabla_y F(x_{k+1}, y_k)$ in (1.3) are replaced by random estimates,
75 $\tilde{\nabla}_x(x_k, y_k)$ and $\tilde{\nabla}_y(x_{k+1}, y_k)$, formed using the gradients of only a few indices $\nabla_x F_j(x_k, y_k)$ and
76 $\nabla_y F_j(x_{k+1}, y_k)$ for $j \in B_k \subset \{1, 2, \dots, n\}$. The mini-batch B_k is chosen uniformly at random from
77 all subsets of $\{1, 2, \dots, n\}$ with cardinality b . We describe SPRING in Algorithm 1.1.

78 Many different gradient estimators can be used in SPRING. The simplest one is the stochastic
79 gradient descent (SGD) estimator [33]

$$80 \quad \tilde{\nabla}_x^{\text{SGD}}(x_k, y_k) = \frac{1}{b} \sum_{j \in B_k} \nabla_x F_j(x_k, y_k),$$

81 which uses the gradient of a randomly sampled batch to represent the full gradient. Another popular
82 choice is SAGA gradient estimator [16], which incorporates the gradient history:

$$83 \quad \begin{aligned} \tilde{\nabla}_x^{\text{SAGA}}(x_k, y_k) &= \frac{1}{b} \sum_{j \in B_k} (\nabla_x F_j(x_k, y_k) - g_{k,j}) + \frac{1}{n} \sum_{i=1}^n g_{k,i}, \\ g_{k+1,i} &= \begin{cases} \nabla_x F_i(x_k, y_k) & \text{if } i \in B_k, \\ g_{k,i} & \text{o.w.} \end{cases} \end{aligned}$$

84 Both SGD and SAGA estimators are *unbiased*. The last gradient estimator we specifically consider in
85 this work is the (loopless) SARAH estimator [24, 27], $\tilde{\nabla}_x^{\text{SARAH}}(x_k, y_k)$, which is *biased*.

$$86 \quad \begin{cases} \nabla_x F(x_k, y_k) & \text{w.p. } \frac{1}{p} \\ \frac{1}{b} \sum_{j \in B_k} (\nabla_x F_j(x_k, y_k) - \nabla_x F_j(x_{k-1}, y_{k-1})) + \tilde{\nabla}_x^{\text{SARAH}}(x_{k-1}, y_{k-1}) & \text{o.w.} \end{cases}$$

87 Here, p is a tuning parameter that is generally set to $\mathcal{O}(n)$. Other variance-reduced estimators can be
88 used in SPRING, including the SAG [34] and SVRG [21] estimators, for example, but we consider only
89 the SAGA and SARAH estimators specifically in this work.

Computing the full gradient is generally n -times more expensive than computing $\nabla_x F_i$, so when n is large and $b \ll n$, each step of SPRING with any of these estimators is significantly less expensive than that of PALM.

Remark 1.1. Although we consider only two variable blocks in (1.1), the results of this paper easily extend to an arbitrary number of blocks to solve problems of the form

$$\min_{x_1, \dots, x_\ell} \left\{ \frac{1}{n} \sum_{i=1}^n F_i(x_1, \dots, x_\ell) + \sum_{t=1}^\ell R_t(x_t) \right\},$$

where each R_t is a (possibly non-smooth) regularizer.

1.3. Contributions. By combining PALM with popular stochastic gradient estimators which are variance reduced, we proposed a novel stochastic algorithm for non-convex and non-smooth optimization. Theoretically, we show that the resulted algorithm matches the convergence rates of PALM given that the gradient estimators $\tilde{\nabla}_x$ and $\tilde{\nabla}_y$ satisfy a *variance-reduced* property (see Definition 2.1). We prove convergence guarantees of two types.

Convergence rate of generalized gradient map. Given a point $z = (x, y)$, the *generalized gradient map* of PALM/SPRING is defined as

$$(1.4) \quad \mathcal{G}_{\gamma_1, \gamma_2}(z) \stackrel{\text{def}}{=} \begin{pmatrix} 1/\gamma_1 (x - \text{prox}_{\gamma_1 J}(x - \gamma_1 \nabla_x F(x, y))) \\ 1/\gamma_2 (y - \text{prox}_{\gamma_2 R}(y - \gamma_2 \nabla_y F(x, y))) \end{pmatrix},$$

where $\gamma_1, \gamma_2 > 0$ are parameters (not necessarily equal to the step-sizes in Algorithm 1.1). If $\text{dist}(0, \mathcal{G}_{\gamma_1, \gamma_2}(z)) = 0$, then by the definition of the proximal operator, $0 \in (\nabla_x F(x, y) + \partial J(x), \nabla_y F(x, y) + \partial R(y)) = \partial \Phi(z)$, meaning z is a critical point. The point z is an ϵ -approximate critical point if it satisfies $\text{dist}(0, \mathcal{G}_{\gamma_1, \gamma_2}(z)) \leq \epsilon$ for some $\gamma_1, \gamma_2 \in (0, \infty)$.² In Section 3, we show that³

$$\mathbb{E}[\text{dist}(0, \mathcal{G}_{\frac{\gamma_{x,\alpha}}{2}, \frac{\gamma_{y,\alpha}}{2}}(z_\alpha))^2] \leq \mathcal{O}\left(\frac{1}{k}\right),$$

where α is chosen uniformly at random from the set $\{1, 2, \dots, k\}$. If Φ satisfies a certain error bound involving the generalized gradient map (see Eq. (3.1)), then SPRING converges linearly to the global optimum. These results generalize almost all existing results for stochastic gradient methods on non-convex, non-smooth objectives [1, 18, 30, 37, 41].

Specializing these convergence guarantees to specific gradient estimators, the constants appearing in these rates scale with the mean-squared error (MSE, see Definition 2.1) of the gradient estimators.

- For the SAGA estimator with $b \leq \mathcal{O}(n^{2/3})$, the iterates of SPRING satisfy

$$\mathbb{E}[\text{dist}(0, \mathcal{G}_{\frac{\gamma_{x,\alpha}}{2}, \frac{\gamma_{y,\alpha}}{2}}(z_\alpha))^2] \leq \mathcal{O}\left(\frac{n^2 L}{b^3 k}\right).$$

- For the SARAH gradient estimator with any batch size, we have

$$\mathbb{E}[\text{dist}(0, \mathcal{G}_{\frac{\gamma_{x,\alpha}}{2}, \frac{\gamma_{y,\alpha}}{2}}(z_\alpha))^2] \leq \mathcal{O}\left(\frac{\sqrt{n} L}{k}\right).$$

²The set of ϵ -critical points depends on the parameters γ_1, γ_2 , with larger parameter values generally increasing the size of the set for fixed ϵ . For fixed and bounded γ_1 and γ_2 , the generalized gradient map provides a notion of distance to a critical point. If $\mathcal{S}(\epsilon)$ is the set of ϵ -critical points, then with γ_1 and γ_2 fixed and bounded, we have $\mathcal{S}(\epsilon_1) \subset \mathcal{S}(\epsilon_2)$ for $\epsilon_1 \leq \epsilon_2$, and as $\epsilon \rightarrow 0$, $\mathcal{S}(\epsilon)$ contains only the set of critical points of Φ .

³We prove bounds on the expectation of the *squared* norm of the generalized gradient map to facilitate comparisons with existing results [1, 30, 31].

These convergence rates imply complexity bounds with respect to a *stochastic first-order oracle* (SFO) which returns the partial gradient of a single component F_i (for example, $\nabla_x F_i(x_k, y_k)$). To find an ϵ -approximate critical point, SAGA with a mini-batch of size $n^{2/3}$ requires no more than $\mathcal{O}(n^{2/3}L/\epsilon^2)$ SFO calls, and SARAH requires no more than $\mathcal{O}(\sqrt{n}L/\epsilon^2)$. The improved dependence on n when using SARAH gradient estimator exists in all of our convergence rates for SPRING. Because most existing works on stochastic optimization for non-smooth, non-convex problems use models that are special cases of (1.1), our results for SPRING capture most existing work as special cases. In particular, in the case $R \equiv J \equiv 0$, our results recover recent results showing that SARAH achieves the *oracle complexity lower-bound* for non-convex problems with a finite-sum structure [18, 28, 37, 40, 41].

Convergence under the Kurdyka–Łojasiewicz property. We also provide convergence guarantees under the Kurdyka–Łojasiewicz property (see Definition 2.4). First, we prove the global convergence of the generated sequence under the assumption that the objective function $\Phi(x, y)$ of (1.1) has the Kurdyka–Łojasiewicz property. Then, under the assumption that Φ is semi-algebraic with KL-exponent θ (see Section 2), we show that the sequence $z_k = (x_k, y_k)$ generated by SPRING converges in expectation to a critical point z^* of problem (1.1) at the following rates:

- If $\theta = 0$, then $\{\mathbb{E}\Phi(z_k)\}_{k \in \mathbb{N}}$ converges to $\mathbb{E}\Phi(z^*)$ in a finite number of steps.
- If $\theta \in (0, 1/2]$, then $\mathbb{E}\|z_k - z^*\| \leq \mathcal{O}(\tau^k)$ for some $\tau \in (0, 1)$.
- If $\theta \in (1/2, 1)$, then $\mathbb{E}\|z_k - z^*\| \leq \mathcal{O}(k^{-\frac{1-\theta}{2\theta-1}})$.

These rates match the rates of the original PALM algorithm.

1.4. Prior Art. SPRING offers several advantages over existing stochastic algorithms for non-smooth non-convex optimization. Reddi *et al.* investigate proximal SAGA and SVRG for solving problems of the form (1.1) when y is constant and J is convex [30]. Using mini-batches of size $b = n^{2/3}$, SAGA and SVRG require $\mathcal{O}(n^{2/3}L/\epsilon^2)$ stochastic gradient evaluations to converge to an ϵ -approximate critical point. Similarly, Aravkin and Davis introduce TSVRG, a stochastic algorithm based on SVRG gradient estimator, for solving another special case of (1.1) [1]. Our work generalizes their results and improves them in many cases. Most importantly, we show that using SARAH gradient estimator allows SPRING to achieve a complexity of $\mathcal{O}(\sqrt{n}L/\epsilon^2)$ even when the mini-batch size is equal to one. Our results for semi-algebraic objectives offer even sharper convergence rates.

The block stochastic gradient method [39] is closely related to SPRING using the (non-variance-reduced) SGD gradient estimator. In a similar work, Davis *et al.* introduce SAPALM, an asynchronous version of PALM that allows stochastic noise in the gradients [15]. The authors prove convergence rates that scale with the variance of the noise in the gradients, with their best complexity bound for finding an ϵ -approximate critical point equal to $\mathcal{O}(nL/\epsilon^2)$. While significant in their own right, these results are not directly related to ours, as these works require an explicit bound on the variance of the noise in the gradients, and the gradient estimators we consider do not admit such a bound [15].

2. Preliminaries. We use the following definitions and notation throughout the manuscript.

Variance Reduction. In our analysis, we mainly focus on stochastic gradient estimators that are variance reduced. We use a general definition of a variance-reduced gradient estimator that includes all existing estimators, for example, SAGA and SARAH, as special cases.

Definition 2.1 (Variance-reduced gradient estimator). Let $\{z_k\}_{k \in \mathbb{N}} = \{(x_k, y_k)\}_{k \in \mathbb{N}}$ be the sequence generated by Algorithm 1.1 with some gradient estimator ∇ . This gradient estimator is variance-reduced with constants $V_1, V_2, V_\Upsilon \geq 0$, and $\rho \in (0, 1]$ if it satisfies the following conditions:

1. **(MSE Bound)** *There exists a sequence of random variables $\{\Upsilon_k\}_{k \geq 1}$ of the form $\Upsilon_k = \sum_{i=1}^s (v_k^i)^2$ for some non-negative random variables $v_k^i \in \mathbb{R}$ such that*

$$(2.1) \quad \mathbb{E}_k[\|\tilde{\nabla}_x(x_k, y_k) - \nabla_x F(x_k, y_k)\|^2 + \|\tilde{\nabla}_y(x_{k+1}, y_k) - \nabla_y F(x_{k+1}, y_k)\|^2] \\ \leq \Upsilon_k + V_1(\mathbb{E}_k\|z_{k+1} - z_k\|^2 + \|z_k - z_{k-1}\|^2),$$

and, with $\Gamma_k = \sum_{i=1}^s v_k^i$,

$$\mathbb{E}_k[\|\tilde{\nabla}_x(x_k, y_k) - \nabla_x F(x_k, y_k)\| + \|\tilde{\nabla}_y(x_{k+1}, y_k) - \nabla_y F(x_{k+1}, y_k)\|] \\ \leq \Gamma_k + V_2(\mathbb{E}_k\|z_{k+1} - z_k\| + \|z_k - z_{k-1}\|).$$

2. **(Geometric Decay)** *The sequence $\{\Upsilon_k\}_{k \geq 1}$ decays geometrically:*

$$(2.2) \quad \mathbb{E}_k \Upsilon_{k+1} \leq (1 - \rho) \Upsilon_k + V_\Upsilon(\mathbb{E}_k\|z_{k+1} - z_k\|^2 + \|z_k - z_{k-1}\|^2).$$

3. **(Convergence of Estimator)** *If $\{z_k\}_{k \in \mathbb{N}}$ satisfies $\lim_{k \rightarrow \infty} \mathbb{E}\|z_k - z_{k-1}\|^2 = 0$, then $\mathbb{E} \Upsilon_k \rightarrow 0$ and $\mathbb{E} \Gamma_k \rightarrow 0$.*

Proposition 2.2. *SAGA gradient estimator is variance-reduced with parameters $V_1 = 6M^2/b$, $V_2 = \sqrt{6}M/\sqrt{b}$, $V_\Upsilon = \frac{134nL^2}{b^2}$, and $\rho = \frac{b}{2n}$. SARAH estimator is variance-reduced with parameters $V_1 = V_\Upsilon = 2L^2$, $V_2 = 2L$, and $\rho = 1/p$.*

Proposition 2.2 is a generalization of existing variance bounds for these estimators. For a derivation of the constants appearing in Proposition 2.2, we refer the reader to our full work [17].

Remark 2.3. Our results allow Algorithm 1.1 to use any variance-reduced gradient estimator, even different estimators for ∇_x and ∇_y . In particular, it is possible to use different mini-batch sizes when approximating the two partial gradients.

Kurdyka–Łojasiewicz property. Let $H : \mathbb{R}^{m_1} \rightarrow \mathbb{R} \cup \{+\infty\}$ be a proper lower semicontinuous function. For ϵ_1, ϵ_2 satisfying $-\infty < \epsilon_1 < \epsilon_2 < +\infty$, define the set $[\epsilon_1 < H < \epsilon_2] \stackrel{\text{def}}{=} \{x \in \mathbb{R}^{m_1} : \epsilon_1 < H(x) < \epsilon_2\}$.

Definition 2.4 (Kurdyka–Łojasiewicz). A function H is said to have the Kurdyka–Łojasiewicz property at $\bar{x} \in \text{dom}(H)$ if there exists $\epsilon \in (0, +\infty]$, a neighborhood U of \bar{x} and a continuous concave function $\varphi : [0, \epsilon) \rightarrow \mathbb{R}_+$ such that

- (i) $\varphi(0) = 0$, φ is C^1 on $(0, \epsilon)$, and for all $r \in (0, \epsilon)$, $\varphi'(r) > 0$;
- (ii) for all $x \in U \cap [H(\bar{x}) < H < H(\bar{x}) + \epsilon]$, the Kurdyka–Łojasiewicz inequality holds:

$$(2.3) \quad \varphi'(H(x) - H(\bar{x})) \text{dist}(0, \partial H(x)) \geq 1.$$

If H satisfies the KL property at each point of $\text{dom}(\partial H)$, then it is called KL functions.

Roughly speaking, KL functions become sharp up to reparameterization via φ , a desingularizing function for H . Typical KL functions include the class of semi-algebraic functions [4, 5]. For instance, the ℓ_0 pseudo-norm and the rank function are KL. Semi-algebraic functions admit desingularizing functions of the form $\varphi(r) = ar^{1-\theta}$ for $a > 0$, and $\theta \in [0, 1)$ is known as the KL exponent of the function [4, 6]. For these functions, the KL inequality reads

$$(2.4) \quad (H(x) - H(\bar{x}))^\theta \leq C\|\zeta\| \quad \forall \zeta \in \partial H(x),$$

for some $C > 0$. In the case $H(x) = H(\bar{x})$, we use the convention $0^0 \stackrel{\text{def}}{=} 0$.

Bounded Iterates. Many of our results require the assumption that the iterates generated by SPRING are bounded, in addition to assumptions (A.1)-(A.4). Because assumption (A.2) only requires ∇F_i to be Lipschitz on bounded sets, assuming the iterates are bounded allows us to say ∇F_i is M -Lipschitz continuous. We also require boundedness of the iterates to ensure that a limit point of this sequence exists during the proof of Lemma 4.3. This assumption is required for the same reasons in the analysis of PALM. It is satisfied, for example, if J and R have bounded domains.

Notation. We use $\{(x_k, y_k)\}_{k \in \mathbb{N}}$ to denote the sequence generated by SPRING. We use $L_x \stackrel{\text{def}}{=} \max_{k \in \mathbb{N}} L_1(y_k)$, and define L_y analogously. We set $\bar{L} \stackrel{\text{def}}{=} \max\{L_x, L_y\}$, $\bar{\gamma}_k \stackrel{\text{def}}{=} \max\{\gamma_{x,k}, \gamma_{y,k}\}$, $\underline{\gamma}_k \stackrel{\text{def}}{=} \min\{\gamma_{x,k}, \gamma_{y,k}\}$, and $\Phi \stackrel{\text{def}}{=} \inf_{(x,y) \in \text{dom}(\Phi)} \Phi(x, y)$. We also use L to denote the maximum of L_x, L_y , and M over the iterates generated by SPRING, so that $\bar{L}, M \leq L$. We use \mathbb{E}_k to denote the expectation conditional on the first k iterations of SPRING. Specifically, $\mathbb{E}_k \equiv \mathbb{E}[\cdot | \mathcal{F}_k]$ where \mathcal{F}_k is the σ -algebra generated by B_0, \dots, B_{k-1} . We require a notion of the expectation of the subdifferential of $\Phi(z_k)$. To define this, let $\bar{n} = \binom{n}{b}$ be the number of possible gradient estimates in one iteration of Algorithm 1.1, and let $\{z_k^i\}_{i=1}^{\bar{n}^k}$ be the set of possible values for z_k .⁴ We use the notation $\mathbb{E} \partial \Phi(z_k) = \partial \mathbb{E} \Phi(z_k) = \{\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \xi_i | \xi_i \in \partial \Phi(z_k^i)\}$. Every subgradient $\xi \in \partial \Phi(z_k)$ is of the form $\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \xi_i$ for $\xi_i \in \partial \Phi(z_k^i)$, and we denote this vector as $\mathbb{E} \xi \in \mathbb{E} \partial \Phi(z_k)$.

2.1. Elementary Lemmas. The following lemmas generalize the sufficient decrease property of proximal gradient descent to the stochastic-gradient setting. They allow us to show that, if the MSE of the stochastic gradient estimator is small enough, then iteratively applying the proximal gradient operator decreases the suboptimality of each iterate in expectation.

Lemma 2.5. *Let $F : \mathbb{R}^m \rightarrow \mathbb{R}$ be a function with L -Lipschitz continuous gradient, $R : \mathbb{R}^m \rightarrow \mathbb{R}$ a proper lower semicontinuous function that is bounded from below, and $z \in \text{prox}_{\eta R}(x - \eta d)$ for some $\eta > 0$ and $x, d \in \mathbb{R}^m$. Then, for all $y \in \mathbb{R}^m$,*

$$(2.5) \quad 0 \leq F(y) + R(y) - F(z) - R(z) + \langle \nabla F(x) - d, z - y \rangle + \left(\frac{L}{2} - \frac{1}{2\eta}\right) \|x - z\|^2 + \left(\frac{L}{2} + \frac{1}{2\eta}\right) \|x - y\|^2.$$

Proof. By the Lipschitz continuity of ∇F , we have the inequalities

$$\begin{aligned} F(x) - F(y) &\leq \langle \nabla F(x), x - y \rangle + \frac{L}{2} \|x - y\|^2, \\ F(z) - F(x) &\leq \langle \nabla F(x), z - x \rangle + \frac{L}{2} \|z - x\|^2. \end{aligned}$$

Furthermore, by the definition of z ,

$$z \in \text{Argmin}_{v \in \mathbb{R}^m} \left\{ \langle d, v - x \rangle + \frac{1}{2\eta} \|v - x\|^2 + R(v) \right\}.$$

Taking $v = y$, we obtain

$$0 \leq R(y) - R(z) + \langle d, y - z \rangle + \frac{1}{2\eta} (\|x - y\|^2 - \|x - z\|^2).$$

Adding these three inequalities completes the proof. ■

⁴When the proximal operator is multi-valued, Algorithm 1.1 requires one element to be chosen for each iterate, so we are not counting “possible” values for z_k that arise from choosing other elements of the proximal operator.

If the full gradient estimator is used, Lemma 2.5 implies the well-known sufficient decrease property of proximal gradient descent. Using a gradient estimator, this decrease is offset by the estimator's MSE. The following lemma quantifies this relationship.

Lemma 2.6 (Sufficient Decrease Property). *Let F , R , and z be defined as in Lemma 2.5. The following inequality holds for any $\lambda > 0$:*

$$(2.6) \quad 0 \leq F(x) + R(x) - F(z) - R(z) + \frac{1}{2L\lambda} \|d - \nabla F(x)\|^2 + \left(\frac{L(\lambda+1)}{2} - \frac{1}{2\eta}\right) \|x - z\|^2.$$

Proof. From Lemma 2.5 with $y = x$, we have

$$0 \leq F(x) + R(x) - F(z) - R(z) + \langle \nabla F(x) - d, z - x \rangle + \left(\frac{L}{2} - \frac{1}{2\eta}\right) \|x - z\|^2.$$

Using Young's inequality $\langle \nabla F(x) - d, z - x \rangle \leq \frac{1}{2L\lambda} \|d - \nabla F(x)\|^2 + \frac{L\lambda}{2} \|x - z\|^2$ we obtain the desired result. ■

As in a related work [14], we use the *supermartingale convergence theorem* to obtain almost sure convergence of sequences generated by SPRING. Below, we present an implication of this result adapted to our context. We refer to [14, Theorem 4.2] and [33, Theorem 1] for more general presentations.

Lemma 2.7 (Supermartingale Convergence). *Let $\{X_k\}_{k=0}^\infty$ and $\{Y_k\}_{k=0}^\infty$ be sequences of bounded non-negative random variables such that X_k, Y_k are functions of only the first k iterations of SPRING. If*

$$(2.7) \quad \mathbb{E}_k X_{k+1} + Y_k \leq X_k,$$

for all k , then $\sum_{k=0}^\infty Y_k < +\infty$ a.s. and X_k converges a.s.

3. Convergence rates of the generalized gradient map. To begin, we present our analysis of the convergence rate of the generalized gradient map defined in (1.4). The following results of Theorem 3.1 generalize many existing convergence guarantees for stochastic gradient methods on non-convex, non-smooth objectives [1, 18, 30, 37, 41]. Recall that $\bar{L} \stackrel{\text{def}}{=} \max\{L_x, L_y\}$, $\bar{\gamma}_k \stackrel{\text{def}}{=} \max\{\gamma_{x,k}, \gamma_{y,k}\}$, $\underline{\gamma}_k \stackrel{\text{def}}{=} \min\{\gamma_{x,k}, \gamma_{y,k}\}$, and $\Phi \stackrel{\text{def}}{=} \inf_{(x,y) \in \text{dom}(\Phi)} \Phi(x, y)$.

Theorem 3.1. *Suppose that assumptions (A.1)-(A.4) hold and that the sequence $\{(x_k, y_k)\}_{k \in \mathbb{N}}$ is bounded. Let $\tilde{\nabla}_x$ and $\tilde{\nabla}_y$ be variance-reduced gradient estimators following Definition 2.1.*

- Suppose $\bar{\gamma}_k$ is non-increasing, and for all k , $\gamma_{y,k} < \frac{1}{4L_y + 2M}$ and

$$\bar{\gamma}_k \leq \frac{1}{16} \sqrt{\frac{(\bar{L}+M)^2}{(V_1+V_T/\rho)^2} + \frac{16}{(V_1+V_T/\rho)}} - \frac{\bar{L}+M}{16(V_1+V_T/\rho)}, \quad 0 < \beta \leq \underline{\gamma}_k, \quad \gamma_{x,k} < \frac{1}{4L_x},$$

With α chosen uniformly at random from $\{0, 1, \dots, T-1\}$,

$$\mathbb{E}[\text{dist}(0, \mathcal{G}_{\frac{\gamma_{x,\alpha}}{2}, \frac{\gamma_{y,\alpha}}{2}}(z_\alpha))^2] \leq \frac{4(\Phi(x_0, y_0) + \frac{2\bar{\gamma}_0}{\rho} \Upsilon_0)}{T\nu\beta^2},$$

where $\nu \stackrel{\text{def}}{=} \min\{\frac{1}{4\gamma_{x,0}} - L_x, \frac{1}{4\gamma_{y,0}} - \frac{M}{2} - L_y\}$.

- If, moreover, Φ satisfies the error bound

$$(3.1) \quad \Phi(x_k, y_k) - \Phi \leq \mu \text{dist}(0, \mathcal{G}_{\frac{\gamma_{x,k}}{2}, \frac{\gamma_{y,k}}{2}}(x_k, y_k))^2,$$

for all $k \in \mathbb{N}$, and $\bar{\gamma}_k$ satisfies

$$\bar{\gamma}_k \leq \frac{1}{20} \sqrt{\frac{(\bar{L}+M)^2}{(V_1+V_T/\rho)^2} + \frac{20}{(V_1+V_T/\rho)}} - \frac{\bar{L}+M}{20(V_1+V_T/\rho)},$$

then the iterates of SPRING converge to the set of global minimizers of Φ , and after T iterations of Algorithm 1.1,

$$\mathbb{E}[\Phi(x_T, y_T) - \underline{\Phi}] \leq (1 - \Theta)^T (\Phi(x_0, y_0) - \underline{\Phi} + \frac{4\bar{\gamma}_0}{\rho} \Upsilon_0),$$

where $\Theta \stackrel{\text{def}}{=} \min\{\frac{\nu\beta^2}{4\mu}, \frac{\rho}{2}\}$.

Remark 3.2. We include convergence guarantees under the error bound (3.1) to compare with related works [1]. This error bound is similar to the Kurdyka–Łojasiewicz property for functions with a KL exponent of $1/2$, as can be seen comparing equation (3.1) to equation (2.4) with $\theta = 1/2$ and $H(\bar{x}) = \underline{\Phi}$. Although objectives satisfying this error bound could be non-convex, this condition ensures that convergence to the global minimum is guaranteed.

Proof of Theorem 3.1, Part 1. Let $\hat{x}_{k+1} \in \text{prox}_{\frac{\gamma_{x,k}}{2}J}(x_k - \frac{\gamma_{x,k}}{2}\nabla_x F(x_k, y_k))$, and let $\hat{y}_{k+1} \in \text{prox}_{\frac{\gamma_{y,k}}{2}R}(y_k - \frac{\gamma_{y,k}}{2}\nabla_y F(x_k, y_k))$. Applying Lemma 2.5 with $z = \hat{x}_{k+1}$, $y = x = x_k$ and $d = \nabla_x F(x_k, y_k)$, we have

$$F(\hat{x}_{k+1}, y_k) + J(\hat{x}_{k+1}) \leq F(x_k, y_k) + J(x_k) + (\frac{L_x}{2} - \frac{1}{\gamma_{x,k}})\|\hat{x}_{k+1} - x_k\|^2.$$

Again, applying Lemma 2.5 with $z = x_{k+1}$, $y = \hat{x}_{k+1}$, $x = x_k$, and $d = \tilde{\nabla}_x(x_k, y_k)$, we obtain

$$F(x_{k+1}, y_k) + J(x_{k+1}) \leq F(\hat{x}_{k+1}, y_k) + J(\hat{x}_{k+1}) + \langle \nabla_x F(x_k, y_k) - \tilde{\nabla}_x(x_k, y_k), x_{k+1} - \hat{x}_{k+1} \rangle + (\frac{L_x}{2} - \frac{1}{2\gamma_{x,k}})\|x_{k+1} - x_k\|^2 + (\frac{L_x}{2} + \frac{1}{2\gamma_{x,k}})\|\hat{x}_{k+1} - x_k\|^2.$$

Adding these two inequalities gives

$$\begin{aligned} & F(x_{k+1}, y_k) + J(x_{k+1}) \\ & \leq F(x_k, y_k) + J(x_k) + (L_x - \frac{1}{2\gamma_{x,k}})\|\hat{x}_{k+1} - x_k\|^2 + (\frac{L_x}{2} - \frac{1}{2\gamma_{x,k}})\|x_{k+1} - x_k\|^2 \\ & \quad + \langle \nabla_x F(x_k, y_k) - \tilde{\nabla}_x(x_k, y_k), x_{k+1} - \hat{x}_{k+1} \rangle \\ & \stackrel{\textcircled{1}}{\leq} F(x_k, y_k) + J(x_k) + (L_x - \frac{1}{2\gamma_{x,k}})\|\hat{x}_{k+1} - x_k\|^2 + (\frac{L_x}{2} - \frac{1}{2\gamma_{x,k}})\|x_{k+1} - x_k\|^2 \\ & \quad + 2\gamma_{x,k}\|\nabla_x F(x_k, y_k) - \tilde{\nabla}_x(x_k, y_k)\|^2 + \frac{1}{8\gamma_{x,k}}\|\hat{x}_{k+1} - x_{k+1}\|^2 \\ & \stackrel{\textcircled{2}}{\leq} F(x_k, y_k) + J(x_k) + (L_x - \frac{1}{4\gamma_{x,k}})\|\hat{x}_{k+1} - x_k\|^2 + (\frac{L_x}{2} - \frac{1}{4\gamma_{x,k}})\|x_{k+1} - x_k\|^2 \\ & \quad + 2\gamma_{x,k}\|\nabla_x F(x_k, y_k) - \tilde{\nabla}_x(x_k, y_k)\|^2. \end{aligned}$$

Inequality ① is Young's, and ② is the standard inequality $\|a - c\|^2 \leq 2\|a - b\|^2 + 2\|b - c\|^2$. For the

updates in y_k , we use Lemma 2.5 with $z = \hat{y}_{k+1}$, $y = x = y_k$, and $d = \nabla_y F(x_k, y_k)$, which gives

$$(3.3) \quad 0 \leq F(x_{k+1}, y_k) + R(y_k) - F(x_{k+1}, \hat{y}_{k+1}) - R(\hat{y}_{k+1}) \\ + \langle \nabla_y F(x_{k+1}, y_k) - \nabla_y F(x_k, y_k), \hat{y}_{k+1} - y_k \rangle + \left(\frac{L_y}{2} - \frac{1}{\gamma_{y,k}}\right) \|y_k - \hat{y}_{k+1}\|^2.$$

Finally, we apply Lemma 2.5 with $z = y_{k+1}$, $y = \hat{y}_{k+1}$, $x = y_k$, and $d = \tilde{\nabla}_y(x_{k+1}, y_k)$

$$(3.4) \quad 0 \leq F(x_{k+1}, \hat{y}_{k+1}) + R(\hat{y}_{k+1}) - F(x_{k+1}, y_{k+1}) - R(y_{k+1}) \\ + \langle \nabla_y F(x_{k+1}, y_k) - \tilde{\nabla}_y(x_{k+1}, y_k), y_{k+1} - \hat{y}_{k+1} \rangle + \left(\frac{L_y}{2} - \frac{1}{2\gamma_{y,k}}\right) \|y_k - y_{k+1}\|^2 \\ + \left(\frac{L_y}{2} + \frac{1}{2\gamma_{y,k}}\right) \|y_k - \hat{y}_{k+1}\|^2.$$

Adding these two inequalities and bounding the result as in (3.2), we obtain

$$(3.5) \quad F(x_{k+1}, y_{k+1}) + R(y_{k+1}) \\ \leq F(x_{k+1}, y_k) + R(y_k) + (L_y - \frac{1}{2\gamma_{y,k}}) \|\hat{y}_{k+1} - y_k\|^2 + \left(\frac{L_y}{2} - \frac{1}{2\gamma_{y,k}}\right) \|y_{k+1} - y_k\|^2 \\ + \langle \nabla_y F(x_{k+1}, y_k) - \nabla_y F(x_k, y_k), \hat{y}_{k+1} - y_k \rangle + \langle \nabla_y F(x_{k+1}, y_k) - \tilde{\nabla}_y(x_{k+1}, y_k), y_{k+1} - \hat{y}_{k+1} \rangle \\ \stackrel{\textcircled{1}}{\leq} F(x_{k+1}, y_k) + R(y_k) + (L_y - \frac{1}{4\gamma_{y,k}}) \|\hat{y}_{k+1} - y_k\|^2 + \left(\frac{L_y}{2} - \frac{1}{4\gamma_{y,k}}\right) \|y_{k+1} - y_k\|^2 \\ + \langle \nabla_y F(x_{k+1}, y_k) - \nabla_y F(x_k, y_k), \hat{y}_{k+1} - y_k \rangle + 2\gamma_{y,k} \|\nabla_y F(x_{k+1}, y_k) - \tilde{\nabla}_y(x_{k+1}, y_k)\|^2 \\ + \frac{1}{8\gamma_{y,k}} \|y_{k+1} - \hat{y}_{k+1}\|^2 \\ \stackrel{\textcircled{2}}{\leq} F(x_{k+1}, y_k) + R(y_k) + (L_y - \frac{1}{4\gamma_{y,k}}) \|\hat{y}_{k+1} - y_k\|^2 + \left(\frac{L_y}{2} - \frac{1}{4\gamma_{y,k}}\right) \|y_{k+1} - y_k\|^2 \\ + \langle \nabla_y F(x_{k+1}, y_k) - \nabla_y F(x_k, y_k), \hat{y}_{k+1} - y_k \rangle + 2\gamma_{y,k} \|\nabla_y F(x_{k+1}, y_k) - \tilde{\nabla}_y(x_{k+1}, y_k)\|^2 \\ \stackrel{\textcircled{3}}{\leq} F(x_{k+1}, y_k) + R(y_k) + (L_y - \frac{1}{4\gamma_{y,k}}) \|\hat{y}_{k+1} - y_k\|^2 + \left(\frac{L_y}{2} - \frac{1}{4\gamma_{y,k}}\right) \|y_{k+1} - y_k\|^2 \\ + \frac{1}{2M} \|\nabla_y F(x_{k+1}, y_k) - \nabla_y F(x_k, y_k)\|^2 + \frac{M}{2} \|\hat{y}_{k+1} - y_k\|^2 \\ + 2\gamma_{y,k} \|\nabla_y F(x_{k+1}, y_k) - \tilde{\nabla}_y(x_{k+1}, y_k)\|^2 \\ \stackrel{\textcircled{4}}{\leq} F(x_{k+1}, y_k) + R(y_k) + (L_y + \frac{M}{2} - \frac{1}{4\gamma_{y,k}}) \|\hat{y}_{k+1} - y_k\|^2 + \left(\frac{L_y}{2} - \frac{1}{4\gamma_{y,k}}\right) \|y_{k+1} - y_k\|^2 \\ + 2\gamma_{y,k} \|\nabla_y F(x_{k+1}, y_k) - \tilde{\nabla}_x(x_{k+1}, y_k)\|^2 + \frac{M}{2} \|x_{k+1} - x_k\|^2.$$

Inequalities ① and ③ are Young's, inequality ② follows from the fact that $\|a - c\|^2 \leq 2\|a - b\|^2 + 2\|b - c\|^2$, and ④ uses the assumptions that the sequence $\{(x_k, y_k)\}_{k \in \mathbb{N}}$ is bounded and ∇F is M -Lipschitz continuous on this bounded set.

Adding inequality (3.2) and inequality (3.5), we have

$$(3.6) \quad \Phi(x_{k+1}, y_{k+1}) \leq \Phi(x_k, y_k) + (L_x - \frac{1}{4\gamma_{x,k}}) \|\hat{x}_{k+1} - x_k\|^2 + (L_y + \frac{M}{2} - \frac{1}{4\gamma_{y,k}}) \|\hat{y}_{k+1} - y_k\|^2 \\ + \left(\frac{L_x}{2} + \frac{M}{2} - \frac{1}{4\gamma_{x,k}}\right) \|x_{k+1} - x_k\|^2 + \left(\frac{L_y}{2} - \frac{1}{4\gamma_{y,k}}\right) \|y_{k+1} - y_k\|^2 \\ + 2\bar{\gamma}_k (\|\nabla_x F(x_k, y_k) - \tilde{\nabla}_x(x_k, y_k)\|^2 + \|\nabla_y F(x_{k+1}, y_k) - \tilde{\nabla}_y(x_{k+1}, y_k)\|^2),$$

where $\bar{\gamma}_k = \max\{\gamma_{x,k}, \gamma_{y,k}\}$. We apply the conditional expectation operator \mathbb{E}_k and bound the MSE terms using (2.1). This gives

$$\begin{aligned} & \mathbb{E}_k[\Phi(x_{k+1}, y_{k+1}) + (-\frac{L_x}{2} - \frac{M}{2} - 2V_1\bar{\gamma}_k + \frac{1}{4\gamma_{x,k}})\|x_{k+1} - x_k\|^2 \\ & \quad + (-\frac{L_y}{2} - 2V_1\bar{\gamma}_k + \frac{1}{4\gamma_{y,k}})\|y_{k+1} - y_k\|^2] \\ (3.7) \quad & \leq \Phi(x_k, y_k) + (L_x - \frac{1}{4\gamma_{x,k}})\|\hat{x}_{k+1} - x_k\|^2 + (L_y + \frac{M}{2} - \frac{1}{4\gamma_{y,k}})\|\hat{y}_{k+1} - y_k\|^2 + 2\bar{\gamma}_k\Upsilon_k \\ & \quad + 2V_1\bar{\gamma}_k\|z_k - z_{k-1}\|^2. \end{aligned}$$

Next, we use (2.2) to say

$$2\bar{\gamma}_k\Upsilon_k \leq \frac{2\bar{\gamma}_k}{\rho}(-\mathbb{E}_k\Upsilon_{k+1} + \Upsilon_k + V_\Upsilon(\mathbb{E}_k\|z_{k+1} - z_k\|^2 + \|z_k - z_{k-1}\|^2)).$$

Adding the previous two inequalities, we have

$$\begin{aligned} & \mathbb{E}_k[\Phi(x_{k+1}, y_{k+1}) + (-\frac{L_x}{2} - \frac{M}{2} - 2V_1\bar{\gamma}_k - \frac{2V_\Upsilon\bar{\gamma}_k}{\rho} + \frac{1}{4\gamma_{x,k}})\|x_{k+1} - x_k\|^2 \\ & \quad + (-\frac{L_y}{2} - 2V_1\bar{\gamma}_k - \frac{2V_\Upsilon\bar{\gamma}_k}{\rho} + \frac{1}{4\gamma_{y,k}})\|y_{k+1} - y_k\|^2 + \frac{2\bar{\gamma}_k}{\rho}\Upsilon_{k+1}] \\ (3.8) \quad & \leq \Phi(x_k, y_k) + (L_x - \frac{1}{4\gamma_{x,k}})\|\hat{x}_{k+1} - x_k\|^2 + (L_y + \frac{M}{2} - \frac{1}{4\gamma_{y,k}})\|\hat{y}_{k+1} - y_k\|^2 + \frac{2\bar{\gamma}_k}{\rho}\Upsilon_k \\ & \quad + 2\bar{\gamma}_k(V_1 + \frac{V_\Upsilon}{\rho})\|z_k - z_{k-1}\|^2. \end{aligned}$$

Let $\bar{L} = \max\{L_x, L_y\}$. To ensure that the coefficients of $\|x_{k+1} - x_k\|^2$ and $\|y_{k+1} - y_k\|^2$ are non-negative, we set

$$(3.8) \quad \bar{\gamma}_k \leq \frac{1}{16} \sqrt{\frac{(\bar{L}+M)^2}{(V_1+V_\Upsilon/\rho)^2} + \frac{16}{(V_1+V_\Upsilon/\rho)}} - \frac{\bar{L}+M}{16(V_1+V_\Upsilon/\rho)},$$

for all $k \in \mathbb{N}$. With this choice,

$$\begin{aligned} & (-\frac{L_x+M}{2} - 2V_1\bar{\gamma}_k - \frac{2V_\Upsilon\bar{\gamma}_k}{\rho} + \frac{1}{4\gamma_{x,k}})\|x_{k+1} - x_k\|^2 + (-\frac{L_y}{2} - 2V_1\bar{\gamma}_k - \frac{2V_\Upsilon\bar{\gamma}_k}{\rho} \\ & \quad + \frac{1}{4\gamma_{y,k}})\|y_{k+1} - y_k\|^2 \\ (3.9) \quad & \geq (-\frac{L_x+M}{2} - 2V_1\bar{\gamma}_k - \frac{2V_\Upsilon\bar{\gamma}_k}{\rho} + \frac{1}{4\bar{\gamma}_k})\|z_{k+1} - z_k\|^2 \\ & \geq 2\bar{\gamma}_k(V_1 + V_\Upsilon/\rho)\|z_{k+1} - z_k\|^2. \end{aligned}$$

The final inequality is due to the bound in (3.8). To ensure that the coefficients of $\|\hat{x}_{k+1} - x_k\|^2$ and $\|\hat{y}_{k+1} - y_k\|^2$ are non-positive, we set $\gamma_{x,k} < \frac{1}{4L_x}$ and $\gamma_{y,k} < \frac{1}{4L_y+2M}$, which yields

$$\begin{aligned} & \mathbb{E}_k[\Phi(x_{k+1}, y_{k+1}) + 2\bar{\gamma}_k(V_1 + V_\Upsilon/\rho)\|z_{k+1} - z_k\|^2 + \frac{2\bar{\gamma}_k}{\rho}\Upsilon_{k+1}] \\ (3.10) \quad & \leq \Phi(x_k, y_k) + (L_x - \frac{1}{4\gamma_{x,k}})\|\hat{x}_{k+1} - x_k\|^2 + (L_y - \frac{1}{4\gamma_{y,k}})\|\hat{y}_{k+1} - y_k\|^2 \\ & \quad + 2\bar{\gamma}_k(V_1 + V_\Upsilon/\rho)\|z_k - z_{k-1}\|^2 + \frac{2\bar{\gamma}_k}{\rho}\Upsilon_k. \end{aligned}$$

Because $\bar{\gamma}_k$ is non-increasing,

$$\begin{aligned} & \mathbb{E}_k[\Phi(x_{k+1}, y_{k+1}) + 2\bar{\gamma}_{k+1}(V_1 + V_\Upsilon/\rho)\|z_{k+1} - z_k\|^2 + \frac{2\bar{\gamma}_{k+1}}{\rho}\Upsilon_{k+1}] \\ (3.11) \quad & \leq \Phi(x_k, y_k) - \nu\|\hat{z}_{k+1} - z_k\|^2 + 2\bar{\gamma}_k(V_1 + V_\Upsilon/\rho)\|z_k - z_{k-1}\|^2 + \frac{2\bar{\gamma}_k}{\rho}\Upsilon_k, \end{aligned}$$

where $\nu = \min\{\frac{1}{4\gamma_{x,0}} - L_x, \frac{1}{4\gamma_{y,0}} - \frac{M}{2} - L_y\}$. Applying the full expectation operator, summing from $k = 0$ to $k = T - 1$, and using the convention $z_{-1} = z_0$ gives

$$\frac{2\bar{\gamma}_T}{\rho}\Upsilon_T + 2\bar{\gamma}_T(V_1 + V_\Upsilon/\rho)\|z_T - z_{T-1}\|^2 + \nu \sum_{k=0}^{T-1} \mathbb{E}\|\hat{z}_{k+1} - z_k\|^2 \leq \Phi(x_0, y_0) + \frac{2\bar{\gamma}_0}{\rho}\Upsilon_0.$$

We drop the first two terms on the left from the inequality as they are non-negative. Let α be drawn uniformly at random from the set $\{0, 1, \dots, T - 1\}$, and recall $\underline{\gamma}_k \geq \beta$. Using the fact that $\|\hat{z}_{k+1} - z_k\|^2 \geq \frac{\beta^2}{4} \text{dist}(0, \mathcal{G}_{\frac{\gamma_{x,k}}{2}, \frac{\gamma_{y,k}}{2}}(z_k))^2$,

$$\mathbb{E} \text{dist}\left(0, \mathcal{G}_{\frac{\gamma_{x,\alpha}}{2}, \frac{\gamma_{y,\alpha}}{2}}(z_\alpha)\right)^2 \leq \frac{4(\Phi(x_0, y_0) + \frac{2\bar{\gamma}_0}{\rho}\Upsilon_0)}{T\nu\beta^2},$$

which completes the proof of the first claim. ■

Combining the same argument with the error bound (3.1), we obtain a linear convergence rate to the global optimum.

Proof of Theorem 3.1, Part 2. We begin with equation (3.7):

$$\begin{aligned} & \mathbb{E}_k[\Phi(x_{k+1}, y_{k+1}) + (-\frac{L_x}{2} - \frac{M}{2} - 2V_1\gamma_{x,k} + \frac{1}{4\gamma_{x,k}})\|x_{k+1} - x_k\|^2 \\ & + (-\frac{L_y}{2} - 2V_1\gamma_{y,k} + \frac{1}{4\gamma_{y,k}})\|y_{k+1} - y_k\|^2] \\ & \leq \Phi(x_k, y_k) - \nu\|\hat{z}_{k+1} - z_k\|^2 + 2\bar{\gamma}_k\Upsilon_k + 2V_1\bar{\gamma}_k\|z_k - z_{k-1}\|^2. \end{aligned}$$

Using (2.2), we can say for any $c > 0$,

$$0 \leq \frac{2c\bar{\gamma}_k}{\rho}(-\mathbb{E}_k\Upsilon_{k+1} + (1 - \rho)\Upsilon_k + V_\Upsilon(\|z_{k+1} - z_k\|^2 + \|z_k - z_{k-1}\|^2)).$$

Adding the previous two inequalities, we have

$$\begin{aligned} & \mathbb{E}_k[\Phi(x_{k+1}, y_{k+1}) + (-\frac{L_x}{2} - \frac{M}{2} - 2V_1\gamma_{x,k} - \frac{2cV_\Upsilon\bar{\gamma}_k}{\rho} + \frac{1}{4\gamma_{x,k}})\|x_{k+1} - x_k\|^2 \\ & + (-\frac{L_y}{2} - 2V_1\gamma_{y,k} - \frac{2cV_\Upsilon\bar{\gamma}_k}{\rho} + \frac{1}{4\gamma_{y,k}})\|y_{k+1} - y_k\|^2 + \frac{2c\bar{\gamma}_k}{\rho}\Upsilon_{k+1}] \\ & \leq \Phi(x_k, y_k) - \nu\|\hat{z}_{k+1} - z_k\|^2 + 2\bar{\gamma}_k(V_1 + \frac{cV_\Upsilon}{\rho})\|z_k - z_{k-1}\|^2 + \frac{2c\bar{\gamma}_k}{\rho}(1 + \frac{\rho}{c} - \rho)\Upsilon_k. \end{aligned}$$

We apply the error bound assumption (3.1) to say

$$-\nu\|\hat{z}_{k+1} - z_k\|^2 \leq -\frac{\nu\gamma_k^2}{4} \text{dist}(0, \mathcal{G}_{\frac{\gamma_{x,k}}{2}, \frac{\gamma_{y,k}}{2}}(z_k))^2 \leq -\frac{\nu\gamma_k^2}{4\mu}(\Phi(x_k, y_k) - \Phi).$$

In total, we have

$$\begin{aligned} & \mathbb{E}_k[\Phi(x_{k+1}, y_{k+1}) - \Phi + (-\frac{L_x}{2} - \frac{M}{2} - 2V_1\gamma_{x,k} - \frac{2cV_\Upsilon\bar{\gamma}_k}{\rho} + \frac{1}{4\gamma_{x,k}})\|x_{k+1} - x_k\|^2 \\ & + (-\frac{L_y}{2} - 2V_1\gamma_{y,k} - \frac{2cV_\Upsilon\bar{\gamma}_k}{\rho} + \frac{1}{4\gamma_{y,k}})\|y_{k+1} - y_k\|^2 + \frac{2c\bar{\gamma}_k}{\rho}\Upsilon_{k+1}] \\ & \leq (1 - \frac{\nu\gamma_k^2}{4\mu})(\Phi(x_k, y_k) - \Phi) + 2\bar{\gamma}_k(V_1 + \frac{cV_\Upsilon}{\rho})\|z_k - z_{k-1}\|^2 + \frac{2c\bar{\gamma}_k}{\rho}(1 + \frac{\rho}{c} - \rho)\Upsilon_k. \end{aligned}$$

326 Choosing $c = 2$, setting the step-sizes so that they satisfy, for all k ,

$$327 \quad \bar{\gamma}_k \leq \frac{1}{20} \sqrt{\frac{(\bar{L}+M)^2}{(V_1+2V_{\Upsilon}/\rho)^2} + \frac{20}{(V_1+2V_{\Upsilon}/\rho)}} - \frac{\bar{L}+M}{20(V_1+2V_{\Upsilon}/\rho)}, \quad \gamma_{x,k} < \frac{1}{4L_x}, \quad \gamma_{y,k} < \frac{1}{4L_y+2M}, \quad 0 < \beta \leq \underline{\gamma}_k,$$

328 and letting $\Theta = \min\{\frac{\nu\beta^2}{4\mu}, \frac{\rho}{2}\}$, we have

$$329 \quad \begin{aligned} & \mathbb{E}_k[\Phi(x_{k+1}, y_{k+1}) - \underline{\Phi} + 2\bar{\gamma}_k(V_1 + \frac{2V_{\Upsilon}}{\rho})\|z_{k+1} - z_k\|^2 + \frac{4\bar{\gamma}_k}{\rho}\Upsilon_{k+1}] \\ & \leq (1 - \Theta)(\Phi(x_k, y_k) - \underline{\Phi} + 2\bar{\gamma}_k(V_1 + \frac{2V_{\Upsilon}}{\rho})\|z_k - z_{k-1}\|^2 + \frac{4\bar{\gamma}_k}{\rho}\Upsilon_k). \end{aligned}$$

330 Because $\bar{\gamma}_k$ is non-increasing,

$$331 \quad \begin{aligned} & \mathbb{E}_k[\Phi(x_{k+1}, y_{k+1}) - \underline{\Phi} + 2\bar{\gamma}_{k+1}(V_1 + \frac{2V_{\Upsilon}}{\rho})\|z_{k+1} - z_k\|^2 + \frac{4\bar{\gamma}_{k+1}}{\rho}\Upsilon_{k+1}] \\ & \leq (1 - \Theta)(\Phi(x_k, y_k) - \underline{\Phi} + 2\bar{\gamma}_k(V_1 + \frac{2V_{\Upsilon}}{\rho})\|z_k - z_{k-1}\|^2 + \frac{4\bar{\gamma}_k}{\rho}\Upsilon_k). \end{aligned}$$

332 Applying the full expectation operator, chaining this inequality over the iterations $k = 0$ to $k = T - 1$,
333 and using the convention $z_{-1} = z_0$,

$$334 \quad \mathbb{E}[\Phi(x_T, y_T) - \underline{\Phi}] \leq (1 - \Theta)^T (\Phi(x_0, y_0) - \underline{\Phi} + \frac{4\bar{\gamma}_0}{\rho}\Upsilon_0),$$

335 which completes the proof. ■

336 Because SAGA and SARAH gradient estimators are variance-reduced, Theorem 3.1 implies specific
337 convergence rates for Algorithm 1.1 when using these estimators.

338 **Corollary 3.3.** *To compute an ϵ -approximate critical point in expectation, Algorithm 1.1 using*

- 339 • *SARAH gradient estimator with $p = n$, $\bar{\gamma}_k \leq \frac{1}{2L\sqrt{30n}}$ and any batch size requires no more than*
340 *$\mathcal{O}(L\sqrt{n}/\epsilon^2)$ SFO calls;*
- 341 • *SAGA gradient estimator with $b = n^{2/3}$ and $\bar{\gamma}_k \leq \frac{1}{2\sqrt{2710}L}$ requires no more than $\mathcal{O}(Ln^{2/3}/\epsilon^2)$*
342 *SFO calls.⁵*

343 *If Φ satisfies the error bound condition (3.1), then to compute an ϵ -suboptimal point in expectation,*
344 *Algorithm 1.1 using*

- 345 • *the SARAH gradient estimator requires no more than $\mathcal{O}((n + L\sqrt{n}/\mu) \log(1/\epsilon))$ SFO calls;*
- 346 • *the SAGA gradient estimator requires no more than $\mathcal{O}((n + Ln^{2/3}/\mu) \log(1/\epsilon))$ SFO calls.*

347 **Remark 3.4.** The improved dependence on n when using SARAH gradient estimator exists in all
348 of our convergence rates for SPRING. Because most existing works on stochastic optimization for non-
349 smooth, non-convex problems use models that are special cases of (1.1), our results for SPRING capture
350 most existing work as special cases. In particular, in the case $R \equiv J \equiv 0$, our results recover recent
351 results showing that SARAH achieves the *oracle complexity lower-bound* for non-convex problems with
352 a finite-sum structure [18, 28, 37, 40, 41].

⁵For ease of exposition, we do not optimize over constants, so these step-sizes (particularly for the SAGA estimator) are not optimal. In general, we find the step-sizes suggested by theory to be conservative in practice (see Section 5 for details regarding practical step-sizes).

4. Convergence Rate under the KL Property. The results from previous section require only assumptions (A.1) to (A.4). To prove convergence of the sequence of the algorithm, and to obtain convergence rates depending on the KL exponent of the objective, we further require that Φ is semi-algebraic. In this section, under these assumptions, we prove convergence of the sequence and extend the convergence rates of PALM to SPRING. To derive these results, we first derive some preparatory results which generalize claims of PALM [6] to the stochastic setting. Given $k \in \mathbb{N}$, define the quantity

$$(4.1) \quad \Psi_k \stackrel{\text{def}}{=} \Phi(z_k) + \frac{1}{2\rho\sqrt{2(V_1+V_Y/\rho)}}\Upsilon_k + \frac{\sqrt{V_1+V_Y/\rho}}{\sqrt{2}}\|z_k - z_{k-1}\|^2.$$

Our first result guarantees that Ψ_k is decreasing in expectation.

Lemma 4.1 (ℓ_2 summability). *Let $\{z_k\}_{k=0}^\infty$ be the sequence generated by SPRING with $\bar{\gamma}_k$ non-increasing and satisfying $\bar{\gamma}_k < \frac{\sqrt{2}}{5(\sqrt{V_1+V_Y/\rho}+\bar{L})}$, $\forall k$, then Ψ_k satisfies*

$$(4.2) \quad \mathbb{E}_k \Psi_{k+1} \leq \Psi_k + \left(\frac{\bar{L}}{2} + \frac{3}{2}\sqrt{2(V_1+V_Y/\rho)} - \frac{1}{2\bar{\gamma}_k}\right)\mathbb{E}_k\|z_{k+1} - z_k\|^2 - \frac{\sqrt{V_1+V_Y/\rho}}{2\sqrt{2}}\|z_k - z_{k-1}\|^2,$$

and the expectation of the squared distance between the iterates is summable:

$$\sum_{k=0}^\infty \mathbb{E} [\|x_{k+1} - x_k\|^2 + \|y_{k+1} - y_k\|^2] = \sum_{k=0}^\infty \mathbb{E} \|z_{k+1} - z_k\|^2 < \infty.$$

Proof. Applying Lemma 2.6 twice, once for the update in x_k and once for the update in y_k , we have

$$F(x_{k+1}, y_k) + J(x_{k+1}) \leq F(x_k, y_k) + J(x_k) + \frac{1}{2L\lambda}\|\tilde{\nabla}_x(x_k, y_k) - \nabla_x F(x_k, y_k)\|^2 + \left(\frac{\bar{L}(\lambda+1)}{2} - \frac{1}{2\gamma_{x,k}}\right)\|x_{k+1} - x_k\|^2,$$

as well as

$$F(x_{k+1}, y_{k+1}) + R(y_{k+1}) \leq F(x_{k+1}, y_k) + R(y_k) + \left(\frac{\bar{L}(\lambda+1)}{2} - \frac{1}{2\gamma_{y,k}}\right)\|y_{k+1} - y_k\|^2 + \frac{1}{2L\lambda}\|\tilde{\nabla}_y(x_{k+1}, y_k) - \nabla_y F(x_{k+1}, y_k)\|^2.$$

Adding these inequalities together,

$$\Phi(x_{k+1}, y_{k+1}) \leq \Phi(x_k, y_k) + \frac{1}{2L\lambda}\|\tilde{\nabla}_x(x_k, y_k) - \nabla_x F(x_k, y_k)\|^2 + \frac{1}{2L\lambda}\|\tilde{\nabla}_y(x_{k+1}, y_k) - \nabla_y F(x_{k+1}, y_k)\|^2 + \left(\frac{\bar{L}(\lambda+1)}{2} - \frac{1}{2\bar{\gamma}_k}\right)\|z_{k+1} - z_k\|^2.$$

Applying the conditional expectation operator \mathbb{E}_k , we can bound the MSE terms using (2.1). This gives

$$(4.3) \quad \mathbb{E}_k [\Phi(z_{k+1}) + \left(-\frac{\bar{L}(\lambda+1)}{2} - \frac{V_1}{2L\lambda} + \frac{1}{2\bar{\gamma}_k}\right)\|z_{k+1} - z_k\|^2] \leq \Phi(z_k) + \frac{1}{2L\lambda}\Upsilon_k + \frac{V_1}{2L\lambda}\|z_k - z_{k-1}\|^2.$$

Next, we use (2.2) to say that

$$\frac{1}{2L\lambda}\Upsilon_k \leq \frac{1}{2L\lambda\rho}(-\mathbb{E}_k \Upsilon_{k+1} + \Upsilon_k + V_Y(\mathbb{E}_k\|z_{k+1} - z_k\|^2 + \|z_k - z_{k-1}\|^2)).$$

376 Combining these inequalities, we have

$$377 \quad \begin{aligned} & \mathbb{E}_k \left[\Phi(z_{k+1}) + \frac{1}{2L\lambda\rho} \Upsilon_{k+1} + \left(-\frac{\bar{L}(\lambda+1)}{2} - \frac{V_1+V_\Upsilon/\rho}{2L\lambda} + \frac{1}{2\bar{\gamma}_k} \right) \|z_{k+1} - z_k\|^2 \right] \\ & \leq \Phi(z_k) + \frac{1}{2L\lambda\rho} \Upsilon_k + \frac{V_1+V_\Upsilon/\rho}{2L\lambda} \|z_k - z_{k-1}\|^2. \end{aligned}$$

378 This is equivalent to

$$379 \quad \begin{aligned} & \mathbb{E}_k \left[\Phi(z_{k+1}) + \frac{1}{2L\lambda\rho} \Upsilon_{k+1} + \left(\frac{V_1+V_\Upsilon/\rho}{2L\lambda} + Z \right) \|z_{k+1} - z_k\|^2 \right. \\ & \quad \left. + \left(-\frac{\bar{L}(\lambda+1)}{2} - \frac{V_1+V_\Upsilon/\rho}{L\lambda} - Z + \frac{1}{2\bar{\gamma}_k} \right) \|z_{k+1} - z_k\|^2 \right] \\ & \leq \Phi(z_k) + \frac{1}{2L\lambda\rho} \Upsilon_k + \left(\frac{V_1+V_\Upsilon/\rho}{2L\lambda} + Z \right) \|z_k - z_{k-1}\|^2 - Z \|z_k - z_{k-1}\|^2, \end{aligned}$$

380 for any constant $Z \geq 0$. We use the choice $Z = \frac{\sqrt{V_1+V_\Upsilon/\rho}}{2\sqrt{2}}$ to simplify later arguments. Setting
 381 $\bar{\gamma}_k \leq (2(\frac{\bar{L}(\lambda+1)}{2} + \frac{V_1+V_\Upsilon/\rho}{L\lambda} + Z))^{-1}$, setting $\lambda = \frac{\sqrt{2(V_1+V_\Upsilon/\rho)}}{\bar{L}}$ to approximately maximize this bound
 382 on $\bar{\gamma}_k$, and using the fact that $\bar{\gamma}_k$ is non-increasing, we have

$$383 \quad (4.4) \quad \mathbb{E}_k \Psi_{k+1} \leq \Psi_k + \left(\frac{\bar{L}(\lambda+1)}{2} + \frac{V_1+V_\Upsilon/\rho}{L\lambda} + Z - \frac{1}{2\bar{\gamma}_k} \right) \mathbb{E}_k \|z_{k+1} - z_k\|^2 - Z \|z_k - z_{k-1}\|^2,$$

384 proving the first claim that Ψ_k is decreasing in expectation.

385 To prove the second claim, we apply the full expectation operator to (4.4) and sum the resulting
 386 inequality from $k = 0$ to $k = T - 1$,

$$387 \quad \mathbb{E} \Psi_T \leq \Psi_0 + \sum_{k=0}^{T-1} \left(\frac{\bar{L}(\lambda+1)}{2} + \frac{V_1+V_\Upsilon/\rho}{L\lambda} + Z - \frac{1}{2\bar{\gamma}_k} \right) \mathbb{E} \|z_{k+1} - z_k\|^2 - Z \mathbb{E} \|z_k - z_{k-1}\|^2.$$

388 Rearranging and using the facts that $\Phi \leq \Psi_T$ and $\bar{\gamma}_k$ is non-increasing,

$$389 \quad (4.5) \quad \sum_{k=0}^{T-1} \left(\frac{1}{2\bar{\gamma}_k} - \frac{\bar{L}(\lambda+1)}{2} - \frac{V_1+V_\Upsilon/\rho}{L\lambda} - Z \right) \mathbb{E} \|z_{k+1} - z_k\|^2 + Z \mathbb{E} \|z_k - z_{k-1}\|^2 \leq \Psi_0 - \Phi.$$

390 Taking the limit $T \rightarrow +\infty$ proves that the sequence $\mathbb{E} \|z_{k+1} - z_k\|^2$ is summable. ■

391 The next lemma establishes a bound on the norm of the subgradients of $\Phi(z_k)$.

392 **Lemma 4.2 (Subgradient Bound).** *Let $\{z_k\}_{k \in \mathbb{N}}$ be the sequence generated by SPRING with*
 393 *step-sizes satisfying $0 < \beta \leq \underline{\gamma}_k$. Define*

$$394 \quad \begin{aligned} A_x^k & \stackrel{\text{def}}{=} 1/\gamma_{x,k}(x_{k-1} - x_k) + \nabla_x F(x_k, y_k) - \tilde{\nabla}_x(x_{k-1}, y_{k-1}) \quad \text{and} \\ A_y^k & \stackrel{\text{def}}{=} 1/\gamma_{y,k}(y_{k-1} - y_k) + \nabla_y F(x_k, y_k) - \tilde{\nabla}_y(x_k, y_{k-1}). \end{aligned}$$

395 Then $(A_x^k, A_y^k) \in \partial\Phi(x_k, y_k)$ and, with $p = 1/\beta + M + L_y + V_2$,

$$396 \quad (4.6) \quad \mathbb{E}_{k-1} \|(A_x^k, A_y^k)\| \leq p(\mathbb{E}_{k-1} \|z_k - z_{k-1}\| + \|z_{k-1} - z_{k-2}\|) + \Gamma_{k-1}.$$

Proof. The fact that $(A_x^k, A_y^k) \in \partial\Phi(x_k, y_k)$ is clear from the definition of the proximal operator:

$$\begin{aligned} \frac{1}{\gamma_{x,k}}(x_{k-1} - x_k) - \tilde{\nabla}_x(x_{k-1}, y_{k-1}) &\in \partial J(x_k), \\ \frac{1}{\gamma_{y,k}}(y_{k-1} - y_k) - \tilde{\nabla}_y(x_k, y_{k-1}) &\in \partial R(y_k). \end{aligned}$$

Combining this with the fact that $\partial\Phi(x_k, y_k) = (\nabla_x F(x_k, y_k) + \partial J(x_k), \nabla_y F(x_k, y_k) + \partial R(y_k))$ makes it clear that $(A_x^k, A_y^k) \in \partial\Phi(x_k, y_k)$. All that remains is to bound the norms of A_x^k and A_y^k . Because ∇F is M -Lipschitz continuous on bounded sets,

$$\begin{aligned} \mathbb{E}_{k-1} \|A_x^k\| &\leq \frac{1}{\gamma_{x,k}} \mathbb{E}_{k-1} \|x_{k-1} - x_k\| + \mathbb{E}_{k-1} \|\nabla_x F(x_k, y_k) - \tilde{\nabla}_x(x_{k-1}, y_{k-1})\| \\ &\leq \frac{1}{\gamma_{x,k}} \mathbb{E}_{k-1} \|x_{k-1} - x_k\| + \mathbb{E}_{k-1} \|\nabla_x F(x_k, y_k) - \nabla_x F(x_{k-1}, y_{k-1})\| \\ &\quad + \mathbb{E}_{k-1} \|\nabla_x F(x_{k-1}, y_{k-1}) - \tilde{\nabla}_x(x_{k-1}, y_{k-1})\| \\ &\leq \left(\frac{1}{\gamma_{x,k}} + M\right) \mathbb{E}_{k-1} \|x_{k-1} - x_k\| + M \mathbb{E}_{k-1} \|y_k - y_{k-1}\| \\ &\quad + \mathbb{E}_{k-1} \|\nabla_x F(x_{k-1}, y_{k-1}) - \tilde{\nabla}_x(x_{k-1}, y_{k-1})\|. \end{aligned} \tag{4.7}$$

A similar argument holds for $\|A_y^k\|$.

$$\begin{aligned} \mathbb{E}_{k-1} \|A_y^k\| &\leq \frac{1}{\gamma_{y,k}} \mathbb{E}_{k-1} \|y_{k-1} - y_k\| + \mathbb{E}_{k-1} \|\nabla_y F(x_k, y_k) - \tilde{\nabla}_y(x_k, y_{k-1})\| \\ &\leq \frac{1}{\gamma_{y,k}} \mathbb{E}_{k-1} \|y_{k-1} - y_k\| + \mathbb{E}_{k-1} \|\nabla_y F(x_k, y_k) - \nabla_y F(x_k, y_{k-1})\| \\ &\quad + \mathbb{E}_{k-1} \|\nabla_y F(x_k, y_{k-1}) - \tilde{\nabla}_y(x_k, y_{k-1})\| \\ &\leq \left(\frac{1}{\gamma_{y,k}} + L_y\right) \mathbb{E}_{k-1} \|y_{k-1} - y_k\| + \mathbb{E}_{k-1} \|\nabla_y F(x_k, y_{k-1}) - \tilde{\nabla}_y(x_k, y_{k-1})\|. \end{aligned}$$

Adding these two inequalities together and using equation (2.1) to bound the MSE terms, we get

$$\mathbb{E}_{k-1} \|(A_x^k, A_y^k)\| \leq \mathbb{E}_{k-1} [\|A_x^k\| + \|A_y^k\|] \leq p(\mathbb{E}_{k-1} \|z_k - z_{k-1}\| + \|z_{k-1} - z_{k-2}\|) + \Gamma_{k-1},$$

where $p = 1/\beta + M + L_y + V_2$. ■

Define the set of limit points of $\{z_k\}_{k=0}^\infty$ as

$$\omega \stackrel{\text{def}}{=} \{z : \exists \text{ an increasing sequence of integers } \{k_\ell\}_{\ell \in \mathbb{N}} \text{ such that } z_{k_\ell} \rightarrow z \text{ as } \ell \rightarrow +\infty\}.$$

The following lemma describes properties of ω .

Lemma 4.3 (Limit points of $\{z_k\}_{k=0}^\infty$). Suppose assumptions (A.1)-(A.4) hold, that the sequence $z_k = (x_k, y_k)$ is bounded, and the step-sizes of Algorithm 1.1 satisfy the following conditions:

$$\gamma_{x,k}, \gamma_{y,k} \in \left[\beta, \frac{\sqrt{2}}{5(\sqrt{V_1 + V_T}/\rho + \bar{L})}\right) \quad \forall k,$$

and $\bar{\gamma}_k$ is non-increasing. Then

- (1). $\sum_{k=1}^\infty \|z_k - z_{k-1}\|^2 < \infty$ a.s., and $\|z_k - z_{k-1}\| \rightarrow 0$ a.s.;
- (2). $\mathbb{E}\Phi(z_k) \rightarrow \Phi^*$, where $\Phi^* \in [\underline{\Phi}, \infty)$;
- (3). $\mathbb{E}\text{dist}(0, \partial\Phi(z_k)) \rightarrow 0$;

(4). The set ω is non-empty, and for all $z^* \in \omega$, $\mathbb{E}\text{dist}(0, \partial\Phi(z^*)) = 0$;

(5). $\text{dist}(z_k, \omega) \rightarrow 0$ a.s.;

(6). ω is a.s. compact and connected;

(7). $\mathbb{E}\Phi(z^*) = \Phi^*$ for all $z^* \in \omega$.

Remark 4.4. The boundedness of z_k is also imposed in the original PALM [6] and asynchronous PALM [14], it can be satisfied automatically if, for instance, each regularizer has bounded domain.

Proof. By Lemma 4.1, we have

$$\mathbb{E}_k \Psi_{k+1} + \mathcal{O}(\|z_k - z_{k-1}\|^2) \leq \Psi_k.$$

The supermartingale convergence theorem implies that $\sum_{k=1}^{\infty} \|z_k - z_{k-1}\|^2 < +\infty$ a.s., and it follows that $\|z_k - z_{k-1}\| \rightarrow 0$ a.s. This proves Claim 1.

The supermartingale convergence theorem also ensures Ψ_k converges a.s. to a finite, positive random variable. Because $\|z_k - z_{k-1}\| \rightarrow 0$ a.s. and $\tilde{\nabla}$ is variance-reduced so $\mathbb{E}\Upsilon_k \rightarrow 0$, we can say $\lim_{k \rightarrow \infty} \mathbb{E}\Psi_k = \lim_{k \rightarrow \infty} \mathbb{E}\Phi(z_k) \in [\underline{\Phi}, \infty)$, implying Claim 2.

Claim 3 holds because, by Lemma 4.2,

$$\mathbb{E}\|(A_x^k, A_y^k)\| \leq p\mathbb{E}[\|z_k - z_{k-1}\| + \|z_{k-1} - z_{k-2}\|] + \mathbb{E}\Gamma_{k-1}.$$

We have that $\mathbb{E}\|z_k - z_{k-1}\| \rightarrow 0$ and $\mathbb{E}\Gamma_k \rightarrow 0$. This ensures that $\mathbb{E}\|(A_x^k, A_y^k)\| \rightarrow 0$.

To prove Claim 4, suppose $z^* = (x^*, y^*)$ is a limit point of the sequence $\{z_k\}_{k=0}^{\infty}$ (a limit point must exist because we suppose the sequence $\{z_k\}_{k=0}^{\infty}$ is bounded). This means there exists a subsequence z_{k_q} satisfying $\lim_{q \rightarrow \infty} z_{k_q} \rightarrow z^*$. Furthermore, by the variance-reduced property of $\tilde{\nabla}_x(x_{k_q}, y_{k_q})$, we have $\mathbb{E}\|\tilde{\nabla}_x(x_{k_q}, y_{k_q}) - \nabla_x F(x_{k_q}, y_{k_q})\|^2 \rightarrow 0$, which implies that there exists a subsequence of $\{z_{k_q}\}_{q \in \mathbb{N}}$ (call it $\{z_{k_q}\}_{q \in \mathcal{I}}$ for some index set $\mathcal{I} \subset \mathbb{N}$) such that $\tilde{\nabla}_x(x_{k_q}, y_{k_q}) - \nabla_x F(x_{k_q}, y_{k_q}) \rightarrow 0$ a.s. Because R and J are lower semicontinuous,

$$(4.8) \quad \liminf_{q \rightarrow \infty} R(x_{k_q}) \geq R(x^*) \quad \text{and} \quad \liminf_{q \rightarrow \infty} J(x_{k_q}) \geq J(x^*).$$

By the update rule for x_{k+1} ,

$$x_{k+1} \in \operatorname{argmin}_x \left\{ \langle x - x_k, \tilde{\nabla}_x(x_k, y_k) \rangle + \frac{1}{2\gamma_{x,k}} \|x - x_k\|^2 + R(x) \right\}.$$

Letting $x = x^*$,

$$\begin{aligned} & \langle x_{k+1} - x_k, \tilde{\nabla}_x(x_k, y_k) \rangle + \frac{1}{2\gamma_{x,k}} \|x_{k+1} - x_k\|^2 + R(x_{k+1}) \\ & \leq \langle x^* - x_k, \nabla_x F(x_k, y_k) \rangle + \langle x^* - x_k, \tilde{\nabla}_x(x_k, y_k) - \nabla_x F(x_k, y_k) \rangle + \frac{1}{2\gamma_{x,k}} \|x^* - x_k\|^2 + R(x^*). \end{aligned}$$

Setting $k = k_q$, taking the expectation, and taking the limit $q \rightarrow \infty$,

$$\begin{aligned} \limsup_{q \rightarrow \infty} R(x_{k_q+1}) & \leq \limsup_{q \rightarrow \infty} \langle x^* - x_{k_q}, \nabla_x F(x_{k_q}, y_{k_q}) \rangle \\ & \quad + \langle x^* - x_{k_q}, \tilde{\nabla}_x(x_{k_q}, y_{k_q}) - \nabla_x F(x_{k_q}, y_{k_q}) \rangle + \frac{1}{2\gamma_{x,k}} \|x^* - x_{k_q}\|^2 + R(x^*). \end{aligned}$$

The first term on the right goes to zero because $x_{k_q} \rightarrow x^*$ and $\nabla_x F(x_{k_q}, y_{k_q})$ is bounded. The second term is zero almost surely because it is bounded above by $\|x_{k_q} - x^*\|^2 + \|\tilde{\nabla}_x(x_{k_q}, y_{k_q}) - \nabla_x F(x_{k_q}, y_{k_q})\|^2$, and we have $\tilde{\nabla}_x(x_{k_q}, y_{k_q}) - \nabla_x F(x_{k_q}, y_{k_q}) \rightarrow 0$ a.s. Therefore, $\limsup_{q \rightarrow \infty} R(x_{k_q+1}) \leq R(x^*)$ a.s., which, together with equation (4.8), implies $R(x_{k_q+1}) \rightarrow R(x^*)$ a.s. The same argument holds for J and y_k , and it follows that

$$(4.9) \quad \lim_{q \rightarrow \infty} \Phi(x_{k_q}, y_{k_q}) = \Phi(x^*, y^*) \quad \text{a.s.}$$

Claim 3 ensures that $\mathbb{E}\|(A_x^k, A_y^k)\| \rightarrow 0$. Combining Claim 3 with (4.9) and the fact that the subdifferential of Φ is closed, we have $\mathbb{E}\text{dist}(0, \partial\Phi(z^*)) = 0$.

Claims 5 and 6 hold for any sequence satisfying $\|z_k - z_{k-1}\| \rightarrow 0$ a.s. (this fact is used in the same context in [6, Remark 5] and [14, Remark 4.1]).

Finally, we must show that Φ has constant expectation over ω . From Claim 2, we have $\mathbb{E}\Phi(z_k) \rightarrow \Phi^*$ which implies $\mathbb{E}\Phi(z_{k_q}) \rightarrow \Phi^*$ for every subsequence $\{z_{k_q}\}_{q=0}^\infty$ converging to some $z^* \in \omega$. In the proof of Claim 4, we show that $\Phi(z_{k_q}) \rightarrow \Phi(z^*)$, so $\mathbb{E}\Phi(z^*) = \Phi^*$ for all $z^* \in \omega$. ■

The following lemma is analogous to the Uniformized Kurdyka–Łojasiewicz Property [6]. It is a slight generalization of the Kurdyka–Łojasiewicz property showing that z_k eventually enters a region of \bar{z} for some \bar{z} satisfying $\Phi(\bar{z}) = \Phi(z^*)$, and in this region, the Kurdyka–Łojasiewicz inequality holds.

Lemma 4.5. *Assume the conditions of Lemma 4.3 hold and that z_k is not a critical point of Φ after a finite number of iterations. Let Φ be a semi-algebraic function with KL exponent θ . Then there exists an index m and a desingularizing function ϕ so that the following bound holds:*

$$\phi'(\mathbb{E}[\Phi(z_k) - \Phi_k^*])\mathbb{E}\text{dist}(0, \partial\Phi(z_k)) \geq 1 \quad \forall k > m,$$

where Φ_k^* is a non-decreasing sequence converging to $\mathbb{E}\Phi(z^*)$ for some $z^* \in \omega$.

Proof. First, we show that $\mathbb{E}\Phi(z_k)$ satisfies the KL property. Recall that b is the mini-batch size. Let $\bar{n} = \binom{n}{b}$ be the number of possible gradient estimates in one iteration, and let $\{z_k^i\}_{i=1}^{\bar{n}^k}$ be the set of possible values for z_k . Considering $\mathbb{E}\Phi$ as a function of $\{z_k^i\}_{i=1}^{\bar{n}^k}$, we have

$$\mathbb{E}\Phi(z_k) = \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i).$$

Because $\mathbb{E}\Phi(z_k)$ can be written as $\sum_i f_i(x_i)$ where f_i are KL functions with exponent θ , $\mathbb{E}\Phi(z_k)$ (as a function of $\{z_k^i\}_{i=1}^{\bar{n}^k}$) is also KL with exponent θ [25, Theorem 3.3]. Hence, $\mathbb{E}\Phi$ satisfies the KL inequality at every point in its domain. Therefore, for every point $(z_k^1, \dots, z_k^{\bar{n}^k})$ in a neighborhood U_k of $(\bar{z}_k^1, \bar{z}_k^2, \dots, \bar{z}_k^{\bar{n}^k})$ and satisfying

$$(4.10) \quad \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(\bar{z}_k^i) < \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) < \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(\bar{z}_k^i) + \epsilon_k$$

for some $\epsilon_k > 0$, the Kurdyka–Łojasiewicz inequality holds with the desingularizing function ϕ_k :

$$(4.11) \quad \phi_k' \left(\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) - \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(\bar{z}_k^i) \right) \text{dist} \left(0, \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \partial\Phi(z_k^i) \right) \geq 1.^6$$

⁶For the subdifferential terms we are taking the Minkowski sum: $\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \partial\Phi(z_k^i) = \{ \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \xi_i \mid \xi_i \in \partial\Phi(z_k^i) \}$.

There always exists a choice of $(\bar{z}_k^1, \bar{z}_k^2, \dots, \bar{z}_k^{\bar{n}^k})$ satisfying (4.10) unless $\mathbb{E}\Phi(z_k)$ is a local minimum. Lemma 4.3 Claim 5 implies $\text{dist}(z_k, \omega) \rightarrow 0$ a.s., and Claims 2 and 7 imply $\mathbb{E}\Phi(z_k) \rightarrow \mathbb{E}\Phi(z^*)$, so we can choose \bar{z}_k such that $\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(\bar{z}_k^i) \rightarrow \mathbb{E}\Phi(z^*)$ as well. To summarize, we have shown that there exists a sequence $(\bar{z}_k^1, \dots, \bar{z}_k^{\bar{n}^k})$ such that

1. The point $(z_k^1, \dots, z_k^{\bar{n}^k})$ lies in a neighborhood U_k of $(\bar{z}_k^1, \dots, \bar{z}_k^{\bar{n}^k})$,
2. The inequality (4.10) is satisfied, and
3. We have $\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(\bar{z}_k^i) \rightarrow \mathbb{E}\Phi(z^*)$.

Points 1.) and 2.) imply the Kurdyka–Łojasiewicz inequality (4.11). This ensures that the Kurdyka–Łojasiewicz inequality holds at every iteration, but the desingularizing function ϕ_k changes every iteration. We now show that the Kurdyka–Łojasiewicz inequality holds using a single function ϕ .

Because Φ is semi-algebraic with KL exponent θ , each desingularizing function is of the form $\phi_k(s) = a_k s^{1-\theta}$. Each a_k is bounded, so $a_{\max} \stackrel{\text{def}}{=} \max\{a_k\}_{k \geq 1}$ is bounded, and inequality (4.11) holds with the desingularizing function $\phi_{\max}(s) = a_{\max} s^{1-\theta}$.

Let $\Phi_k^* \stackrel{\text{def}}{=} \min_{j \geq k} \frac{1}{\bar{n}^j} \sum_{i=1}^{\bar{n}^j} \Phi(\bar{z}_j^i)$. It is clear that Φ_k^* is non-decreasing and $\Phi_k^* \rightarrow \mathbb{E}\Phi(z^*)$. From point 3, we can say there exists an index m and a constant a such that for all $k \geq m$,

$$(4.12) \quad a \left(\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) - \Phi_k^* \right)^{-\theta} \geq a_{\max} \left(\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) - \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(\bar{z}_k^i) \right)^{-\theta}.$$

The constant a exists; we can take a to be

$$(4.13) \quad \max_{k \geq 1} \left\{ \left(\frac{\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) - \Phi_k^*}{\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) - \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(\bar{z}_k^i)} \right)^{\theta} \right\}_{k \geq 1},$$

which is bounded. To see this, we acknowledge that this ratio is bounded for every k , and

$$(4.14) \quad \lim_{k \rightarrow \infty} \left(\frac{\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) - \Phi_k^*}{\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) - \frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(\bar{z}_k^i)} \right) = \lim_{k \rightarrow \infty} \left(\frac{\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) - \mathbb{E}\Phi(z^*)}{\frac{1}{\bar{n}^k} \sum_{i=1}^{\bar{n}^k} \Phi(z_k^i) - \mathbb{E}\Phi(z^*)} \right) = 1.$$

Therefore, with $\phi(s) = a s^{1-\theta}$, we have

$$(500) \quad \phi'(\mathbb{E}[\Phi(z_k) - \Phi_k^*]) \text{dist}(0, \mathbb{E}\partial\Phi(z_k)) \geq \phi'_{\max}(\mathbb{E}[\Phi(z_k) - \Phi_k^*]) \text{dist}(0, \mathbb{E}\partial\Phi(z_k)) \geq 1, \forall k > m,$$

The desired inequality follows from Jensen's inequality and the convexity of $x \mapsto \text{dist}(0, x)$. ■

We now show that the iterates of SPRING have finite length in expectation.

Lemma 4.6 (Finite Length). *Suppose Φ is a semi-algebraic function with KL exponent $\theta \in [0, 1)$. Let $\{z_k\}_{k=0}^{\infty}$ be a bounded sequence of iterates of SPRING using a variance-reduced gradient estimator and step-sizes satisfying the hypotheses of Lemma 4.3.*

- (1). *Either z_k is a critical point after a finite number of iterations, or $\{z_k\}_{k=0}^{\infty}$ satisfies the finite length property in expectation:*

$$(508) \quad \sum_{k=0}^{\infty} \mathbb{E} \|z_{k+1} - z_k\| < \infty,$$

and there exists an iteration m so that for all $i > m$,

$$\sum_{k=m}^i \mathbb{E} \|z_{k+1} - z_k\| + \mathbb{E} \|z_k - z_{k-1}\| \leq \sqrt{\mathbb{E} \|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E} \|z_{m-1} - z_{m-2}\|^2} \\ + \frac{2\sqrt{s}}{K_1\rho} \sqrt{\mathbb{E}\Upsilon_{m-1}} + K_3\Delta_{m,i+1},$$

where

$$K_1 \stackrel{\text{def}}{=} p + 2\sqrt{sV_\Upsilon}/\rho, \quad K_2 \stackrel{\text{def}}{=} \frac{1}{2\gamma_0} - \frac{\bar{L}}{2} - \frac{3\sqrt{2}}{4} \sqrt{V_1 + V_\Upsilon/\rho}, \quad K_3 \stackrel{\text{def}}{=} \frac{2K_1(K_2 + Z)}{K_2Z},$$

p is as in Lemma 4.2, and $\Delta_{p,q} \stackrel{\text{def}}{=} \phi(\mathbb{E}[\Psi_p - \Phi_p^*]) - \phi(\mathbb{E}[\Psi_q - \Phi_q^*])$.

(2). The iterates of SPRING $\{z_k\}_{k=0}^\infty$ converge to a critical point of Φ in expectation.

Remark 4.7. Our analysis for SPRING requires Φ to be semi-algebraic for the finite-length property to hold, but in the analysis of PALM, the finite-length property requires only that Φ is KL (and not necessarily semi-algebraic) [6, Thm. 1]. This difference arises because SPRING does not necessarily decrease the objective every iteration (even in expectation), but PALM does [6, Lem. 3]. Instead, we prove that the iterates of SPRING decrease Ψ_k in expectation. Related works [14] solve this problem by requiring an analog of Ψ_k to be KL, but this is not a straightforward approach for SPRING because of the complex variance bounds required to analyze variance-reduced gradient estimators.

Proof. We begin with a proof of Claim 1. If $\theta \in (0, 1/2)$, then Φ satisfies the KL property with exponent $1/2$, so we consider only the case $\theta \in [1/2, 1)$. By Lemma 4.5, there exists a function $\phi_0(r) = ar^{1-\theta}$ such that

$$\phi'_0(\mathbb{E}[\Phi(z_k) - \Phi_k^*])\mathbb{E}\text{dist}(0, \partial\Phi(z_k)) \geq 1 \quad \forall k > m.$$

Lemma 4.2 provides a bound on $\mathbb{E}\text{dist}(0, \partial\Phi(z_k))$.

(4.15)

$$\mathbb{E}\text{dist}(0, \partial\Phi(z_k)) \leq \mathbb{E}\|(A_x^k, A_y^k)\| \leq p\mathbb{E}[\|z_k - z_{k-1}\| + \|z_{k-1} - z_{k-2}\|] + \mathbb{E}\Gamma_{k-1} \\ \leq p(\sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} + \sqrt{\mathbb{E}\|z_{k-1} - z_{k-2}\|^2}) + \sqrt{s\mathbb{E}\Upsilon_{k-1}}.$$

The final inequality is Jensen's. Because $\Gamma_k = \sum_{i=1}^s v_k^i$ for some non-negative random variables v_k^i , we can say $\mathbb{E}\Gamma_k = \mathbb{E}\sum_{i=1}^s v_k^i \leq \mathbb{E}\sqrt{s\sum_{i=1}^s (v_k^i)^2} \leq \sqrt{s\mathbb{E}\Upsilon_k}$. We can bound the term $\sqrt{\mathbb{E}\Upsilon_k}$ using (2.2):

$$\sqrt{\mathbb{E}\Upsilon_k} \leq \sqrt{(1-\rho)\mathbb{E}\Upsilon_{k-1} + V_\Upsilon\mathbb{E}[\|z_k - z_{k-1}\|^2 + \|z_{k-1} - z_{k-2}\|^2]} \\ \leq \sqrt{(1-\rho)}\sqrt{\mathbb{E}\Upsilon_{k-1}} + \sqrt{V_\Upsilon}(\sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} + \sqrt{\mathbb{E}\|z_{k-1} - z_{k-2}\|^2}) \\ \leq (1 - \frac{\rho}{2})\sqrt{\mathbb{E}\Upsilon_{k-1}} + \sqrt{V_\Upsilon}(\sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} + \sqrt{\mathbb{E}\|z_{k-1} - z_{k-2}\|^2}).$$

The final inequality uses the fact that $\sqrt{1-\rho} = 1 - \rho/2 - \rho^2/8 - \dots$. This allows us to say

(4.17)

$$\mathbb{E}\text{dist}(0, \partial\Phi(z_k)) \leq K_1\sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} + K_1\sqrt{\mathbb{E}\|z_{k-1} - z_{k-2}\|^2} + \frac{2\sqrt{s}}{\rho}(\sqrt{\mathbb{E}\Upsilon_{k-1}} - \sqrt{\mathbb{E}\Upsilon_k}),$$

where $K_1 \stackrel{\text{def}}{=} p + 2\sqrt{sV_\Upsilon}/\rho$. Define C_k to be the right side of this inequality:

$$C_k \stackrel{\text{def}}{=} K_1\sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} + K_1\sqrt{\mathbb{E}\|z_{k-1} - z_{k-2}\|^2} + \frac{2\sqrt{s}}{\rho}(\sqrt{\mathbb{E}\Upsilon_{k-1}} - \sqrt{\mathbb{E}\Upsilon_k}).$$

535 We then have

$$536 \quad (4.18) \quad \phi'_0(\mathbb{E}[\Phi(z_k) - \Phi_k^*])C_k \geq 1 \quad \forall k > m.$$

537 By the definition of ϕ_0 , this is equivalent to

$$538 \quad (4.19) \quad \frac{a(1-\theta)C_k}{(\mathbb{E}[\Phi(z_k) - \Phi_k^*])^\theta} \geq 1 \quad \forall k > m.$$

539 We would like the inequality above to hold for Ψ_k rather than $\Phi(z_k)$. Replacing $\mathbb{E}\Phi(z_k)$ with $\mathbb{E}\Psi_k$
 540 introduces a term of $\mathcal{O}((\mathbb{E}[\|z_k - z_{k-1}\|^2 + \Upsilon_k])^\theta)$ in the denominator. We show that inequality (4.19)
 541 still holds after this adjustment because these terms are small compared to C_k .

542 The quantity $C_k \geq c_1(\sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} + \sqrt{\mathbb{E}\|z_{k-1} - z_{k-2}\|^2} + \sqrt{\mathbb{E}\Upsilon_{k-1}})$ for some constant
 543 $c_1 > 0$, and because $\mathbb{E}\|z_k - z_{k-1}\|^2, \mathbb{E}\Upsilon_k \rightarrow 0$, and $\theta \geq 1/2$, there exists an index m and a constants
 544 $c_2, c_3 > 0$ such that

$$545 \quad \left(\mathbb{E} \left[\frac{1}{2\rho\sqrt{2(V_1+V_\Upsilon/\rho)}} \Upsilon_k + \frac{\sqrt{V_1+V_\Upsilon/\rho}}{\sqrt{2}} \|z_k - z_{k-1}\|^2 \right] \right)^\theta \\ \leq c_2 \left(\mathbb{E} [\Upsilon_{k-1} + \|z_k - z_{k-1}\|^2 + \|z_{k-1} - z_{k-2}\|^2] \right)^\theta \leq c_3 C_k \quad \forall k > m.$$

546 The first inequality uses (2.2). Because the terms above are small compared to C_k , there exists a constant
 547 $+\infty > d > c_3$ such that

$$548 \quad \frac{ad(1-\theta)C_k}{(\mathbb{E}[\Phi(z_k) - \Phi_k^*])^\theta + \left(\mathbb{E} \left[\frac{1}{2\rho\sqrt{2(V_1+V_\Upsilon/\rho)}} \Upsilon_k + \frac{\sqrt{V_1+V_\Upsilon/\rho}}{\sqrt{2}} \|z_k - z_{k-1}\|^2 \right] \right)^\theta} \geq 1,$$

549 for all $k > m$. Using the fact that $(a+b)^\theta \leq a^\theta + b^\theta$ for all $a, b \geq 0$ because $\theta \in [1/2, 1)$, we have

$$550 \quad \frac{ad(1-\theta)C_k}{(\mathbb{E}[\Psi_k - \Psi_k^*])^\theta} = \frac{ad(1-\theta)C_k}{\left(\mathbb{E} [\Phi(z_k) - \Phi_k^* + \frac{1}{2\rho\sqrt{2(V_1+V_\Upsilon/\rho)}} \Upsilon_k + \frac{\sqrt{V_1+V_\Upsilon/\rho}}{\sqrt{2}} \|z_k - z_{k-1}\|^2] \right)^\theta} \\ \geq \frac{ad(1-\theta)C_k}{(\mathbb{E}[\Phi(z_k) - \Phi_k^*])^\theta + \left(\mathbb{E} \left[\frac{1}{2\rho\sqrt{2(V_1+V_\Upsilon/\rho)}} \Upsilon_k + \frac{\sqrt{V_1+V_\Upsilon/\rho}}{\sqrt{2}} \|z_k - z_{k-1}\|^2 \right] \right)^\theta} \geq 1 \quad \forall k > m.$$

551 Therefore, with $\phi(r) = adr^{1-\theta}$,

$$552 \quad \phi'(\mathbb{E}[\Psi_k - \Phi_k^*])C_k \geq 1 \quad \forall k > m.$$

553 By the concavity of ϕ ,

$$554 \quad (4.20) \quad \phi(\mathbb{E}[\Psi_k - \Phi_k^*]) - \phi(\mathbb{E}[\Psi_{k+1} - \Phi_{k+1}^*]) \geq \phi'(\mathbb{E}[\Psi_k - \Phi_k^*])(\mathbb{E}[\Psi_k - \Phi_k^* + \Phi_{k+1}^* - \Psi_{k+1}]) \\ \geq \phi'(\mathbb{E}[\Psi_k - \Phi_k^*])(\mathbb{E}[\Psi_k - \Psi_{k+1}]),$$

555 where the last inequality follows from the fact that Φ_k^* is non-decreasing. With $\Delta_{p,q} \stackrel{\text{def}}{=} \phi(\mathbb{E}[\Psi_p -$
 556 $\Phi_p^*]) - \phi(\mathbb{E}[\Psi_q - \Phi_q^*])$, we have shown

$$557 \quad \Delta_{k,k+1}C_k \geq \mathbb{E}[\Psi_k - \Psi_{k+1}].$$

Using Lemma 4.1, we can bound $\mathbb{E}[\Psi_k - \Psi_{k+1}]$ below by both $\mathbb{E}\|z_{k+1} - z_k\|^2$ and $\mathbb{E}\|z_k - z_{k-1}\|^2$. Specifically,

$$(4.21) \quad \Delta_{k,k+1} C_k \geq Z \mathbb{E}[\|z_k - z_{k-1}\|^2],$$

as well as

$$(4.22) \quad \Delta_{k,k+1} C_k \geq K_2 \mathbb{E}[\|z_{k+1} - z_k\|^2],$$

where $K_2 \stackrel{\text{def}}{=} -\left(\frac{\bar{L}(\lambda+1)}{2} + \frac{V_1+V_T/\rho}{L\lambda} + Z - \frac{1}{2\gamma_0}\right)$ and λ and Z are set as in Lemma 4.1. Let us use the first of these inequalities to begin. Applying Young's inequality to (4.21) yields

$$(4.23) \quad 2\sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \leq 2\sqrt{C_k \Delta_{k,k+1} Z^{-1}} \leq \frac{C_k}{2K_1} + \frac{2K_1 \Delta_{k,k+1}}{Z}$$

Summing inequality (4.23) from $k = m$ to $k = i$,

$$(4.24) \quad \begin{aligned} 2 \sum_{k=m}^i \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} &\leq \sum_{k=m}^i \frac{C_k}{2K_1} + \frac{2K_1 \Delta_{m,i+1}}{Z} \\ &\leq \sum_{k=m}^i \frac{1}{2} \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} + \frac{1}{2} \sqrt{\mathbb{E}\|z_{k-1} - z_{k-2}\|^2} \\ &\quad - \frac{\sqrt{s}}{K_1 \rho} (\sqrt{\mathbb{E}\Upsilon_i} - \sqrt{\mathbb{E}\Upsilon_{m-1}}) + \frac{2K_1 \Delta_{m,i+1}}{Z}, \end{aligned}$$

Dropping the non-positive term $-\sqrt{\mathbb{E}\Upsilon_i}$, this shows that

$$\sum_{k=m}^i \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \leq \frac{1}{2} \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2} + \frac{\sqrt{s}}{K_1 \rho} \sqrt{\mathbb{E}\Upsilon_{m-1}} + \frac{2K_1 \Delta_{m,i+1}}{Z}.$$

Applying the same argument using inequality (4.22) instead of (4.21), we obtain

$$\begin{aligned} \sum_{k=m}^i \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} \\ \leq \frac{1}{2} \sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \frac{1}{2} \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2} + \frac{\sqrt{s}}{K_1 \rho} \sqrt{\mathbb{E}\Upsilon_{m-1}} + \frac{2K_1 \Delta_{m,i+1}}{K_2}. \end{aligned}$$

Adding these inequalities together, we have

$$\begin{aligned} \sum_{k=m}^i \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} &\leq \frac{1}{2} \sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2} \\ &\quad + \frac{2\sqrt{s}}{K_1 \rho} \sqrt{\mathbb{E}\Upsilon_{m-1}} + \frac{2K_1(K_2+Z)\Delta_{m,i+1}}{K_2 Z}. \end{aligned}$$

For easier analysis, we add $\frac{1}{2} \sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2}$ to the right side:

$$(4.25) \quad \begin{aligned} \sum_{k=m}^i \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \\ \leq \sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2} + \frac{2\sqrt{s}}{K_1 \rho} \sqrt{\mathbb{E}\Upsilon_{m-1}} + \frac{2K_1(K_2+Z)\Delta_{m,i+1}}{K_2 Z}. \end{aligned}$$

Applying Jensen's inequality to the terms on the left gives

$$\begin{aligned} \sum_{k=m}^i \mathbb{E}\|z_{k+1} - z_k\| + \mathbb{E}\|z_k - z_{k-1}\| \\ \leq \sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2} + \frac{2\sqrt{s}}{K_1 \rho} \sqrt{\mathbb{E}\Upsilon_{m-1}} + \frac{2K_1(K_2+Z)\Delta_{m,i+1}}{K_2 Z}. \end{aligned}$$

578 The term $\lim_{i \rightarrow \infty} \Delta_{m,i+1}$ is bounded because $\mathbb{E}\Psi_k$ is bounded due to Lemma 4.1, so letting $i \rightarrow \infty$
 579 proves the assertion.

580 An immediate consequence of Claim 1 is that the sequence $\mathbb{E}\|z_{k+1} - z_k\|$ is Cauchy, so the sequence
 581 $\{z_k\}_{k=0}^\infty$ converges in expectation to a critical point. This is because, for any $p, q \in \mathbb{N}$ with $p \geq q$,
 582 $\mathbb{E}\|z_p - z_q\| = \mathbb{E}\|\sum_{k=q}^{p-1} z_{k+1} - z_k\| \leq \sum_{k=q}^{p-1} \mathbb{E}\|z_{k+1} - z_k\|$, and the finite length property implies this
 583 final sum converges to zero. This proves Claim 2. \blacksquare

584 Finally, we prove convergence rates for SPRING depending on the KL exponent of the objective
 585 function, demonstrating that the full convergence theory of PALM extends to SPRING.

586 **Theorem 4.8 (Convergence Rates).** *Suppose Φ is a semi-algebraic function with KL exponent*
 587 *$\theta \in [0, 1)$. Let $\{z_k\}_{k=0}^\infty$ be a bounded sequence of iterates of SPRING using a variance-reduced gradient*
 588 *estimator and step-sizes satisfying the hypotheses of Lemma 4.3. The following convergence rates hold:*

- 589 (1). *If $\theta \in (0, 1/2]$, then there exists $d_1 > 0$ and $\tau \in [1 - \rho, 1)$ such that $\mathbb{E}\|z_k - z^*\| \leq d_1 \tau^k$.*
 590 (2). *If $\theta \in (1/2, 1)$, then there exists a constant $d_2 > 0$ such that $\mathbb{E}\|z_k - z^*\| \leq d_2 k^{-\frac{1-\theta}{2\theta-1}}$.*
 591 (3). *If $\theta = 0$, then there exists an $m \in \mathbb{N}$ such that $\mathbb{E}\Phi(z_k) = \mathbb{E}\Phi(z^*)$ for all $k \geq m$.*

592 **Proof.** As in the proof of the previous lemma, if $\theta \in (0, 1/2)$, then Φ satisfies the KL property
 593 with exponent $1/2$, so we consider only the case $\theta \in [1/2, 1)$.

594 Substituting the desingularizing function $\phi(r) = ar^{1-\theta}$ into (4.25),

$$595 \quad (4.26) \quad \sum_{k=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \\ \leq \sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2} + \frac{2\sqrt{s}}{K_1\rho} \sqrt{\mathbb{E}\Upsilon_{m-1}} + aK_3(\mathbb{E}[\Psi_m - \Phi_m^*])^{1-\theta}.$$

596 Because $\Psi_m = \Phi(z_m) + \mathcal{O}(\|z_m - z_{m-1}\|^2 + \Upsilon_m)$, we can rewrite the final term as $\Phi(z_m) - \Phi_m^*$.

$$597 \quad (\mathbb{E}[\Psi_m - \Phi_m^*])^{1-\theta} = (\mathbb{E}[\Phi(z_m) - \Phi_m^* + \frac{1}{2L\lambda\rho}\Upsilon_m + \frac{V_1+V_\Upsilon/\rho}{2L\lambda}\|z_m - z_{m-1}\|^2])^{1-\theta} \\ \stackrel{\textcircled{1}}{\leq} (\mathbb{E}[\Phi(z_m) - \Phi_m^*])^{1-\theta} + \left(\frac{1}{2L\lambda\rho}\mathbb{E}\Upsilon_m\right)^{1-\theta} + \left(\frac{V_1+V_\Upsilon/\rho}{2L\lambda}\mathbb{E}\|z_m - z_{m-1}\|^2\right)^{1-\theta}.$$

598 Inequality $\textcircled{1}$ is due to the fact that $(a + b)^{1-\theta} \leq a^{1-\theta} + b^{1-\theta}$. This yields the inequality

$$599 \quad \sum_{k=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \\ \leq \sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2} + \frac{2\sqrt{s}}{K_1\rho} \sqrt{\mathbb{E}\Upsilon_{m-1}} + aK_3(\mathbb{E}[\Phi(z_m) - \Phi_m^*])^{1-\theta} \\ + aK_3\left(\frac{1}{2L\lambda\rho}\mathbb{E}\Upsilon_m\right)^{1-\theta} + aK_3\left(\frac{V_1+V_\Upsilon/\rho}{2L\lambda}\mathbb{E}\|z_m - z_{m-1}\|^2\right)^{1-\theta}.$$

600 Applying the Kurdyka–Łojasiewicz inequality (2.4),

$$601 \quad (4.27) \quad aK_3(\mathbb{E}[\Phi(z_m) - \Phi_m^*])^{1-\theta} \leq aK_3(\mathbb{E}\|\zeta_m\|)^{\frac{1-\theta}{\theta}},$$

602 for all $\zeta_m \in \partial\Phi(z_m)$ and we have absorbed the constant C into a . Equation (4.15) provides a bound on
 603 the norm of the subgradient:

$$604 \quad (\mathbb{E}\|\zeta_m\|)^{\frac{1-\theta}{\theta}} \leq \left(p(\sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2}) + \sqrt{s\mathbb{E}\Upsilon_{m-1}}\right)^{\frac{1-\theta}{\theta}}.$$

605 Denote the right side of this inequality $\Theta_m^{\frac{1-\theta}{\theta}}$. Therefore,

$$\begin{aligned}
 & \sum_{k=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \\
 606 \quad (4.28) \quad & \leq \sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2} + \frac{2\sqrt{s}}{K_1\rho} \sqrt{\mathbb{E}\Upsilon_{m-1}} + aK_3\Theta_m^{\frac{1-\theta}{\theta}} \\
 & \quad + aK_3\left(\frac{1}{2L\lambda\rho}\mathbb{E}\Upsilon_m\right)^{1-\theta} + aK_3\left(\frac{V_1+V_{\Upsilon}/\rho}{2L\lambda}\mathbb{E}\|z_m - z_{m-1}\|^2\right)^{1-\theta}.
 \end{aligned}$$

607 Suppose $\theta \in (1/2, 1)$. Each of the terms on the right side of this inequality are converging to zero, but
 608 at different rates. Because $\Theta_m = \mathcal{O}(\sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2} + \sqrt{\mathbb{E}\Upsilon_{m-1}})$, and θ
 609 satisfies $\frac{1-\theta}{\theta} < 1$, the term $\Theta_m^{\frac{1-\theta}{\theta}}$ dominates the first three terms on the right side of this inequality for
 610 large m . Also, because $\frac{1-\theta}{2\theta} \leq 1 - \theta$, $\Theta_m^{\frac{1-\theta}{\theta}}$ dominates the final two terms as well. Combining these
 611 facts, there exists a natural number M_1 such that for all $m \geq M_1$,

$$612 \quad (4.29) \quad \left(\sum_{k=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \right)^{\frac{\theta}{1-\theta}} \leq P\Theta_m,$$

613 for some constant $P > (aK_3)^{\frac{\theta}{1-\theta}}$. The bound of (4.16) implies

$$614 \quad 2\sqrt{s\mathbb{E}\Upsilon_{m-1}} \leq \frac{4\sqrt{s}}{\rho} (\sqrt{\mathbb{E}\Upsilon_{m-1}} - \sqrt{\mathbb{E}\Upsilon_m} + \sqrt{V_{\Upsilon}}(\sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2})).$$

615 Therefore,

$$\begin{aligned}
 \Theta_m &= p(\sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2}) + (2\sqrt{s\mathbb{E}\Upsilon_{m-1}} - \sqrt{s\mathbb{E}\Upsilon_m}) \\
 616 \quad (4.30) \quad & \leq \left(p + \frac{4\sqrt{sV_{\Upsilon}}}{\rho}\right) (\sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2}) \\
 & \quad + \frac{4\sqrt{s}}{\rho} (\sqrt{\mathbb{E}\Upsilon_{m-1}} - \sqrt{\mathbb{E}\Upsilon_m}) - \sqrt{s\mathbb{E}\Upsilon_{m-1}}.
 \end{aligned}$$

617 Furthermore, because $\frac{\theta}{1-\theta} > 1$ and $\mathbb{E}\Upsilon_m \rightarrow 0$, for large enough m , we have $(\sqrt{\mathbb{E}\Upsilon_m})^{\frac{\theta}{1-\theta}} \ll \sqrt{\mathbb{E}\Upsilon_m}$.
 618 This ensures that there exists a natural number M_2 such that for every $m \geq M_2$,

$$619 \quad (4.31) \quad \left(\frac{4\sqrt{s}(1-\rho/4)}{\rho(p+4\sqrt{sV_{\Upsilon}}/\rho)} \sqrt{\mathbb{E}\Upsilon_m} \right)^{\frac{\theta}{1-\theta}} \leq P\sqrt{s\mathbb{E}\Upsilon_m}.$$

620 The constant appearing on the left was chosen to simplify later arguments. Therefore, (4.29) implies

$$\begin{aligned}
 & \left(\sum_{k=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} + \frac{4\sqrt{s}(1-\rho/4)}{\rho(p+4\sqrt{sV_{\Upsilon}}/\rho)} \sqrt{\mathbb{E}\Upsilon_m} \right)^{\frac{\theta}{1-\theta}} \\
 & \stackrel{\textcircled{1}}{\leq} \frac{2^{\frac{\theta}{1-\theta}}}{2} \left(\sum_{k=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \right)^{\frac{\theta}{1-\theta}} \\
 & \quad + \frac{2^{\frac{\theta}{1-\theta}}}{2} \left(\frac{4\sqrt{s}(1-\rho/4)}{\rho(p+4\sqrt{sV_{\Upsilon}}/\rho)} \sqrt{\mathbb{E}\Upsilon_m} \right)^{\frac{\theta}{1-\theta}} \\
 621 \quad & \stackrel{\textcircled{2}}{\leq} \frac{2^{\frac{\theta}{1-\theta}}}{2} \left(\sum_{k=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \right)^{\frac{\theta}{1-\theta}} + \frac{2^{\frac{\theta}{1-\theta}}}{2} (P\sqrt{s\mathbb{E}\Upsilon_m}) \\
 & \stackrel{\textcircled{3}}{\leq} \frac{2^{\frac{\theta}{1-\theta}}}{2} \left(P(p+4\sqrt{sV_{\Upsilon}}/\rho) (\sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2}) \right. \\
 & \quad \left. + \frac{4\sqrt{s}P(1-\rho/4)}{\rho} (\sqrt{\mathbb{E}\Upsilon_{m-1}} - \sqrt{\mathbb{E}\Upsilon_m}) \right).
 \end{aligned}$$

Here, ① follows by convexity of the function $x^{\frac{\theta}{1-\theta}}$ for $\theta \in [1/2, 1)$ and $x \geq 0$, ② is (4.31), and ③ is (4.29) combined with (4.30). We absorb the constant $\frac{2^{\frac{\theta}{1-\theta}}}{2}$ into P . Define

$$S_m \stackrel{\text{def}}{=} \sum_{k=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} + \frac{4\sqrt{sP(1-\rho/4)}}{\rho(p+4\sqrt{sV_{\Upsilon}}/\rho)} \sqrt{\mathbb{E}\Upsilon_m}.$$

S_m is bounded for all m because $\sum_{k=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2}$ is bounded by equation (4.26). Hence, we have shown

$$(4.32) \quad S_m^{\frac{\theta}{1-\theta}} \leq P(p+4\sqrt{sV_{\Upsilon}}/\rho)(S_{m-1} - S_m).$$

The rest of the proof follows the proof of [2, Theorem 5]. Let $h(r) \stackrel{\text{def}}{=} r^{-\frac{\theta}{1-\theta}}$. First, suppose that $h(S_m) \leq Rh(S_{m-1})$ for some $R \in (1, \infty)$. Then (4.32) ensures that

$$\begin{aligned} 1 &\leq P(p+4\sqrt{sV_{\Upsilon}}/\rho)(S_{m-1} - S_m)h(S_m) \leq RP(p+4\sqrt{sV_{\Upsilon}}/\rho)(S_{m-1} - S_m)h(S_{m-1}) \\ &\leq RP(p+4\sqrt{sV_{\Upsilon}}/\rho) \int_{S_m}^{S_{m-1}} h(r)dr \\ &= \frac{RP(p+4\sqrt{sV_{\Upsilon}}/\rho)(1-\theta)}{1-2\theta} \left[S_{m-1}^{\frac{1-2\theta}{1-\theta}} - S_m^{\frac{1-2\theta}{1-\theta}} \right]. \end{aligned}$$

Hence,

$$0 < -\frac{1-2\theta}{RP(p+4\sqrt{sV_{\Upsilon}}/\rho)(1-\theta)} \leq S_m^{\frac{1-2\theta}{1-\theta}} - S_{m-1}^{\frac{1-2\theta}{1-\theta}}.$$

Now suppose $h(S_m) > Rh(S_{m-1})$, so that $S_m < R^{-\frac{1-\theta}{\theta}} S_{m-1}$ and $S_m^{\frac{1-2\theta}{1-\theta}} > q^{\frac{1-2\theta}{1-\theta}} S_{m-1}^{\frac{1-2\theta}{1-\theta}}$ where $q = R^{-\frac{1-\theta}{\theta}}$. This implies that

$$(q^{\frac{1-2\theta}{1-\theta}} - 1)S_{m-1}^{\frac{1-2\theta}{1-\theta}} \leq S_m^{\frac{1-2\theta}{1-\theta}} - S_{m-1}^{\frac{1-2\theta}{1-\theta}},$$

and the quantity on the left is clearly bounded away from zero because $q < 1$, $\frac{1-2\theta}{1-\theta} < 0$, and $S_{m-1} \rightarrow 0$. This shows that in either case, there exists a $\mu' > 0$ such that

$$\mu' \leq S_m^{\frac{1-2\theta}{1-\theta}} - S_{m-1}^{\frac{1-2\theta}{1-\theta}}.$$

Summing this inequality from $m = M_2$ to $m = M_3$, we obtain $(M_3 - M_2)\mu' \leq S_{M_3}^{\frac{1-2\theta}{1-\theta}} - S_{M_2-1}^{\frac{1-2\theta}{1-\theta}}$, and because the function $x \mapsto x^{\frac{1-\theta}{1-2\theta}}$ is decreasing, this implies

$$S_{M_3} \leq (S_{M_2-1}^{\frac{1-2\theta}{1-\theta}} + (M_3 - M_2)\mu')^{\frac{1-\theta}{1-2\theta}} \leq dM_3^{\frac{1-\theta}{1-2\theta}},$$

for some constant d . By Jensen's inequality, we can say $\sum_{k=M_3}^{\infty} \mathbb{E}\|z_k - z_{k-1}\| \leq S_{M_3} \leq dM_3^{-\frac{1-\theta}{2\theta-1}}$. Using the fact that $\mathbb{E}\|z_m - z^*\| = \mathbb{E}\|\sum_{k=m+1}^{\infty} z_k - z_{k-1}\| \leq \mathbb{E}\sum_{k=m}^{\infty} \|z_k - z_{k-1}\|$ proves Claim 1.

If $\theta = 1/2$, then $(\mathbb{E}\|\zeta_m\|)^{\frac{1-\theta}{\theta}} = \mathbb{E}\|\zeta_m\|$. Equation (4.28) then gives

$$\begin{aligned} & \sum_{i=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \\ & \leq \left(1 + aK_3(p + \sqrt{\frac{V_1 + V_Y/\rho}{2L\lambda}})\right) (\sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2}) \\ & \quad + \left(\frac{2\sqrt{s}}{K_1\rho} + aK_3\sqrt{s}\right) \sqrt{\mathbb{E}\Upsilon_{m-1}} + aK_3\sqrt{\frac{1}{2L\lambda\rho}} \sqrt{\mathbb{E}\Upsilon_m}, \end{aligned} \quad (4.33)$$

where we have added the non-negative term $aK_3\sqrt{\frac{V_1 + V_Y/\rho}{2L\lambda}}\sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2}$ to the right to simplify the presentation. Using equation (4.16), we have that, for any constant $c > 0$,

$$0 \leq -c\sqrt{\mathbb{E}\Upsilon_m} + c(1 - \frac{\rho}{2})\sqrt{\mathbb{E}\Upsilon_{m-1}} + c\sqrt{V_Y}(\sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2}).$$

Combining this inequality with (4.33),

$$\begin{aligned} & \sum_{i=m}^{\infty} \sqrt{\mathbb{E}\|z_{k+1} - z_k\|^2} + \sqrt{\mathbb{E}\|z_k - z_{k-1}\|^2} \\ & \leq \left(1 + aK_3(p + \sqrt{\frac{V_1 + V_Y/\rho}{2L\lambda}}) + c\sqrt{V_Y}\right) (\sqrt{\mathbb{E}\|z_m - z_{m-1}\|^2} + \sqrt{\mathbb{E}\|z_{m-1} - z_{m-2}\|^2}) \\ & \quad + c(1 - \frac{\rho}{2} + \frac{2\sqrt{s}}{cK_1\rho} + \frac{aK_3\sqrt{s}}{c})\sqrt{\mathbb{E}\Upsilon_{m-1}} - c(1 - aK_3c^{-1}\sqrt{\frac{1}{2L\lambda\rho}})\sqrt{\mathbb{E}\Upsilon_m}. \end{aligned}$$

Defining

$$T_m \stackrel{\text{def}}{=} \sum_{i=m}^{\infty} \sqrt{\mathbb{E}\|z_{i+1} - z_i\|^2} + \sqrt{\mathbb{E}\|z_i - z_{i-1}\|^2},$$

and $P_2 = 1 + aK_3(p + 4\sqrt{sV_Y}/\rho + \sqrt{\frac{V_1 + V_Y/\rho}{2L\lambda}}) + c\sqrt{V_Y}$, we have shown

$$\begin{aligned} & T_m + c(1 - aK_3c^{-1}\sqrt{\frac{1}{2L\lambda\rho}})\sqrt{\mathbb{E}\Upsilon_m} \\ & \leq P_2(T_{m-1} - T_m) + c(1 - \frac{\rho}{2} + \frac{2\sqrt{s}}{cK_1\rho} + \frac{aK_3\sqrt{s}}{c})\sqrt{\mathbb{E}\Upsilon_{m-1}}. \end{aligned}$$

Rearranging,

$$(1 + P_2)T_m + c(1 - aK_3c^{-1}\sqrt{\frac{1}{2L\lambda\rho}})\sqrt{\mathbb{E}\Upsilon_m} \leq P_2T_{m-1} + c(1 - \frac{\rho}{2} + \frac{2\sqrt{s}}{cK_1\rho} + \frac{aK_3\sqrt{s}}{c})\sqrt{\mathbb{E}\Upsilon_{m-1}}.$$

This implies

$$\begin{aligned} & T_m + \sqrt{\mathbb{E}\Upsilon_m} \\ & \leq \max\left\{\frac{P_2}{1 + P_2}, (1 - \frac{\rho}{2} + \frac{2\sqrt{s}}{cK_1\rho} + \frac{aK_3\sqrt{s}}{c})(1 - aK_3c^{-1}\sqrt{\frac{1}{2L\lambda\rho}})^{-1}\right\}(T_{m-1} + \sqrt{\mathbb{E}\Upsilon_{m-1}}). \end{aligned}$$

For large c , the second coefficient in the above expression approaches $1 - \rho/2$. This proves the linear rate of Claim 2.

When $\theta = 0$, the KL property (2.4) implies that exactly one of the following two scenarios holds: either $\mathbb{E}\Phi(z_k) \neq \Phi_k^*$ and

$$0 < C \leq \mathbb{E}\|\zeta_k\| \quad \forall \zeta_k \in \partial\Phi(z_k), \quad (4.34)$$

664 or $\Phi(z_k) = \Phi_k^*$. We show that the above inequality can only hold for a finite number of iterations.

665 Using the subgradient bound, the first scenario implies

$$\begin{aligned}
 C^2 &\leq (\mathbb{E}\|\zeta_k\|)^2 \leq (p\mathbb{E}\|z_k - z_{k-1}\| + p\mathbb{E}\|z_{k-1} - z_{k-2}\| + \mathbb{E}\Gamma_{k-1})^2, \\
 &\leq 3p^2(\mathbb{E}\|z_k - z_{k-1}\|)^2 + 3p^2(\mathbb{E}\|z_{k-1} - z_{k-2}\|)^2 + 3(\mathbb{E}\Gamma_{k-1})^2, \\
 &\leq 3p^2\mathbb{E}\|z_k - z_{k-1}\|^2 + 3p^2\mathbb{E}\|z_{k-1} - z_{k-2}\|^2 + 3s\mathbb{E}\Upsilon_{k-1}.
 \end{aligned}$$

667 where we have used the inequality $(a_1 + a_2 + \dots + a_s)^2 \leq s(a_1^2 + \dots + a_s^2)$ and Jensen's inequality.

668 Applying this inequality to the decrease of Ψ_k (4.2), we obtain

$$\begin{aligned}
 \mathbb{E}\Psi_k &\leq \mathbb{E}\Psi_{k-1} + \left(\frac{\bar{L}(\lambda+1)}{2} + \frac{V_1+V_\Gamma/\rho}{2L\lambda} + Z - \frac{1}{2\eta}\right)\mathbb{E}\|z_k - z_{k-1}\|^2 - Z\mathbb{E}\|z_{k-1} - z_{k-2}\|^2 \\
 &\leq \mathbb{E}\Psi_{k-1} - C^2 + \mathcal{O}(\mathbb{E}\|z_k - z_{k-1}\|^2) + \mathcal{O}(\mathbb{E}\|z_{k-1} - z_{k-2}\|^2) + \mathcal{O}(\mathbb{E}\Upsilon_{k-1}),
 \end{aligned}$$

670 for some constant C^2 .⁷ Because the final three terms go to zero as $k \rightarrow \infty$, there exists an index M_4 so
 671 that the sum of these three terms is bounded above by $C^2/2$ for all $k \geq M_4$. Therefore,

$$672 \quad \mathbb{E}\Psi_k \leq \mathbb{E}\Psi_{k-1} - \frac{C^2}{2}, \quad \forall k \geq M_4.$$

673 Because Ψ_k is bounded below for all k , this inequality can only hold for $N < \infty$ steps. After N steps, it
 674 is no longer possible for the bound (4.34) to hold, so it must be that $\Phi(z_k) = \Phi_k^*$. Because $\Phi_k^* < \Phi(z^*)$,
 675 $\Phi_k^* < \mathbb{E}\Phi(z_k)$, and both $\mathbb{E}\Phi(z_k)$, Φ_k^* converge to $\mathbb{E}\Phi(z^*)$, we must have $\Phi_k^* = \mathbb{E}\Phi(z_k) = \mathbb{E}\Phi(z^*)$. ■

676 The main difference between these convergence rates and those of PALM occurs when $\theta \in (0, 1/2]$.
 677 In this case, the linear convergence rate cannot be faster than the geometric decay of the MSE of the
 678 gradient estimator, which is of order $(1 - \rho)^k$ after k iterations. Without mini-batching (i.e. $b = 1$), this
 679 rate is approximately $(1 - 1/n)^k$ for the SAGA estimator and $(1 - 1/p)^k$ for the SARAH estimator.

680 **5. Numerical Experiments.** To demonstrate the advantages of SPRING, we compare SPRING
 681 using the SAGA and SARAH gradient estimators to PALM [6] and inertial PALM [29]. We also
 682 present results for SPRING using the (non-variance-reduced) SGD estimator (a case studied by Xu and
 683 Yin [39]). We refer to SPRING using the SGD, SAGA, and SARAH gradient estimators as SPRING-
 684 SGD, SPRING-SAGA, and SPRING-SARAH, respectively. Two applications are considered here for
 685 comparison: sparse non-negative matrix factorization (Sparse-NMF) and blind image-deblurring (BID)⁸.
 686 We also provide in the appendix additional results on sparse principal component analysis (Sparse-PCA).

687 **Sparse-NMF:** Given a data-matrix A , we seek a factorization $A \approx XY$ where $X \in \mathbb{R}^{n \times r}$, $Y \in$
 688 $\mathbb{R}^{r \times d}$ are non-negative with $r \leq d$ and X sparse. Sparse-NMF has the following formulation:

$$689 \quad (5.1) \quad \min_{X,Y} \|A - XY\|_F^2, \quad \text{s.t. } X, Y \geq 0, \quad \|X_i\|_0 \leq s, \quad i = 1, \dots, r.$$

690 Here, X_i denotes the i^{th} column of X . In dictionary learning and sparse coding, X is called the learned
 691 dictionary with coefficients Y . In this formulation, the sparsity on X is strictly enforced using the
 692 non-convex ℓ_0 constraint, restricting 75% of the entries to be 0.

⁷We have ignored extraneous constants in the final three terms for clarity.

⁸The implementations are available at <https://junqitang.com/>

Blind Image-Deblurring: Let Z be a blurred image. The problem of blind deconvolution reads:

$$(5.2) \quad \min_{X,Y} \|Z - X \odot Y\|_F^2 + \lambda \sum_{r=1}^{2d} \Phi([D(X)]_r) \quad \text{s.t.} \quad 0 \leq X \leq 1, 0 \leq Y \leq 1, \|Y\|_1 \leq 1,$$

where \odot is the 2D convolution operator, X is the image to recover, and Y is the blur-kernel to estimate. We choose a classic smooth edge-preserving regularizer in the image domain, with $D(\cdot)$ being the 2D differential operator computing the horizontal and vertical gradients for each pixel. For the potential function $\Phi(\cdot)$, we choose $\Phi(v) := \log(1 + \theta v^2)$ as in [29]. This potential function promotes sparsity in image gradients, hence yielding sharp images. We choose $\theta = 10^3$ and $\lambda = 5 \times 10^{-5}$ in our experiments

One of the benefits of SPRING and PALM is that the two step-sizes, $\gamma_{X,k}$ and $\gamma_{Y,k}$, depend separately on the Lipschitz constants $\hat{L}_X(Y_k)$ and $\hat{L}_Y(X_k)$. The practical performance of these algorithms depends significantly on the step-size choices. The following section describes how we use adaptive step-sizes in our experiments.

5.1. Parameter choices and on-the-fly estimation of Lipschitz constants. The global Lipschitz constants of the partial gradients of F are usually unknown and difficult to estimate. In practice, adaptive step-size choices based on estimating local Lipschitz constants are needed for PALM and inertial PALM [29]. In our experiments, we use the power method to estimate the Lipschitz constants on-the-fly in every iteration of the compared algorithms. For SPRING-SGD, SPRING-SAGA, and SPRING-SARAH, we find that it is sufficient to randomly sub-sample a mini-batch and run 5 iterations of the power method to get an estimate of the Lipschitz constants of the stochastic gradients. For PALM, we run 5 iterations of the power method in each iteration on the full batch to get an estimate of the Lipschitz constants of the full partial gradients.

For example, consider estimating the Lipschitz constants of the gradients corresponding to the objective function of Sparse-NMF (5.1). Let X_k and Y_k be the updates of k -th iteration, then $L_Y(X_k) = \|X_k\|^2$, which is the largest squared singular value of X_k , and can be computed via power iteration:

$$v_i = \frac{X_k^T(X_k v_{i-1})}{\|X_k^T(X_k v_{i-1})\|_2},$$

with a random initialization satisfying $\|v_0\|_2 = 1$. We find that using 5 iterations is sufficient to provide good estimates, so we approximate $L_Y(X_k)$ by $\|X_k^T(X_k v_5)\|_2$. We use the same strategy for $L_X(Y_k)$.

Denote the estimated Lipschitz constants of the full gradients as $\hat{L}_X(Y_k)$ and $\hat{L}_Y(X_k)$, and denote the estimated Lipschitz constants of the stochastic estimates as $\tilde{L}_X(Y_k)$ and $\tilde{L}_Y(X_k)$. We set the step-sizes of the compared algorithms as follows:

- **PALM:** $\gamma_{X,k} = \frac{1}{\hat{L}_X(Y_k)}$ and $\gamma_{Y,k} = \frac{1}{\hat{L}_Y(X_k)}$ (these are the standard step-sizes [6]).
- **Inertial PALM:** $\gamma_{X,k} = \frac{0.9}{\tilde{L}_X(Y_k)}$, $\gamma_{Y,k} = \frac{0.9}{\tilde{L}_Y(X_k)}$, and we set the momentum parameter to $\frac{k-1}{k+2}$, where k denotes the number of iterations. Pock and Sabach [29] assert that this dynamic momentum parameter achieves the best practical performance.⁹

⁹The dynamic choice of momentum parameter is not theoretically analyzed by Pock and Sabach [29], but it appears to be superior to the constant inertial parameter choice. Pock and Sabach suggest the aggressive step-sizes $\gamma_{X,k} = \frac{1}{\hat{L}_X(Y_k)}$ and $\gamma_{Y,k} = \frac{1}{\hat{L}_Y(X_k)}$ for the dynamic scheme, but we find these choices sometimes lead to unstable/divergent behavior in the late iterations. Hence, we use the slightly smaller step-sizes $\gamma_{X,k} = \frac{0.9}{\tilde{L}_X(Y_k)}$ and $\gamma_{Y,k} = \frac{0.9}{\tilde{L}_Y(X_k)}$ instead. These choices ensure the algorithm is stable, and we observe that they do not compromise the convergence rate in practice.

- **SPRING-SGD:** $\gamma_{X,k} = \frac{1}{\sqrt{\lceil kb/n \rceil} \tilde{L}_X(Y_k)}$ and $\gamma_{Y,k} = \frac{1}{\sqrt{\lceil kb/n \rceil} \tilde{L}_Y(X_k)}$. It is well-known in the literature that a shrinking step-size is necessary for SGD to converge to a critical point [7, 22, 26, 39].
- **SPRING-SAGA:** $\gamma_{X,k} = \frac{1}{3\tilde{L}_X(Y_k)}$ and $\gamma_{Y,k} = \frac{1}{3\tilde{L}_Y(X_k)}$.
- **SPRING-SARAH:** $\gamma_{X,k} = \frac{1}{2\tilde{L}_X(Y_k)}$ and $\gamma_{Y,k} = \frac{1}{2\tilde{L}_Y(X_k)}$.

Remark 5.1 (Practical step-sizes for SPRING-SAGA and SPRING-SARAH). While the step-sizes suggested in Sections 3 and 4 lead to state-of-the-art convergence rates for (1.1), we observe that those step-size choices are conservative for SPRING-SAGA and SPRING-SARAH in practice. Hence, we adopt the suggested step-size choices in the original works with scale factors 1/3 for SAGA [16, Section 2] and 1/2 for SARAH [27, Corollary 3]. For all tested methods, the step-sizes we use are optimal in practice while ensuring convergence in all experiments with extensive tests.

The same random initialization is used for all of the compared algorithms in our Sparse-NMF experiments, while for BID we initialize the image estimate with the blurred image and the kernel estimate with all ones. We observe that SPRING with variance-reduced gradients can be sensitive to poor initialization, and this may initially compromise convergence. However, this initialization issue can be effectively resolved if we use plain stochastic gradient without variance-reduction in the first epoch of SPRING-SARAH/SPRING-SAGA as a warm-start, which is suggested in [23].

In all the convergence plots for our experiments, we report the average results for stochastic methods with 10 independent runs. We comment here that from our numerical observations, the final objective values achieved by the stochastic algorithms vary very little from the average.

5.2. Sparse-NMF. We consider the extended Yale-B dataset and ORL dataset, which are standard facial recognition benchmarks consisting of human face images.¹⁰ The extended Yale-B dataset contains 2414 cropped images of size 32×32 , while the ORL dataset contains 400 images sized 64×64 . In the experiment for Yale dataset, we extract 49 sparse basis-images for the dataset. For ORL dataset we extract 25 sparse basis-images. In each iteration of the stochastic algorithms, we randomly sub-sample 5% of the full batch as a mini-batch. Here for SPRING-SARAH we set $p = \frac{1}{20}$. To reflect the effect of the algorithmic randomness within our methods, we report the average results (over 10 independent runs) of objective values in Figure 1. Meanwhile we also report the variance of the objective value at termination in Table 1. The obtained results are shown in Figures 1 and 2, from which we observe:

- Overall, SPRING using SAGA and SARAH estimators achieves superior performance compared to PALM, inertial PALM, and SPRING using the vanilla SGD gradient estimator.
- PALM has the worst performance in the considered Sparse-NMF tasks, which is not surprising since PALM is the baseline method in this comparison. Incorporating inertia can offer considerable acceleration for PALM. We believe that such inertial schemes can also be extended to accelerate SPRING and leave it as an important direction of future research (see [19] for some work in this direction).
- SPRING using the vanilla SGD gradient estimator achieves fast convergence initially, but gradually slows its convergence due to the shrinking step-size that is necessary to combat the non-reducing variance. However, using variance-reduced gradient estimators SAGA and SARAH, SPRING is able to overcome this issue and achieve the best overall convergence rates.

¹⁰Preprocessed versions [8, 9] can be found in: <http://www.cad.zju.edu.cn/home/dengcai/Data/FaceData.html>

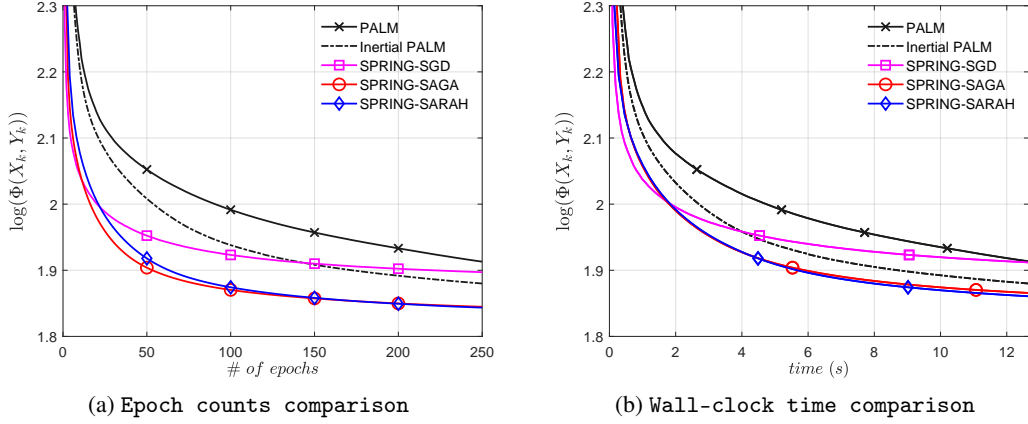


Figure 1: Objective decrease comparison of Sparse-NMF on Yale dataset.

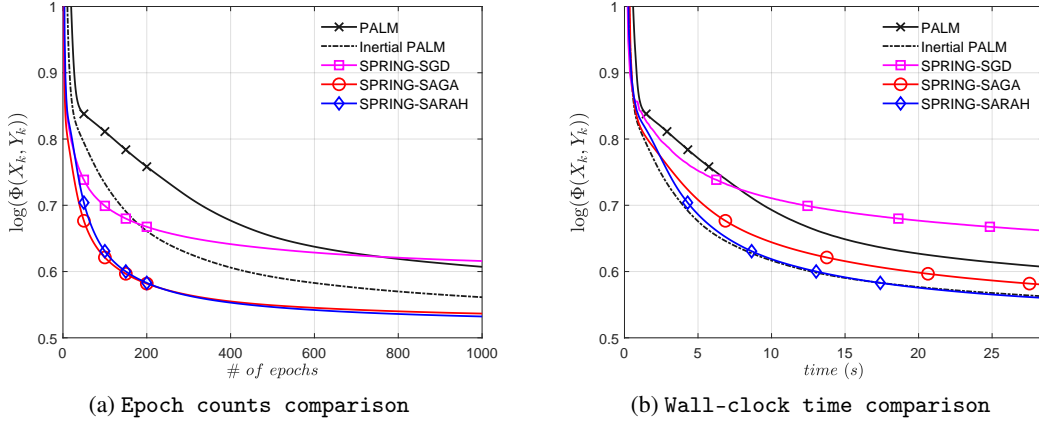


Figure 2: Objective decrease comparison of Sparse-NMF on ORL dataset.

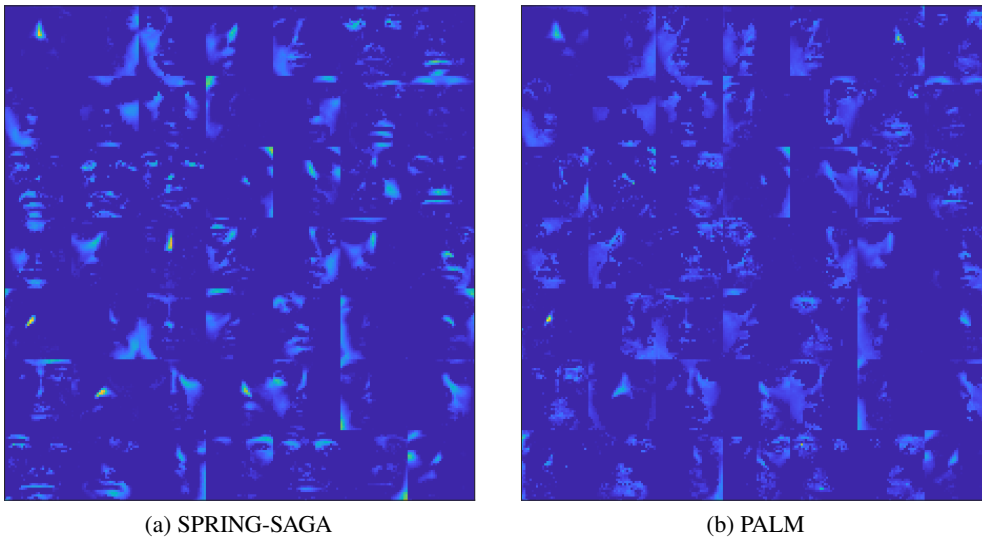
Remark 5.2 (Computational overheads for stochastic gradient methods). While the epoch count metric measures the gradient complexities of the algorithms, it does not reflect the computation overheads of the stochastic algorithms. The most important overhead for stochastic gradient methods in our setting would be the multiple calls to the proximal operator [35, 36]. Even though the proximal operators in our settings are not computationally expensive, computing such an operation many times still accumulates to a non-negligible overhead. Although our epoch count results confirm the complexity advantage predicted by theory, we can only observe compromised benefits from the wall-clock time comparison.

Remark 5.3 (The effect of algorithmic randomness). In order to reflect the algorithmic randomness of our stochastic methods, we present in log-scale the averaged convergence curves over 10 independent runs (in Figure 1 and 2). We also report that the variation of these results are virtually negligible, as we show in Table 1. The variances of the objective values at termination (250^{th} epoch for Yale dataset, and 1000^{th} epoch for ORL dataset) in the same log-scale are very small.

Table 1: The variation of the objective value (log-scale) at termination for randomized methods

DATASET/ALGORITHM	SPRING-SGD	SPRING-SAGA	SPRING-SARAH
YALE	1.8711×10^{-5}	7.5532×10^{-6}	8.3064×10^{-6}
ORL	9.9082×10^{-5}	1.6723×10^{-5}	1.2961×10^{-5}

As a visual illustration we also present in Figure 3 the basis images generated by SPRING-SAGA and PALM for the Yale dataset at the 50^{th} epoch. It is clear that the basis images generated by SPRING-SAGA appear natural and smooth quickly at an early epoch, while PALM’s results at the same epoch appear noisy and distorted.

Figure 3: Basis images from the Sparse-NMF experiment generated by SPRING-SAGA and PALM on the 50^{th} epoch for the Yale dataset.

5.3. Blind Image-Deblurring. For blind image-deconvolution, we choose to compare SPRING-SARAH, PALM and inertial PALM. We use two images, *Kodim08* and *Kodim15*, of size 256×256 for testing. For each image, two blur kernels—linear motion blur and out-of-focus blur—are considered with additional additive Gaussian noise. For SPRING-SARAH, the mini-batch size is $1/64$ of the full batch (and also we set $p = \frac{1}{64}$). For this mini-batch size, we choose smaller step sizes $\gamma_{X,k} = \frac{1}{8\bar{L}_X(Y_k)}$, $\gamma_{Y,k} = \frac{1}{3\bar{L}_Y(X_k)}$ than the default choices to encourage stability. As above, we present results of SPRING in terms of an average of 10 independent runs in Figures 6 and 7, and we report that the variance due to the algorithmic randomness evaluated at termination is also negligible (on the order of 10^{-6}).

For both images with motion blur, the convergence comparisons of the algorithms are provided in Figures 4 and 5, from which we observe SPRING-SARAH is faster than the other two methods in both cases. Figures 6 and 7 provide comparisons of the recovered image and blur kernel. We observe superior performance of SPRING-SARAH over PALM in these figures as well. In particular, when comparing

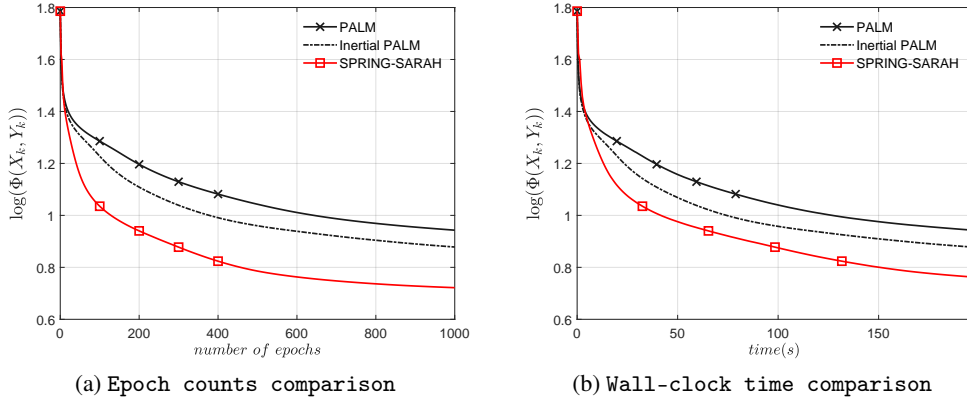


Figure 4: Objective decrease comparison (epoch counts) of blind image-deconvolution experiment on Kodim08 image using an 11×11 motion-blur kernel.

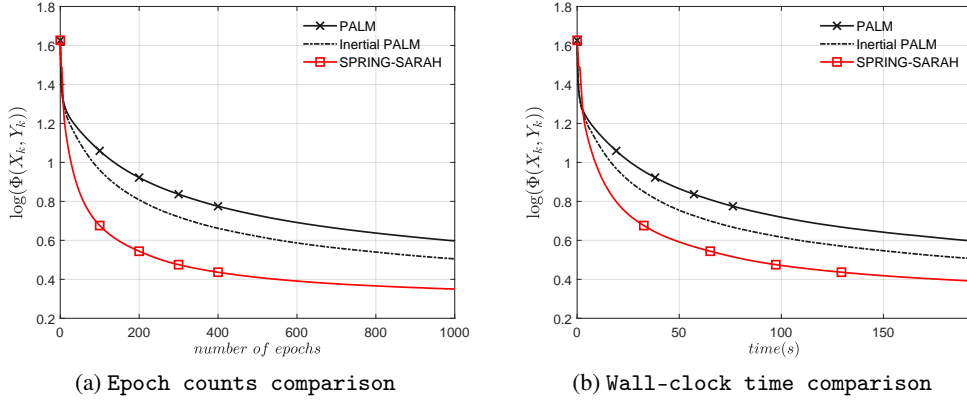


Figure 5: Objective decrease comparison (wall-clock time) of blind image-deconvolution experiment on Kodim15 image using an 11×11 motion-blur kernel.

the estimated blur kernels of the two algorithms every 100 epochs, we clearly see that SPRING-SARAH more quickly recovers more accurate solutions than PALM. It is worth noting that, although stochastic gradient methods have been shown to be inherently inefficient for non-blind and non-uniform deblurring task where the blur kernels are known or estimated beforehand [36], SPRING still offers significant acceleration over PALM in blind-deblurring tasks. Additional experiments using out-of-focus blur kernels are provided in the appendix.

6. Conclusion. We propose SPRING, a stochastic extension of the PALM algorithm for solving a class of structured non-smooth and non-convex optimization problems. We analyze the convergence properties of SPRING when using a variety of variance-reduced gradient estimators, and we prove specific convergence rates using the SAGA and SARAH estimators. For generic optimization problems of the form (1.1), we show that SPRING-SAGA (with $b \leq \mathcal{O}(n^{2/3})$) and SPRING-SARAH return an

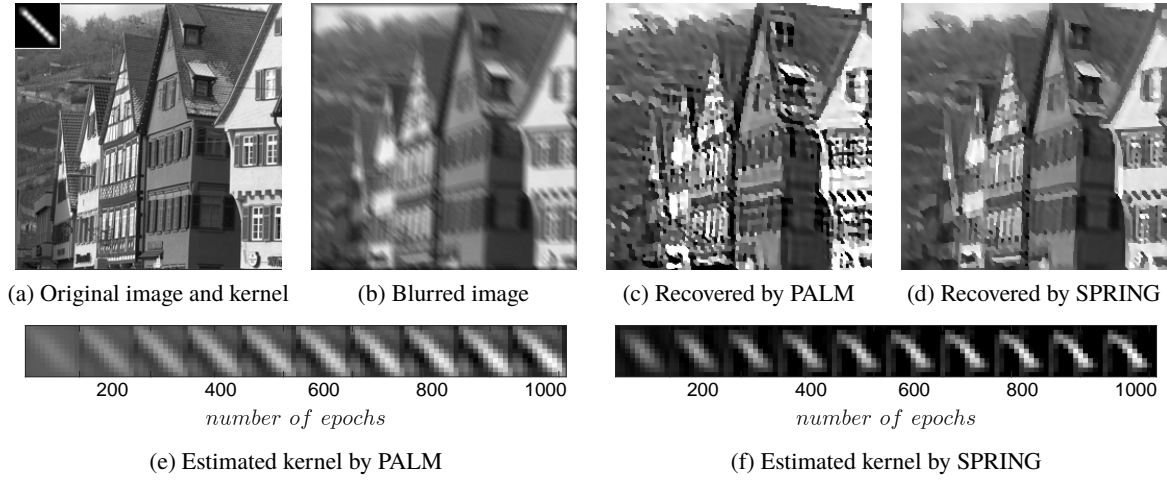


Figure 6: Image and kernel reconstructions from the blind image-deconvolution experiment on the Kodim08 image using an 11×11 motion blur kernel.

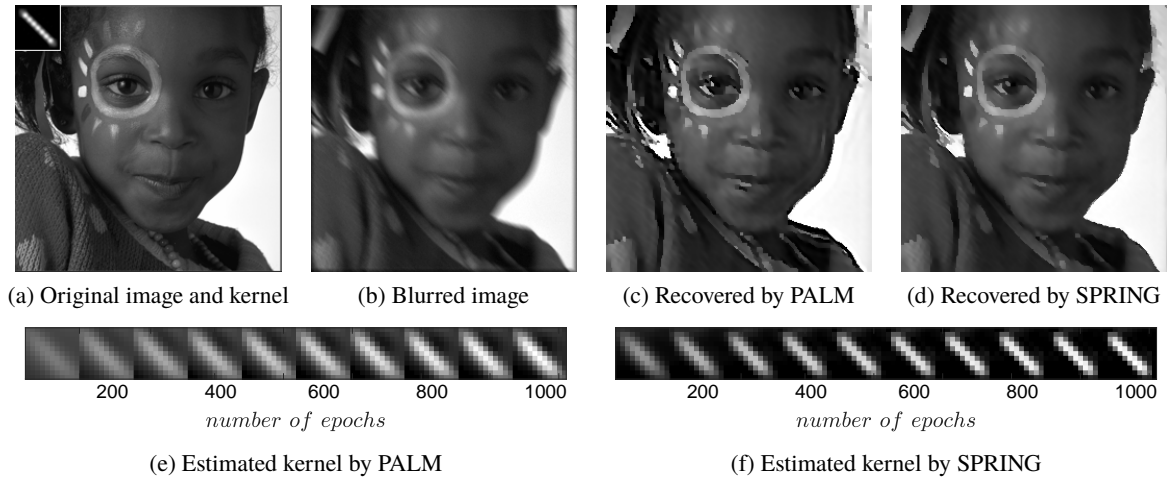


Figure 7: Image and kernel reconstructions from the blind image-deconvolution experiment on the Kodim15 image using an 11×11 motion blur kernel.

805 ϵ -approximate critical point in expectation in no more than $O(\frac{n^2 L}{b^3 \epsilon^2})$ and $O(\frac{\sqrt{n} L}{\epsilon^2})$ SFO calls, respectively,
 806 showing that SPRING-SARAH achieves the complexity lower bound for stochastic non-convex opti-
 807 mization. For objectives satisfying an error bound, we further demonstrate that our methods converge
 808 linearly to the global optimum. Because of the generality of our results, they contain almost all existing
 809 results for stochastic non-convex optimization as special cases, and they improve on them in many
 810 settings. Most importantly, we extend the full convergence theory of PALM to the stochastic setting,
 811 showing that SPRING achieves the same convergence rates as PALM on semialgebraic objectives.

812 **Acknowledgements.** JT and MD acknowledge support from the ERC Advanced grant, project
 813 694888, C-SENSE. JL acknowledges support from the Leverhulme Trust. CBS acknowledges support

from the Leverhulme Trust project on Breaking the Non-Convexity Barrier, and on Unveiling the Invisible, the Philip Leverhulme Prize, the EPSRC grant No. EP/S026045/1, EPSRC grant No. EP/M00483X/1, and EPSRC Centre No. EP/N014588/1, the European Union Horizon 2020 research and innovation programmes under the Marie Skłodowska-Curie grant agreement No. 691070 CHiPS and the Marie Skłodowska-Curie grant agreement No 777826, the Cantab Capital Institute for the Mathematics of Information, and the Alan Turing Institute.

REFERENCES

- [1] ARAVKIN, A., AND DAVIS, D. Trimmed statistical estimation via variance reduction. *Mathematics of Operations Research* (2019).
- [2] ATTOUCH, H., AND BOLTE, J. On the convergence of the proximal algorithm for nonsmooth functions involving analytic features. *Mathematical Programming B* 116, 1 (2007), 5–16.
- [3] ATTOUCH, H., BOLTE, J., REDONT, P., AND SOUBEYRAN, A. Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Lojasiewicz inequality. *Mathematics of Operations Research* 35, 2 (2010), 438–457.
- [4] BOLTE, J., DANIILIDIS, A., AND LEWIS, A. The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimization* 17, 4 (2007), 1205–1223.
- [5] BOLTE, J., DANIILIDIS, A., LEY, O., AND MAZET, L. Characterizations of Łojasiewicz inequalities: subgradient flows, talweg, convexity. *Transactions of the American Mathematical Society* 362, 6 (2010), 3319–3363.
- [6] BOLTE, J., SABACH, S., AND TEBoulLE, M. Proximal alternating linearised minimization for nonconvex and nonsmooth problems. *Mathematical Programming* 146, 1-2 (2014), 459–494.
- [7] BOTTOU, L. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT 2010*. Springer, 2010, pp. 177–186.
- [8] CAI, D., HE, X., AND HAN, J. Spectral regression for efficient regularized subspace learning. In *IEEE 11th International Conference on Computer Vision* (2007), IEEE, pp. 1–8.
- [9] CAI, D., HE, X., HU, Y., HAN, J., AND HUANG, T. Learning a spatially smooth subspace for face recognition. In *IEEE Conference on Computer Vision and Pattern Recognition* (2007), IEEE, pp. 1–7.
- [10] CAMPISI, P., AND EGIAZARIAN, K. *Blind image deconvolution: theory and applications*. CRC press, 2016.
- [11] CANDÈS, E. J., LI, X., MA, Y., AND WRIGHT, J. Robust principal component analysis? *Journal of the ACM (JACM)* 58, 3 (2011), 11.
- [12] CHAMBOLLE, A., EHRHARDT, M. J., RICHTÁRIK, P., AND SCHÖNLIEB, C.-B. Stochastic primal-dual hybrid gradient algorithm with arbitrary sampling and imaging applications. *SIAM J. Optim.* 28, 4 (2018), 2783–2808.
- [13] D’ASPREMONT, A., GHAOUI, L. E., JORDAN, M. I., AND LANCKRIET, G. R. A direct formulation for sparse pca using semidefinite programming. In *Advances in neural information processing systems* (2005), pp. 41–48.
- [14] DAVIS, D. The asynchronous palm algorithm for nonsmooth nonconvex problems. *arXiv:1604.00526* (2016).
- [15] DAVIS, D., EDMUNDS, B., AND UDELL, M. The sound of APALM clapping: Faster nonsmooth nonconvex optimization with stochastic asynchronous palm. In *Advances in Neural Information Processing Systems* (2016), pp. 226–234.
- [16] 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* (2014), pp. 1646–1654.
- [17] DRIGGS, D., TANG, J., LIANG, J., DAVIES, M., AND SCHÖNLIEB, C. Spring: A fast stochastic proximal alternating method for non-smooth non-convex optimization. *arXiv preprint arXiv:2002.12266* (2020).
- [18] FANG, C., LI, C. J., LIN, Z., AND ZHANG, T. Spider: Near-optimal non-convex optimization via stochastic path integrated differential estimator. In *32nd Conference on Neural Information Processing Systems* (2018).
- [19] HERTRICH, J., AND STEIDL, G. Inertial stochastic palm and its application for learning student-t mixture models. *arXiv:2005.02204* (2020).
- [20] HOYER, P. O. Non-negative matrix factorization with sparseness constraints. *Journal of machine learning research* 5, Nov (2004), 1457–1469.
- [21] JOHNSON, R., AND ZHANG, T. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems* (2013), pp. 315–323.
- [22] KONEČNÝ, J., LIU, J., RICHTÁRIK, P., AND TAKÁČ, M. Mini-batch semi-stochastic gradient descent in the proximal setting. *IEEE Journal of Selected Topics in Signal Processing* 10, 2 (2015), 242–255.

- [23] KONEČNÝ, J., AND RICHÁRIK, P. Semi-stochastic gradient descent methods. *Frontiers in Applied Mathematics and Statistics* 3 (2017), 9.
- [24] LI, B., MA, M., AND GIANNAKIS, G. B. On the convergence of sarah and beyond. *arXiv:1906.02351v2* (2020).
- [25] LI, G., AND PONG, T. K. Calculus of the exponent of kurdyka–Łojasiewicz inequality and its applications to linear convergence of first-order methods. *Foundations of Computational Mathematics* 18 (2018), 1199–1232.
- [26] MOULINES, E., AND BACH, F. R. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems* (2011), pp. 451–459.
- [27] NGUYEN, L. M., LIU, J., SCHEINBERG, K., AND TAKÁČ, M. SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *Proceedings of the 34th International Conference on Machine Learning* (2017), vol. 70, pp. 2613–2621.
- [28] PHAM, N. H., NGUYEN, L. M., PHAN, D. T., AND TRAN-DINH, Q. ProxSARAH: An efficient algorithmic framework for stochastic composite nonconvex optimization. *arXiv:1902.05679* (2019).
- [29] POCK, T., AND SABACH, S. Inertial proximal alternating linearized minimization (iPALM) for nonconvex and nonsmooth problems. *SIAM Journal on Imaging Sciences* 9, 4 (2016), 1756–1787.
- [30] REDDI, S. J., HEFNY, A., SRA, S., PÓCZOS, B., AND SMOLA, A. Stochastic variance reduction for nonconvex optimization. In *Proc. 33rd International Conference on Machine Learning* (2016).
- [31] REDDI, S. J., SRA, S., PÓCZOS, B., AND SMOLA, A. Fast stochastic methods for nonsmooth nonconvex optimization. In *Proc. 30th Annual Conference on Neural Information Processing Systems* (2016).
- [32] ROBBINS, H., AND MONRO, S. A stochastic approximation method. *Annals of Mathematical Statistics* 22, 3 (1951), 400–407.
- [33] ROBBINS, H., AND SIEGMUND, D. A convergence theorem for non-negative almost supermartingales and some applications. *Optimizing Methods in Statistics* (1971), 233–257.
- [34] SCHMIDT, M., ROUX, N. L., AND BACH, F. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming* 162 (2017), 83–112.
- [35] TANG, J., EGAZARIAN, K., AND DAVIES, M. The limitation and practical acceleration of stochastic gradient algorithms in inverse problems. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (2019), IEEE, pp. 7680–7684.
- [36] TANG, J., EGAZARIAN, K., GOLBABAEE, M., AND DAVIES, M. The practicality of stochastic optimization in imaging inverse problems. *IEEE Transactions on Computational Imaging* 6 (2020), 1471–1485.
- [37] WANG, Z., JI, K., ZHOU, Y., LIANG, Y., AND TAROKH, V. SpiderBoost: A class of faster variance-reduced algorithms for nonconvex optimization. *arXiv:1810.10690* (2018).
- [38] XIAO, L., AND ZHANG, T. A proximal stochastic gradient method with progressive variance reduction. *Technical report, Microsoft Research* (2014).
- [39] XU, Y., AND YIN, W. Block stochastic gradient iteration for convex and nonconvex optimization. *SIAM Journal on Optimization* 25, 3 (2015), 1686–1716.
- [40] ZHOU, D., AND GU, Q. Lower bounds for smooth nonconvex finite-sum optimization. *arXiv preprint arXiv:1901.11224* (2019).
- [41] ZHOU, Y., WANG, Z., JI, K., LIANG, Y., AND TAROKH, V. Momentum schemes with stochastic variance reduction for nonconvex composite optimization. *arXiv:1902.02715* (2019).
- [42] ZOU, H., HASTIE, T., AND TIBSHIRANI, R. Sparse principal component analysis. *Journal of computational and graphical statistics* 15, 2 (2006), 265–286.

Appendix A. Additional numerical experiments. This section contains additional numerical experiments demonstrating the superiority of SPRING over PALM.

We first present our additional results on the Sparse-PCA example for the Yale and ORL datasets. The problem of Sparse-PCA with r principal components can be written as:

$$(A.1) \quad \min_{X,Y} \|A - XY\|_F^2 + \lambda_1 \|X\|_1 + \lambda_2 \|Y\|_1.$$

where $X \in \mathbb{R}^{n \times r}$, $Y \in \mathbb{R}^{r \times d}$. We use ℓ_1 regularization on both X and Y to promote sparsity with $\lambda_1 = 10^{-3}$ and $\lambda_2 = 5 \times 10^{-3}$, and $r = 25$. We compare SPRING-SAGA, SPRING-SARAH, SPRING-SGD and PALM. We choose the mini-batch size to be $\frac{1}{40}$ of the full batch (for SPRING-SARAH we set $p = \frac{1}{40}$). We report the results of 10 independent runs of the stochastic methods in Figure 8 and 9, and we denote that the variance due to the algorithmic randomness evaluated at termination is also negligible (on the order of 10^{-5}). Similar to what we observe in the Sparse-NMF experiments, our results in Figure 8 and 9 show that SPRING with stochastic variance-reduced gradient estimators achieves the fastest convergence.

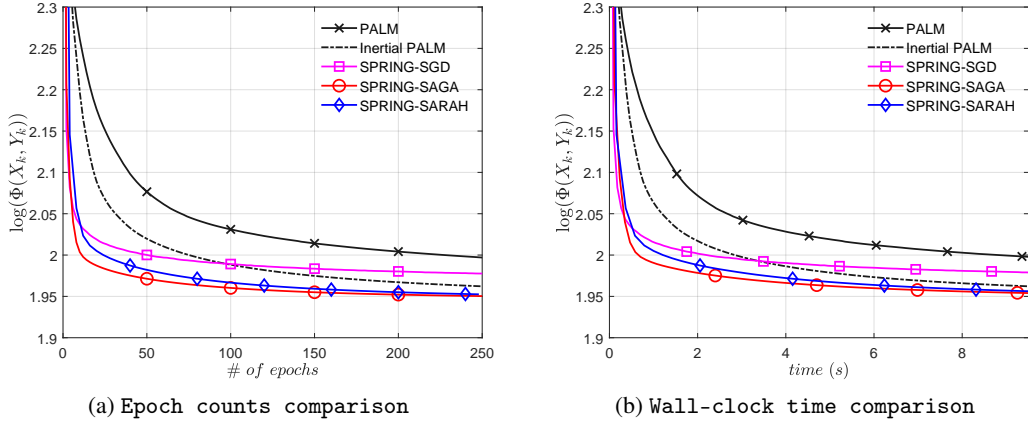


Figure 8: Objective decrease comparison of Sparse-PCA on Yale dataset.

Figures 10 to 12 show additional comparisons for blind image-deblurring where the images are blurred with an out-of-focus kernel. We choose the regularization parameter $\lambda = 1 \times 10^{-4}$ and the other settings are the same for the BID experiments presented in the main text. Again, we observe that our SPRING-SARAH algorithm outperforms PALM and inertial-PALM.

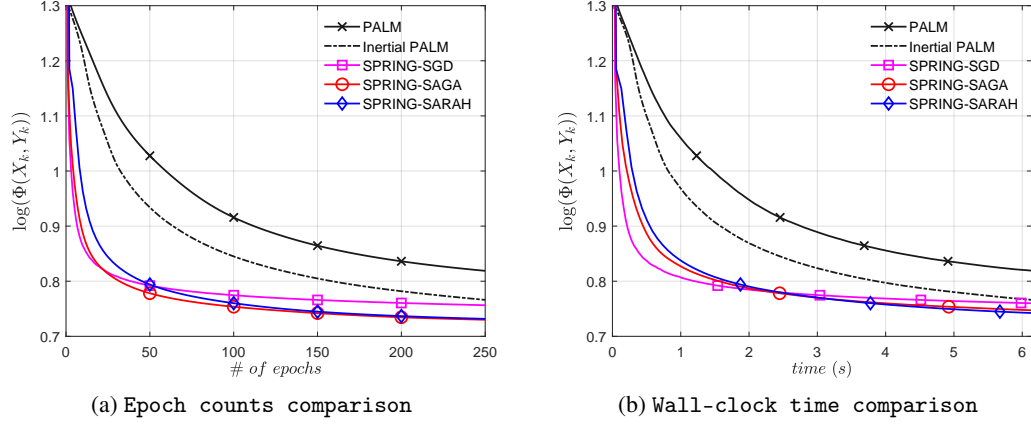


Figure 9: Objective decrease comparison of Sparse-PCA on ORL dataset.

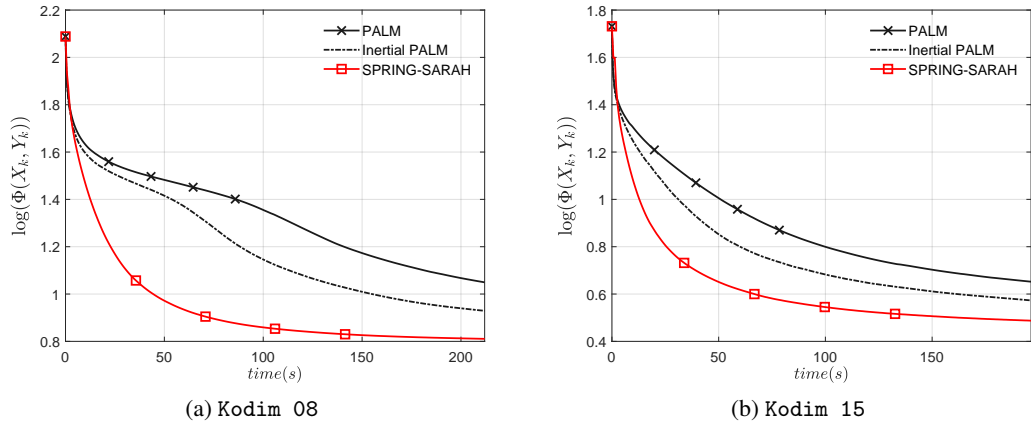


Figure 10: Objective decrease comparison (versus run time) of blind image-deconvolution experiment on Kodim08 and Kodim15 images using an out-of-focus blur kernel.

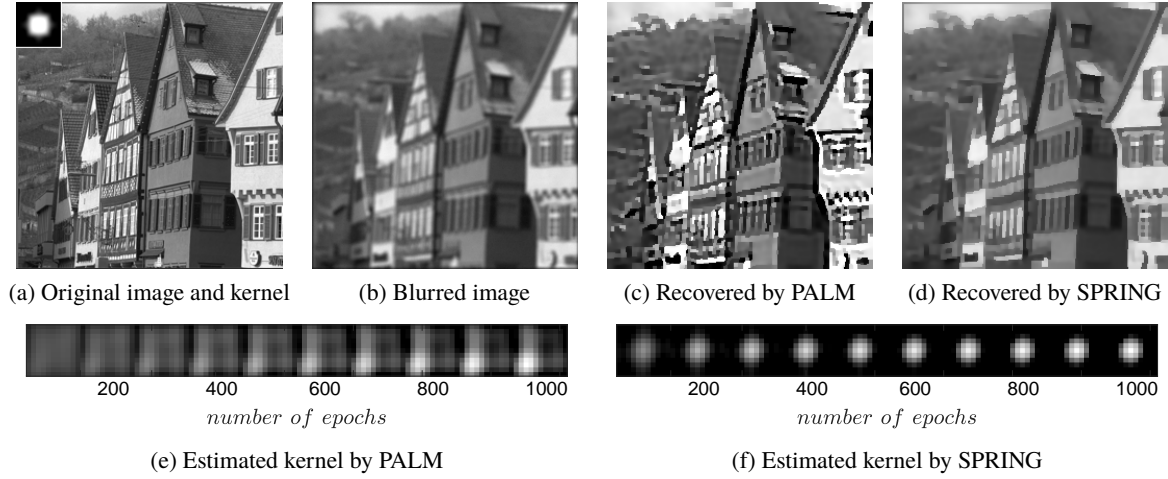


Figure 11: Image and kernel reconstructions from the blind image-deconvolution experiment on the Kodim08 image using an out-of-focus blur kernel.

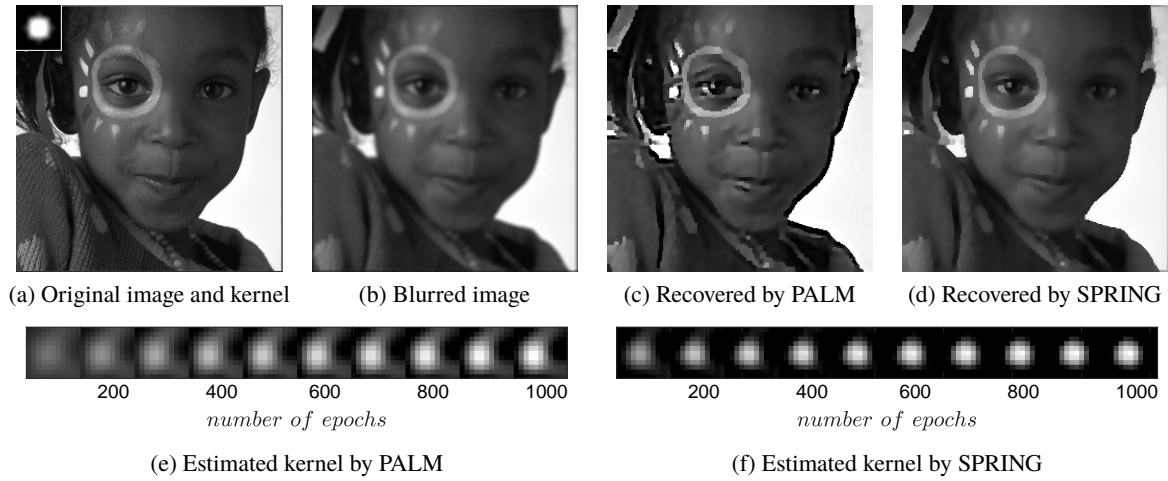


Figure 12: Image and kernel reconstructions from the blind image-deconvolution experiment on the Kodim15 image using an out-of-focus blur kernel.