# Revisiting Inexact Fixed-Point Iterations for Min-Max Problems: Stochasticity and Structured Nonconvexity

Ahmet Alacaoglu 12 Donghwan Kim 3 Stephen J. Wright 1

# **Abstract**

We focus on constrained, L-smooth, potentially stochastic and nonconvex-nonconcave min-max problems either satisfying  $\rho$ -cohypomonotonicity or admitting a solution to the  $\rho$ -weakly Minty Variational Inequality (MVI), where larger values of the parameter  $\rho > 0$  correspond to a greater degree of nonconvexity. These problem classes include examples in two player reinforcement learning, interaction dominant min-max problems, and certain synthetic test problems on which classical min-max algorithms fail. It has been conjectured that first-order methods can tolerate a value of  $\rho$ no larger than  $\frac{1}{L}$ , but existing results in the literature have stagnated at the tighter requirement  $\rho < \frac{1}{2L}$ . With a simple argument, we obtain optimal or best-known complexity guarantees with cohypomonotonicity or weak MVI conditions for  $\rho < \frac{1}{L}$ . First main insight for the improvements in the convergence analyses is to harness the recently proposed conic nonexpansiveness property of operators. Second, we provide a refined analysis for inexact Halpern iteration that relaxes the required inexactness level to improve some state-of-the-art complexity results even for constrained stochastic convex-concave min-max problems. Third, we analyze a stochastic inexact Krasnosel'skiĭ-Mann iteration with a multilevel Monte Carlo estimator when the assumptions only hold with respect to a solution.

# 1. Introduction

We consider the problem

$$\min_{u \in U} \max_{v \in V} f(u, v), \tag{1}$$

Proceedings of the 41<sup>st</sup> International Conference on Machine Learning, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

where  $U \subseteq \mathbb{R}^m$ ,  $V \subseteq \mathbb{R}^n$  are closed convex sets admitting efficient projection operators and  $f: \mathbb{R}^m \times \mathbb{R}^n \to \mathbb{R}$  is a function such that  $\nabla_u f(u,v)$  and  $\nabla_v f(u,v)$  are Lipschitz continuous. The general setting where f(u, v) is allowed to be nonconvex-nonconcave is extremely relevant in machine learning (ML), with applications in generative adversarial networks (GANs) (Goodfellow et al., 2014) and adversarial ML (Madry et al., 2018). Yet, at the same time, such problems are extremely challenging to solve, with documented hardness results, see e.g., (Daskalakis et al., 2021). As a result, an extensive literature has arisen about special cases of the nonconvex-nonconcave problem (1) for which algorithms with good convergence and complexity properties can be derived (Diakonikolas et al., 2021; Bauschke et al., 2021; Lee & Kim, 2021; Pethick et al., 2022; 2023a;b; Gorbunov et al., 2023; Böhm, 2022; Cai et al., 2022b; Cai & Zheng, 2022; Hajizadeh et al., 2023; Kohlenbach, 2022; Lee & Kim, 2024; Fan et al., 2024; Grimmer et al., 2023; Tran-Dinh & Luo, 2023).

To describe these special cases of (1), we state the following *nonmonotone* inclusion problem, which generalizes (1):

Find 
$$x^* \in \mathbb{R}^d$$
 such that  $0 \in F(x^*) + G(x^*)$ , (2)

where  $F: \mathbb{R}^d \to \mathbb{R}^d$  is L-Lipschitz and  $G: \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  is maximally monotone. Mapping this problem to finding stationary points of (1) is standard by setting  $x = \binom{u}{v}$ ,  $F(x) = \binom{\nabla_u f(u,v)}{\nabla_v f(u,v)}$  and  $G(x) = \binom{\partial_{\iota_U}}{\partial \iota_V}$ , where  $\iota_U$  is the indicator function for set U. The nonmonotonicity in problem (2) is due to nonconvex-nonconcavity of problem (1).

The main additional assumption we make is that F+G is  $\rho$ -cohypomonotone. Recalling the standard definition  $\operatorname{gra}(F+G)=\{(x,u)\in\mathbb{R}^d\times\mathbb{R}^d\colon\ u\in(F+G)(x)\},$   $\rho$ -cohypomonotonicity is defined as

$$\langle u - v, x - y \rangle \ge -\rho \|u - v\|^2$$
  
 
$$\forall (x, u) \in \operatorname{gra}(F + G) \text{ and } \forall (y, v) \in \operatorname{gra}(F + G),$$
 (3)

for  $\rho>0$ , see (Bauschke et al., 2021, Def. 2.4). When (3) holds only for  $y=x^\star$ , it is also called the *weak MVI condition* or  $\rho$ -star-cohypomonotonicity, due to (Diakonikolas et al., 2021). For  $\rho>0$ , the weak MVI condition requires the existence of a solution  $x^\star$  to the  $\rho$ -weakly MVI:

$$\langle u, x - x^* \rangle \ge -\rho \|u\|^2 \quad \forall (x, u) \in \operatorname{gra}(F + G).$$
 (4)

<sup>&</sup>lt;sup>1</sup>University of Wisconsin–Madison, USA <sup>2</sup>University of British Columbia, Canada <sup>3</sup>KAIST, Republic of Korea. Correspondence to: Ahmet Alacaoglu <a href="mailto:</a> <a href="mailto:ahmet.ed">ahmet Alacaoglu <a href="mailto:ahmet.ed">alacaoglu@math.ubc.ca</a>.

For standard monotone operators (corresponding to convex-concave instances of (1)), the inner product in (3) is lower bounded by 0. The assumption (3) allows the right-hand side to be negative, allowing nonmonotonicity of F+G or nonconvex-nonconcavity of f(u,v), while the limit of nonmonotonicity is determined by  $\rho>0$ . These two assumptions, cohypomonotonicity or weak MVI, are required in the extensive literature cited above.

As the first contribution of this paper, we extend the range of  $\rho$ , doubling the upper limit of  $\frac{1}{2L}$  considered in the previous works, thus allowing a wider range of nonconvex problems of the form (1) to be solved by first-order algorithms, while ensuring optimal or best-known complexity guarantees.

Motivation. Cohypomonotonicity and weak MVI conditions, defined in (3) and (4), allowed progress to be made in understanding the behavior of first-order algorithms for structured nonconvex-nonconcave problems, in a wide variety of works cited at the end of first paragraph. On the one hand, these assumptions are not as general as one might desire: They have not been shown to hold for problems arising in generative or adversarial ML. On the other hand, they have been proven to hold for other relevant problems in ML.

Examples where cohypomonotonicity holds include the *interaction dominant min-max problems* (Example 1) and some stylized worst-case nonconvex-nonconcave instances (Hsieh et al., 2021; Pethick et al., 2023b) (see also (Bauschke et al., 2021, Sections 5, 6)). The relaxed assumption of having a weak MVI solution is implied by star (and quasi-strong) monotonicity (Loizou et al., 2021) or existence of a solution to MVI (Dang & Lan, 2015), the latter being relevant in the context of policy gradient algorithms for reinforcement learning (RL) (Lan, 2023). Weak MVI condition is satisfied in the context of an RL problem described in Example 2.

Example 1. Interaction dominant min-max problems (Grimmer et al., 2023): We say that f in (1) is  $\alpha(r)$ -interaction dominant if it satisfies for all  $z = \binom{u}{v} \in \mathbb{R}^{n+m}$  that

$$\nabla_{uu}^{2} f(z) + \nabla_{uv}^{2} f(z) (r^{-1} \operatorname{Id} - \nabla_{vv}^{2} f(z))^{-1} \nabla_{vu}^{2} f(z)$$

$$\succeq \alpha(r) \operatorname{Id},$$

$$-\nabla_{vv}^{2} f(z) + \nabla_{vu}^{2} f(z) (r^{-1} \operatorname{Id} + \nabla_{uu}^{2} f(z))^{-1} \nabla_{uv}^{2} f(z)$$

$$\succeq \alpha(r) \operatorname{Id}.$$

Interaction is captured by the second terms on the left-hand side of each condition. The problem is called (nonnegative) interaction dominant if these terms dominate the smallest eigenvalue of  $\nabla^2_{uu}f$  and largest eigenvalue of  $\nabla^2_{vv}f$ , i.e.,  $\alpha(r) \geq 0$ . This is equivalent to the r-cohypomonotonicity of F (Hajizadeh et al., 2023, Proposition 1).

Example 2. Instances of von Neumann's ratio game: This is a simple two player stochastic game (Neumann, 1945;

Daskalakis et al., 2020; Diakonikolas et al., 2021). Using the standard definition of the simplex  $\Delta^d=\{x\in\mathbb{R}^d\colon,x\geq0,\sum_{i=1}^dx_i=1\}$ , the problem is

$$\min_{x \in \Delta^m} \max_{y \in \Delta^n} \frac{\langle x, Ry \rangle}{\langle x, Sy \rangle},$$

where  $R \in \mathbb{R}^{m \times n}$ ,  $S \in \mathbb{R}_+^{m \times n}$  and  $\langle x, Sy \rangle > 0 \ \forall (x, y) \in \Delta^m \times \Delta^n$ . As described in (Diakonikolas et al., 2021), it is easy to construct instances of this problem where it satisfies  $\rho$ -weakly MVI condition, but not cohypomonotonicity.  $\blacklozenge$ 

The limit for the parameter  $\rho$  in (3) and (4) for which convergence first-order complexity results are proven seems to have stagnated at  $\rho < \frac{1}{2L}$ . Two exceptions exist for a special case of our setting when  $G \equiv 0$ , which corresponds in view of (1) to an unconstrained problem. First is the recent work (Fan et al., 2024) that claimed to improve the limit of  $\rho$  for weak MVI to  $\approx \frac{0.63}{L}$  with a rather complicated analysis. The rate obtained is also suboptimal under cohypomonotonicity. This work conjectured (but did not prove)  $\frac{1}{L}$  as the maximum limit for  $\rho$  and also did not provide any algorithm achieving this. For an unconstrained cohypomonotone problem, (Cai et al., 2023, Corollary 4.5) also showed possibility of obtaining guarantees with  $\rho < \frac{1}{\sqrt{2L}} \approx \frac{0.7}{L}$ . Relevant citations and discussions appear in Table 1 and Appendix D.

**First-order oracles.** As standard in the operator splitting literature (see e.g., (Bauschke & Combettes, 2017)), a first-order oracle call for (2) consists of one evaluation of F and one resolvent of G (see (5)). In the context of the min-max problem (1), this requires computation of gradients  $\nabla_u f(u,v)$ ,  $\nabla_v f(u,v)$  together with projections on sets U,V. (All works in Table 1 have the same oracle access.) See Assumption 4 for the oracles in the stochastic case.

Contributions. We show how to increase the range of the cohypomonotonicity parameter to  $\rho < \frac{1}{L}$  while maintaining first-order oracle complexity  $\widetilde{O}(\varepsilon^{-1})$  for finding a point x such that  $\mathrm{dist}(0,(F+G)(x)) \leq \varepsilon$ , in Section 2. Such a complexity is optimal (up to a log factor) even for monotone problems (Yoon & Ryu, 2021, Section 3). In Section 3, with weak MVI and the improved range of  $\rho < \frac{1}{L}$ , we show complexity  $\widetilde{O}(\varepsilon^{-2})$  for  $\mathrm{dist}(0,(F+G)(x)) \leq \varepsilon$  which is the best-known (up to a log factor) under this assumption. Table 1 summarizes known results on complexity and the upper bound of  $\rho$ .

Thanks to the modularity of our approach, we extend our results to the stochastic case where F is accessed via unbiased oracles  $\widetilde{F}(\cdot)$  (that is,  $\mathbb{E}[\widetilde{F}(x)] = F(x)$ ). These extensions require the development of further tools for stochastic minmax problems. First, in Section 2.2.1, we tighten the analysis of Halpern iteration with inexact resolvent computations. This leads to improvements for the existing complexities even for some classes of convex-concave problems, see Sec-

Assumption	Reference	Upper bound of $\rho$	Constraints	Oracle complexity
cohypomonotone	(Cai & Zheng, 2022)	$\frac{1}{60L}$	<b>√</b>	$O(\varepsilon^{-1})$
	(Cai et al., 2022b), (Pethick et al., 2023b) (Lee & Kim, 2021), (Tran-Dinh, 2023) (Gorbunov et al., 2023)	$rac{1}{2L}$	✓	$O(\varepsilon^{-1})$
	(Cai et al., 2023)	$\frac{0.7}{L}$	×	$\widetilde{O}(\varepsilon^{-1})$
	Theorem 2.1	$\frac{1}{L}$	$\checkmark$	$\widetilde{O}(\varepsilon^{-1})$
weak MVI	(Diakonikolas et al., 2021) <sup>‡</sup>	$\frac{1}{8L}$	×	$O(\varepsilon^{-2})$
	(Böhm, 2022) <sup>‡</sup>	$\frac{1}{2L}$	×	$O(\varepsilon^{-2})$
	(Cai & Zheng, 2022)	$\frac{1}{12\sqrt{3}L}$	$\checkmark$	$O(\varepsilon^{-2})$
	(Lee & Kim, 2024) <sup>‡</sup>	$\frac{1}{3L}$	$\checkmark$	$\widetilde{O}(\varepsilon^{-2})$
	(Pethick et al., 2022)	$\frac{1}{2L}$	$\checkmark$	$O(\varepsilon^{-2})$
	(Fan et al., 2024)	$\frac{0.63}{L}$	×	$O(\varepsilon^{-2})$ $\widetilde{O}(\varepsilon^{-2})$
	Theorem 3.1	$\frac{1}{L}$	✓	$\widetilde{O}(\varepsilon^{-2})$

Table 1. Comparison of first-order algorithms for deterministic problems. Complexity refers to the number of oracle calls to get  $\operatorname{dist}(0,(F+G)(x)) \leq \varepsilon$ . See also Remark 2.3. <sup>‡</sup>These works defined weak MVI as  $\langle F(x),x-x^{\star}\rangle \geq -\frac{\gamma}{2}\|F(x)\|^2$ , i.e.,  $\gamma=2\rho$ .

tion 4.1. Second, to obtain the best-known complexity for stochastic problems under weak MVI, we incorporate the multilevel Monte Carlo estimator to KM iteration to control the bias in subproblem solutions, see Section 4.2.

# 1.1. Preliminaries

**Notation.** We denote the  $\ell_2$  norm as  $\|\cdot\|$ . Given  $G \colon \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ , we use standard definitions  $\operatorname{gra} G = \{(x,u) \in \mathbb{R}^d \times \mathbb{R}^d \colon u \in G(x)\}$  and  $\operatorname{dist}(0,G(x)) = \min_{u \in G(x)} \|u\|$ . Domain of an operator is defined as  $\operatorname{dom} G = \{x \in \mathbb{R}^d \colon G(x) \neq \emptyset\}$ . The operator G is *maximally* monotone (resp. cohypomonotone or hypomonotone) if its graph is not strictly contained in the graph of any other monotone (resp. cohypomonotone or hypomonotone) operator.

An operator  $F: \mathbb{R}^d \rightrightarrows \mathbb{R}^d$ , given  $(x,u) \in \operatorname{gra} F$  and  $(y,v) \in \operatorname{gra} F$ , is (i)  $\gamma$ -strongly monotone if  $\langle u-v,x-y \rangle \geq \gamma \|x-y\|^2$  with  $\gamma>0$  and monotone if the inequality holds with  $\gamma=0$ ; (ii)  $\rho$ -hypomonotone if  $\langle u-v,x-y \rangle \geq -\rho \|x-y\|^2$  with  $\rho>0$ . An operator  $F: \mathbb{R}^d \to \mathbb{R}^d$  is (iii) L-Lipschitz if  $\|F(x)-F(y)\| \leq L\|x-y\|$ ; (iv) non-expansive if F is 1-Lipschitz; (v)  $\gamma$ -cocoercive if  $\langle F(x)-F(y),x-y \rangle \geq \gamma \|F(x)-F(y)\|^2$  with  $\gamma>0$ . We refer to star variants of these properties (e.g., star-cocoercive) when they are required only at  $(y,v)=(x^\star,0)$  where  $0\in F(x^\star)$ . Since it is a standard notion, we use quasi-nonexpansive instead of star-nonexpansive.

The *resolvent* of an operator  $F : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  is defined as

$$J_F = (\mathrm{Id} + F)^{-1}.$$
 (5)

The resolvent generalizes the well-known proximal operator that has been ubiquitous in optimization and ML, where F is typically the subdifferential of a regularizer function, e.g.,  $\ell_1$  norm. Favorable properties of the resolvent are well-known when F is monotone (Bauschke & Combettes, 2017). Meanwhile, in our nonmonotone case, immense care must be taken in utilizing this object, as it might even be undefined. A comprehensive reference for the properties of resolvent of a nonmonotone operator is (Bauschke et al., 2021). We review and explain the results relevant to our work in the sequel.

The algorithms we analyze are based on the classical Halpern (Halpern, 1967) and Krasnosel'skiĭ-Mann (KM) (Krasnosel'skii, 1955; Mann, 1953) iterations. Given an operator  $T: \mathbb{R}^d \to \mathbb{R}^d$ , Halpern iteration is defined as

$$x_{k+1} = \beta_k x_0 + (1 - \beta_k) T(x_k), \tag{6}$$

for a decreasing sequence  $\{\beta_k\} \in (0,1)$  and initial point  $x_0$ . The KM iteration, with a fixed  $\beta \in (0,1)$ , is defined as

$$x_{k+1} = \beta x_k + (1 - \beta)T(x_k). \tag{7}$$

Conic nonexpansiveness. The key to relaxing the range of  $\rho$  parameter for both assumptions is to harness the algorithmic consequences of *conic nonexpansiveness*, the notion introduced by the influential work of Bauschke et al. (2021) that also inspired our developments. We say that  $T: \mathbb{R}^d \to \mathbb{R}^d$  is  $\lambda$ -conically nonexpansive with  $\lambda > 0$  when there exists a nonexpansive operator  $N: \mathbb{R}^d \to \mathbb{R}^d$  such that  $T = (1 - \lambda)\mathrm{Id} + \lambda N$ , see (Bauschke et al., 2021, Def. 3.1).

This equivalently means that a particular combination of Id and T is nonexpansive:  $\|((1-\lambda^{-1})\mathrm{Id} + \lambda^{-1}T)(x-y)\| \le \|x-y\|$ . An important characterization of this property given in (Bauschke et al., 2021, Cor. 3.5(iii)) is that T is  $\lambda$ -conically nonexpansive if and only if  $\mathrm{Id} - T$  is  $\frac{1}{2\lambda}$ -cocoercive. We also consider the *star* variants (in the sense defined in the Notation paragraph) of these properties and characterizations, which are detailed in Appendix B.1.1.

**Assumption 1.** The operator  $F : \mathbb{R}^d \to \mathbb{R}^d$  is L-Lipschitz and  $G : \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  is maximally monotone. The solution set for the problem (2) is nonempty.

Assumption 1 is standard, see (Facchinei & Pang, 2003), and is required throughout the text. Monotonicity is not assumed for F. Lipschitzness of F corresponds to smoothness of f in context of (1) and maximal monotonicity of G is satisfied when we have constraint sets given in (1) but also when we have convex regularizers added on (1) (e.g.,  $\|\cdot\|_1$ ).

**Assumption 2.** The operator F + G is maximally  $\rho$ -cohypomonotone (see (3) for the definition).

Assumption 2 is abundant in the recent literature for nonconvex-nonconcave optimization (Lee & Kim, 2021; Bauschke et al., 2021; Cai et al., 2022b; Cai & Zheng, 2022; Gorbunov et al., 2023; Pethick et al., 2023b). An instance is provided in Example 1 with further pointers to related problems given in Section 1. Assumption 2 is required only for the results in Sections 2 and 4.1.

**Assumption 3.** There exists a nonempty subset of solutions of (2) whose elements satisfy (4).

Assumption 3 is weaker than Assumption 2 as it is only required with respect to a solution, see also Example 2. Assumption 3, used in Sections 3 and 4.2, is also widespread in the recent literature for nonconvex-nonconcave optimization (Diakonikolas et al., 2021; Pethick et al., 2022; 2023a; Cai et al., 2022b; Lee & Kim, 2024; Fan et al., 2024; Böhm, 2022).

# 2. Algorithm and Analysis under Cohypomonotonicity

# 2.1. Algorithm Construction and Analysis Ideas

Recall the definitions of resolvent (5) and cohypomonotonicity (3). We sketch the algorithmic construction and analysis ideas which will be expanded on in Section 2.2.

(I) We know that Halpern iteration in (6) with  $\beta_k = \frac{1}{k+2}$  has optimal rate when T is nonexpansive, see (Sabach & Shtern, 2017; Lieder, 2021; Kim, 2021). That is, one gets  $||x_k - T(x_k)|| \le \varepsilon$  with  $O(\varepsilon^{-1})$  evaluations of T.

(II) When F+G is maximally  $\rho$ -cohypomonotone (per Assumption 2), we know from (Bauschke et al., 2021) (with precise pointers in Fact A.1) that  $J_{\eta(F+G)}$  is  $\frac{1}{2\alpha}$ -conically nonexpansive where  $\alpha=1-\frac{\rho}{\eta}$ , its domain is  $\mathbb{R}^d$  and it is single-valued when  $\frac{\rho}{\eta}<1$ . Consequently,  $T=(1-\alpha)\mathrm{Id}+\alpha J_{\eta(F+G)}$  is firmly nonexpansive (see Fact A.1). Then, one can use the result in (I).

We next see a high level discussion on the approximate computation of  $J_{\eta(F+G)}$ .

(III) Since F is L-Lipschitz, we have that F is L-hypomonotone by Cauchy-Schwarz inequality, i.e.,

$$\langle F(x) - F(y), x - y \rangle \ge -L||x - y||^2.$$

Hence,  $\mathrm{Id} + \eta F$  is  $(1 - \eta L)$ -strongly monotone.

By definition, we have  $x_k^\star = J_{\eta(F+G)}(x_k) = (\mathrm{Id} + \eta(F+G))^{-1}(x_k)$ . Existence and uniqueness of  $x_k^\star$  is guaranteed by (II) when  $\rho < \eta$  (see Fact A.1). By definition,  $x_k^\star$  is the solution of the problem

$$0 \in (\mathrm{Id} + \eta(F+G))(x_k^{\star}) - x_k.$$

Hence, computation of the resolvent is a strongly monotone inclusion problem where  $\mathrm{Id} + \eta F$  is  $(1 - \eta L)$ -strongly monotone and  $(\eta L + 1)$ -Lipschitz, and G is maximally monotone. In view of (1) this also corresponds to a strongly convex-strongly concave problem.

(IV) Any optimal algorithm for monotone inclusions, such as forward-backward-forward (FBF) (Tseng, 2000), gives  $\hat{x}_k$  with  $\|\hat{x}_k - J_{\eta(F+G)}(x_k)\|^2 \le \varepsilon_k^2$  with complexity  $\widetilde{O}\left(\frac{1+\eta L}{1-\eta L}\right)$ .

In summary, our requirements are  $\frac{\rho}{\eta} < 1$  for ensuring well-definedness of the resolvent, as per (II), and  $1 - \eta L > 0$  for ensuring strong monotonicity for efficient approximation of the resolvent, as per (III). Hence, we need  $\rho < \eta < \frac{1}{L}$ , leading to the claimed improved range on  $\rho$ .

Item (II) refers to the resolvent of  $\eta(F+G)$ , which cannot be evaluated exactly in general with standard first-order oracles. We approximate  $J_{\eta(F+G)}$ , which leads to the inexact Halpern iteration, similar to (Diakonikolas et al., 2021; Cai et al., 2023). Note that in the context of problem (1), approximating the resolvent corresponds to computing approximation of *proximal operator* for function f which is a strongly convex-strongly concave min-max problem.

In the next section, by extending the arguments in (Diakonikolas, 2020, Lemma 12) and (Cai et al., 2023, Lemma C.3) to accommodate conic nonexpansiveness, we

#### Algorithm 1 Inexact Halpern iteration for problems with cohypomonotonicity

Input: Parameters 
$$\beta_k = \frac{1}{k+2}, \eta > 0, L, \rho, \alpha = 1 - \frac{\rho}{\eta}, K \geq 1$$
, initial iterate  $x_0 \in \mathbb{R}^d$ , subroutine FBF in Algorithm 2 for  $k = 0, 1, 2, \dots, K - 1$  do 
$$\widetilde{J}_{\eta(F+G)}(x_k) = \text{FBF}\left(x_k, N_k, \eta G, \text{Id} + \eta F, 1 + \eta L\right) \text{ where } N_k = \left\lceil \frac{4(1+\eta L)}{1-\eta L} \log(98\sqrt{k+2}\log(k+2)) \right\rceil$$

$$x_{k+1} = \beta_k x_0 + (1-\beta_k)((1-\alpha)x_k + \alpha \widetilde{J}_{\eta(F+G)}(x_k))$$

Algorithm 2 FBF 
$$(z_0, N, A, B_{\mathrm{in}}, L_B)$$
 from (Tseng, 2000)

Input: Parameter  $\tau = \frac{1}{2L_B}$ , initial iterate  $z_0 \in \mathbb{R}^d$ ,  $B(\cdot) = B_{\mathrm{in}}(\cdot) - z_0$ 

for  $t = 0, 1, 2, \dots, N-1$  do

 $z_{t+1/2} = J_{\tau A}(z_t - \tau B(z_t))$ 
 $z_{t+1} = z_{t+1/2} + \tau B(z_t) - \tau B(z_{t+1/2})$ 
end for

show that  $\eta^{-1}\|x_k-J_{\eta(F+G)}(x_k)\|\leq \varepsilon$ , where the number of (outer) Halpern iterations is  $O\left(\frac{\|x_0-x^\star\|}{(\eta-\rho)\varepsilon}\right)$ , when we approximate the resolvent to an accuracy of poly  $\left(\frac{1}{k}\right)$ . To achieve this, we can run a subsolver as per (IV), with  $\widetilde{O}\left(\frac{1+\eta L}{1-\eta L}\right)$  calls to evaluations of F and resolvents of G. By combining the complexities at outer and inner levels, we obtain the optimal first-order complexity under  $\rho<\frac{1}{L}$ .

**Discussion.** From the construction (I)-(IV), we see that the ingredients of our approach are based on known results. This raises the question: what insight makes it possible to go beyond the  $\rho < \frac{1}{2L}$  barrier? The key is conic nonexpansiveness, the critical notion introduced by Bauschke et al. (2021). In particular, previous results on first-order complexity for nonmonotone problems (including (Pethick et al., 2023b) who utilized a similar algorithmic construction based on KM as ours in Section 3) used nonexpansiveness of the resolvent, which asks for the stringent requirement  $\rho \leq \frac{\eta}{2} < \frac{1}{2L}$ . This allows Halpern or KM iteration to be analyzed in a standard way.

Our main starting insight is that, from the viewpoint of the analysis of Halpern iteration, we do not necessarily need nonexpansiveness of  $J_{\eta(F+G)}$ . We can apply the Halpern iteration to the operator  $T=(1-\alpha)\mathrm{Id}+\alpha J_{\eta(F+G)}$  where  $\alpha=1-\frac{\rho}{\eta}$ , which is *firmly nonexpansive* for  $\rho<\eta$  (see Fact A.1). Hence, as long as  $\rho<\eta$ , Halpern iteration can be analyzed with  $\rho<\eta<\frac{1}{L}$ , at essentially no cost. We see later how *firm nonexpansiveness* is essential for improving the inexactness criterion in approximating the resolvent.

#### 2.2. Analysis

We now analyze the construction described in the previous section, given as Algorithm 1. We start with the main result,

see Section 2.2.3 and Appendix A.4 for its proof.

**Theorem 2.1.** Let Assumptions 1 and 2 hold. Let  $\eta < \frac{1}{L}$  in Algorithm 1 and suppose  $\rho < \eta$ . For any k = 1, ..., K, we have that  $(x_k)$  from Algorithm 1 satisfies

$$\frac{1}{\eta^2} \|x_k - J_{\eta(F+G)}(x_k)\|^2 \le \frac{16 \|x_0 - x^*\|^2}{(\eta - \rho)^2 (k+1)^2}.$$

The number of first-order oracles used at iteration k is upper bounded by  $2N_k$  where  $N_k$  is defined in Algorithm 1.

**Corollary 2.2.** Under the setting of Theorem 2.1, for any  $\varepsilon > 0$ , we have  $\eta^{-1} \| (\operatorname{Id} - J_{\eta(F+G)})(x_K) \| \le \varepsilon$ , for  $K \le \left\lceil \frac{4 \|x_0 - x^*\|}{(\eta - \rho)\varepsilon} \right\rceil$  and first-order oracle complexity

$$\widetilde{O}\left(\frac{(1+\eta L)\|x_0-x^{\star}\|}{\varepsilon(\eta-\rho)(1-\eta L)}\right).$$

**Remark 2.3.** The definition of  $x^*$  gives that  $(\operatorname{Id} - J_{\eta(F+G)})(x^*) = 0$  and  $(\operatorname{Id} - J_{\eta(F+G)})(x_k)$  is indeed the fixed point residual, which is a standard way to measure optimality for fixed point iterations, see e.g., (Ryu & Yin, 2022, Section 2.4.2). Based on Cor. 2.2, it is straightforward to produce  $x^{\text{out}}$  with  $\operatorname{dist}(0, (F+G)(x^{\text{out}})) \leq \varepsilon$  as claimed in Table 1, with no change in the worst-case complexity. This is clear when  $G \equiv 0$ . In the general case, see (Cai et al., 2023, Lemma C.4).

Remark 2.4. The constant in our complexity deteriorates as  $\rho$  gets close to  $\eta$  which is the same as most of the works included in Table 1. It is straightforward to make our bound  $\rho$ -independent in view of (Pethick et al., 2023b) by simply expressing  $\rho$  as a fraction of  $\eta$ , e.g. assume  $\rho < \frac{9\eta}{10}$ . Then, at the expense of a constant multiple of 10, we have the complexity  $\widetilde{O}\left(\frac{(1+\eta L)\|x_0-x^\star\|}{\varepsilon\eta(1-\eta L)}\right)$ , valid for the range  $\rho < \frac{9}{10L}$ . In comparison, the  $\rho$ -independent complexity result in (Pethick et al., 2023b) had  $\widetilde{O}(\varepsilon^{-2})$  for  $\rho < \frac{1}{2L}$ . A similar reasoning by slightly restricting the range of  $\rho$  can also make the algorithms agnostic to the knowledge of  $\rho$ .

**Outline of the analysis.** We follow the steps sketched in Section 2.1. First, we analyze Halpern iteration with inexactness using the tools mentioned in (I), (II). Second, we analyze the inner loop (Algorithm 2) as mentioned in (IV). Finally we piece together these ingredients.

#### 2.2.1. OUTER-LOOP COMPLEXITY

We now analyze Halpern iteration with inexactness in the resolvent computation. See Appendix A.2 for the proof.

**Lemma 2.5.** Let Assumptions 1 and 2 hold. Suppose that the iterates  $(x_k)$  of Algorithm 1 satisfy  $||J_{\eta(F+G)}(x_i) - \widetilde{J}_{\eta(F+G)}(x_i)|| \le \varepsilon_i$  for some  $\varepsilon_i > 0$  and  $\rho < \eta$ . Let  $R = \operatorname{Id} - J_{\eta(F+G)}$ . Then, we have for any  $K \ge 1$  that

$$\begin{split} & \frac{K(K+1)}{4} \|R(x_K)\|^2 - \frac{K+1}{K\alpha^2} \|x^* - x_0\|^2 \\ & \leq \sum_{k=0}^{K-1} \left( \frac{(k+1)(k+2)\varepsilon_k^2}{2} + (k+1) \|R(x_k)\|\varepsilon_k \right), \end{split}$$

where 
$$||x_k - x^*|| \le ||x_0 - x^*|| + \frac{\alpha}{k+1} \sum_{i=0}^{k-1} (i+1)\varepsilon_i$$
.

In (8) below, we define appropriate values for  $\varepsilon_k$ , and show that the number of inner iterations  $N_k$  selected for FBF in Algorithm 1 suffices to achieve the inexactness level  $\varepsilon_k$ .

This analysis extends Diakonikolas (2020), who studied monotone inclusions, in two aspects. First, we analyze the convergence of the method under conic nonexpansiveness which is the relevant property when the parameter  $\rho$  lies in the range  $\left[\frac{1}{2L}, \frac{1}{L}\right)$ . Second, and more importantly, we conduct a tighter error analysis that allows the inexactness on the error in resolvent computation  $(\varepsilon_k)$  to be  $\widetilde{O}(k^{-3/2})$  instead of the tolerance  $\widetilde{O}(k^{-3})$  used in (Diakonikolas, 2020; Yoon & Ryu, 2022; Cai et al., 2023). Even though it is not immediately obvious, this is because the bottleneck term on the bound in Lemma 2.5 is  $\sum_{k=0}^{K-1} (k+1)(k+2)\varepsilon_k^2$  which sums to a log with  $\varepsilon_k = \widetilde{O}(k^{-3/2})$ .\textstyle{1} This tightening becomes important in the stochastic case in Section 4, where the inner loop does not have a linear rate of convergence.

The improvement derives from applying Halpern to the *firmly nonexpansive* operator  $(1-\alpha)\mathrm{Id} + \alpha J_{\eta(F+G)}$ , which helps avoid the main source of *looseness* in the previous analysis which only uses nonexpansiveness. We discuss this further following (18). See Remark 2.6 for a discussion from the viewpoint of nonexpansive operators.

**Remark 2.6.** By  $\frac{1}{2\alpha}$ -conic nonexpansiveness of  $J_{\eta(F+G)}$  (see Fact A.1(ii)), we have nonexpansiveness of  $T'=(1-2\alpha)\mathrm{Id}+2\alpha J_{\eta(F+G)}$ . If we were to apply Halpern iteration to this operator, we would still obtain results with  $\rho<\frac{1}{L}$  but we would need a stricter inexactness requirement as the analyses in (Diakonikolas, 2020; Cai et al., 2023) dictate.

This can be viewed as a Cayley (or reflection) operator of a firmly nonexpansive operator  $T=(1-\alpha)\mathrm{Id}+\alpha J_{\eta(F+G)}$ , since  $T'=2T-\mathrm{Id}$ . Our algorithm applies Halpern iteration to T which helps us relax the inexactness requirement.

On the other hand, as shown in (Ryu & Yin, 2022, Section 12.2), while solving monotone inclusions with *exact* evaluations of the resolvent, applying Halpern to the Cayley operator of the resolvent gives a better constant, by a factor of 4. Our analysis brings to light a tradeoff between the constant in the convergence bound and the allowed inexactness in the computation of the resolvent.

#### 2.2.2. Inner-Loop Complexity

The seminal FBF algorithm of (Tseng, 2000) is optimal for solving the resolvent subproblem, which is a strongly monotone inclusion. We provide the derivation of the precise constants appearing in the statement in Appendix A.3.

**Theorem 2.7.** (See (Tseng, 2000, Theorem 3.4)) Let B be  $\mu$ -strongly monotone with  $\mu > 0$  and  $L_B$ -Lipschitz; A be maximally monotone, and  $z^* = (A+B)^{-1}(0) \neq \emptyset$ . For any  $\zeta > 0$ , running Algorithm 2 with  $\tau = \frac{1}{2L_B}$  and initial point  $z_0$  for  $N = \left\lceil \frac{4L_B}{\mu} \log \frac{\|z_0 - z^*\|}{\zeta} \right\rceil$  iterations give

$$||z_N - z^\star|| \le \zeta,$$

where the number of calls to evaluations of B and resolvents of A is upper bounded by 2N.

#### 2.2.3. TOTAL COMPLEXITY

Section 2.1 already shows the key steps in our analysis, but we combine the preliminary results above into a proof sketch here, to highlight the simplicity of our approach. Full proof is given in Appendix A.4.

Proof sketch of Theorem 2.1. Denote  $R = \operatorname{Id} - J_{\eta(F+G)}$  for brevity. Suppose that  $\varepsilon_k$  in Lemma 2.5 satisfies

$$\varepsilon_k = \frac{\gamma \|R(x_k)\|}{\sqrt{k+2}\log(k+2)}, \text{ with } \gamma = \frac{1}{98}.$$
 (8)

We justify this supposition further below. Then we have by Lemma 2.5 (after multiplying both sides by  $\alpha$ ) that

$$\frac{\alpha K(K+1)}{4} \|R(x_K)\|^2 - \frac{K+1}{K\alpha} \|x_0 - x^*\|^2$$

$$\leq \sum_{k=0}^{K-1} \|R(x_k)\|^2 \left(\frac{\alpha \gamma^2 (k+1)}{2 \log^2 (k+2)} + \frac{\alpha \gamma \sqrt{k+2}}{\log (k+2)}\right).$$

We can show by induction from this bound that

$$||R(x_k)|| \le \frac{4||x_0 - x^*||}{\alpha(k+1)} \quad \forall k \ge 1.$$

We see that for  $K \leq \lceil \frac{4\|x_0 - x^\star\|}{\eta \alpha \varepsilon} \rceil$ , we are guaranteed to have  $\eta^{-1} \|R(x_K)\| \leq \varepsilon$ .

We now calculate the number of inner iterations to reach the accuracy  $\varepsilon_k$  (see (8)). At iteration k, as per the setup in

<sup>&</sup>lt;sup>1</sup>A similar insight appeared in a different context in the independent work (Liang et al., 2024), which came out on arXiv at the same time as our paper.

Theorem 2.7, we set

$$A \equiv \eta G, \quad B(\cdot) \equiv (\mathrm{Id} + \eta F)(\cdot) - x_k, \quad z_0 \equiv x_k,$$
  
$$z_N \equiv \widetilde{J}_{\eta(F+G)}(x_k), \quad z^* \equiv J_{\eta(F+G)}(x_k), \quad \zeta \equiv \varepsilon_k,$$

hence  $z_0-z^\star=(\operatorname{Id}-J_{\eta(F+G)})(x_k)=R(x_k).$  B is  $L_B\equiv(1+\eta L)$ -Lipschitz and  $(1-\eta L)$ -strongly monotone due to Fact A.1(iv). Existence of  $z^\star$  is guaranteed by Fact A.1(i).

By matching these definitions with Algorithm 1, we see by invoking Theorem 2.7 that the number of inner iterations used at step k to obtain  $\|J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)\| \le \varepsilon_k$  is

$$N_k \equiv \left\lceil \frac{4(1+\eta L)}{1-\eta L} \log \frac{\|R(x_k)\|}{\varepsilon_k} \right\rceil,$$

by the settings of  $z_0$ ,  $z^*$ ,  $R(x_k)$ , and  $\zeta$  above, along with  $\varepsilon_k$  defined in (8). This value is precisely  $N_k$  used in Algorithm 1 (by the definition of  $\varepsilon_k$ ), which justifies our application of Lemma 2.5 and the cost at iteration k.

# 3. Algorithm and Analysis under weak MVI

#### 3.1. Algorithm Construction and Analysis Ideas

We turn to the *weak MVI condition* of Assumption 3, which (as mentioned in Section 1.1) is weaker than cohypomonotonicity. The best-known complexity under this assumption is  $O(\varepsilon^{-2})$ : the lower part of Table 1 outlines existing results. Our aim is to obtain  $\widetilde{O}(\varepsilon^{-2})$  complexity for the extended range  $\rho < \frac{1}{L}$ . The steps of our construction are as follows.

- (i) KM iteration (7), when  $\operatorname{Id} T$  is star-cocoercive, gets  $\eta^{-1} \| x_k T(x_k) \| \le \varepsilon$  with  $O(\varepsilon^{-2})$  evaluations of T (Groetsch, 1972; Browder & Petryshyn, 1967).
- (ii) We get from (Bauschke et al., 2021) that  $J_{\eta(F+G)}$  has domain  $\mathbb{R}^d$  and is single-valued when F is L-Lipschitz and  $\eta < \frac{1}{L}$ . Lemma B.3 gives that  $J_{\eta(F+G)}$  is  $\frac{1}{2\alpha}$ -conically quasi-nonexpansive, with  $\alpha = 1 \frac{\rho}{\eta}$ , leading to  $\mathrm{Id} J_{\eta(F+G)}$  being  $\alpha$ -star-cocoercive.
  - Thus, we require  $\rho < \eta$ . As per (i), KM applied to  $\mathrm{Id} J_{\eta(F+G)}$  requires  $O(\varepsilon^{-2})$  evaluations of  $J_{\eta(F+G)}$  to find x such that  $\eta^{-1} \|x J_{\eta(F+G)}(x)\| \leq \varepsilon$ .
- (iii) Since F is Lipschitz and G is maximally monotone, we can estimate  $J_{\eta(F+G)}$  as before (via (III) and (IV) of Section 2), with a linear rate of convergence when  $\eta < \frac{1}{L}$ . The existence of a solution to the subproblem is guaranteed by item (ii). The inner iterations introduce a logarithmic factor into the total complexity. As a result, the range for  $\rho$  is again  $\rho < \eta < \frac{1}{L}$ .

Even with inexactness, Alg. 3 is classical; see (Facchinei & Pang, 2003, Theorem 12.3.7), (Combettes, 2001) and (Combettes & Pennanen, 2002). We analyze this scheme

for problems with weak MVI solutions and characterize the first-order oracle complexity. Pethick et al. (2023b) recently analyzed a similar scheme under cohypomonotonicity, by using quasi-nonexpansiveness of the resulting operator.<sup>2</sup> Our main difference regarding the results in this section is that we harness the milder property of conic quasi-nonexpansiveness to improve the range of  $\rho$  (see also (Bartz et al., 2022) for a similar idea by using exact resolvent). We also approximate the resolvent slightly differently. FBF can be replaced with other optimal algorithms like (Malitsky & Tam, 2020), showing the modularity of our approach.

The key insight for extending the upper bound of  $\rho$  to  $\frac{1}{L}$  is similar to that of Section 2. The difference is that the analysis of Halpern iteration requires conic nonexpansiveness between any pair of points in the space, making it unsuitable with weak MVI. In contrast, the KM iteration can be analyzed with conic nonexpansiveness holding only with respect to a solution, a property that is a consequence of weak MVI. Conic quasi-nonexpansiveness, while not defined explicitly in (Bauschke et al., 2021), directly follows by adapting the corresponding results therein by using  $\rho$ -weak MVI condition instead of cohypomonotonicity; see Appendix B.1.1 for the details.

#### 3.2. Analysis

Similar to Section 2, we start with the main complexity result, under weak MVI. Its proof appears in Appendix B.4.

**Theorem 3.1.** Let Assumptions 1 and 3 hold. Let  $\eta < \frac{1}{L}$  in Algorithm 3 and suppose  $\rho < \eta$ . For any  $K \ge 1$ , we have

$$\frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{\eta^2} \|x_k - J_{\eta(F+G)}(x_k)\|^2 \le \frac{11 \|x_0 - x^\star\|^2}{(\eta - \rho)^2 K}.$$

The number of first-order oracles used at iteration k is upper bounded by  $2N_k$  where  $N_k$  is defined in Algorithm 3.

**Corollary 3.2.** Under the setting of Theorem 3.1, for any  $\varepsilon > 0$ , we have for some  $x^{\text{out}} \in \{x_0, \dots, x_{K-1}\}$  that  $\eta^{-1} \| (\operatorname{Id} - J_{\eta(F+G)})(x^{\text{out}}) \| \le \varepsilon$  for  $K \le \left\lceil \frac{11 \| x_0 - x^* \|}{(\eta - \rho)^2 \varepsilon^2} \right\rceil$  with first-order oracle complexity

$$\widetilde{O}\left(\frac{(1+\eta L)\|x_0-x^{\star}\|^2}{\varepsilon^2(\eta-\rho)^2(1-\eta L)}\right).$$

See Remark 2.3 for details to convert this result to produce a point with  $\operatorname{dist}(0, (F+G)(x^{\operatorname{out}})) \leq \varepsilon$  as in Table 1.

**Remark 3.3.** This result is for the *best iterate*, that is,  $x^{\text{out}} = \arg\min_{x \in \{x_0, \dots, x_{k-1}\}} \| (\operatorname{Id} - J_{\eta(F+G)})(x) \|$ , consistent with existing results for weak MVI, see (Diakonikolas et al., 2021; Pethick et al., 2022; Cai & Zheng, 2022).

<sup>&</sup>lt;sup>2</sup>This work claimed that some of their results extend to accommodate weak MVI condition as well.

#### Algorithm 3 Inexact KM iteration for problems with weak MVI

Input: Parameters  $\eta>0, L, \rho, \alpha_k=\alpha=1-\frac{\rho}{\eta}, K>0$ , initial iterate  $x_0\in\mathbb{R}^d$ , subroutine FBF in Algorithm 2

for 
$$k=0,1,2,\ldots,K-1$$
 do 
$$\widetilde{J}_{\eta(F+G)}(x_k)=\operatorname{FBF}(x_k,N_k,\eta G,\operatorname{Id}+\eta F,1+\eta L), \text{ where } N_k=\left\lceil\frac{4(1+\eta L)}{1-\eta L}\log(8(k+1)\log^2(k+2))\right\rceil$$
  $x_{k+1}=(1-\alpha_k)x_k+\alpha_k\widetilde{J}_{\eta(F+G)}(x_k)$  end for

Remark 3.4. Note that  $x^{\text{out}}$  as defined in Remark 3.3 is not computable since we do not have access to  $J_{\eta(F+G)}(x_k)$ . For the unconstrained case, i.e.,  $G \equiv 0$ , we can show the result with  $x^{\text{out}} = \arg\min_{x \in \{x_0, \dots, x_{K-1}\}} \|Fx\|^2$ , which is computable. For the constrained problem (1), we can handle this issue by slightly changing how  $\widetilde{J}_{\eta(F+G)}$  is calculated and requiring the knowledge of the target accuracy  $\varepsilon$ , with no change in the order of complexity bounds. We present Algorithm 3 in its current form so that it is *anytime*, not requiring the target accuracy as an input. The details for making  $x^{\text{out}}$  computable are in Appendix B.5. We can also present this result as an *expected* bound for a *randomly selected*  $x^{\text{out}}$ , like (Diakonikolas et al., 2021, Thm. 3.2(ii)).

**Outer-loop complexity.** We analyze the iteration complexity of the outer loop; see Appendix B.2 for a proof which is a modification of (Combettes, 2001) and (Bartz et al., 2022) to accommodate conic quasi-nonexpansiveness and inexact resolvent computations.

**Lemma 3.5.** Let Assumptions 1 and 3 hold. Suppose that the iterates  $(x_k)$  of Algorithm 3 satisfy  $||J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)|| \le \varepsilon_k$  for some  $\varepsilon_k > 0$  and  $\rho < \eta$ . Then, we have for  $K \ge 1$  that

$$\sum_{k=0}^{K-1} \| (\operatorname{Id} - J_{\eta(F+G)})(x_k) \|^2 - \frac{2\eta^2}{(\eta - \rho)^2} \| x_0 - x^* \|^2$$

$$\leq 6 \sum_{k=0}^{K-1} \varepsilon_k^2 + \frac{4\eta}{\eta - \rho} \sum_{k=0}^{K-1} \| x_k - x^* \| \varepsilon_k,$$

where 
$$||x_k - x^*|| \le ||x_{k-1} - x^*|| + \alpha \varepsilon_{k-1}$$
.

**Total Complexity.** The sketch of the proof of Theorem 3.1 follows Section 2.2.3 closely. We use Lemma 3.5 instead of Lemma 2.5. The choice of  $\varepsilon_k$  is slightly different, as can be noticed by the number of inner iterations  $N_k$  in Algorithm 3. However, with the same argument in Section 2.2.3, we can show that this  $N_k$  is sufficient to attain the inexactness required by  $\varepsilon_k$ .

### 4. Algorithms and Analyses with Stochasticity

In this case, F in (2) is accessed via unbiased oracles.

**Assumption 4.** The stochastic first-order oracle (SFO)

$$\widetilde{F} \colon \mathbb{R}^d \to \mathbb{R}^d$$
 satisfies

$$F(x) = \mathbb{E}[\widetilde{F}(x)]$$
 and  $\mathbb{E}||\widetilde{F}(x) - F(x)||^2 \le \sigma^2$ .

In view of (1), this corresponds to using *stochastic gradients*  $\widetilde{F}(x) = \begin{pmatrix} \widetilde{\nabla}_u f(u,v) \\ -\widetilde{\nabla}_v f(u,v) \end{pmatrix}$  where  $\mathbb{E}[\widetilde{\nabla}_u f(u,v)] = \nabla_u f(u,v)$  (and similarly for the v component). Table 2, with comparisons for stochastic problems, is in Appendix C. The variance assumption could be relaxed by using, e.g., an argument similar to (Wright & Recht, 2022, Section 5.4.3).

#### 4.1. Cohypomonotone Case

For this setup, Algorithm 1 will call FBF with stochastic oracles  $\widetilde{F}(x_t)$  as per Assumption 4 to approximate  $\widetilde{J}_{n(F+G)}$ :

$$\widetilde{J}_{\eta(F+G)}(x_k) = \text{FBF}(x_k, N_k, \eta G, \text{Id} + \eta \widetilde{F}, 1 + \eta L),$$
 (9)

where 
$$N_k = \lceil 1734(k+2)^3 \log^2(k+2)(1-\eta L)^{-2} \rceil$$
.

**Corollary 4.1.** Let Assumptions 1, 2 and 4 hold. Let  $\eta < \frac{1}{L}$  in Alg. 1,  $\rho < \eta$  and use (9) for computing  $\widetilde{J}_{\eta(F+G)}$  (see Alg. 4). Then we have for  $k \geq 1$  that

$$\eta^{-2} \mathbb{E} \|x_k - J_{\eta(F+G)}(x_k)\|^2 = O(k^{-2}).$$

For any  $\varepsilon > 0$ , we have  $\eta^{-1}\mathbb{E}\|(\operatorname{Id} - J_{\eta(F+G)})(x_K)\| \le \varepsilon$  for the last iterate, with SFO complexity  $\widetilde{O}(\varepsilon^{-4})$ .

The proof, provided in Appendix C.2.1 is the stochastic adaptation of Section 2. Our tighter analysis for the level of inexactness (which is highlighted after Lemma 2.5) is the main reason we could get the  $\widetilde{O}(\varepsilon^{-4})$  complexity. The inexactness level required by following the existing analyses in (Diakonikolas, 2020; Cai et al., 2023) would instead result in a  $\widetilde{O}(\varepsilon^{-7})$  complexity.

Remark 4.2. The previous *last iterate* result for constrained, cohypomonotone, stochastic problems by (Pethick et al., 2023b, Corollary E.3(ii)) was  $\widetilde{O}(\varepsilon^{-16})$  (in fact we are not aware of another last iterate result even for stochastic and constrained convex-concave problems). This result also required increasing batch sizes in the inner loop and  $\rho < \frac{1}{2L}$ . For unconstrained problems, Chen & Luo (2022) showed an improved  $\widetilde{O}(\varepsilon^{-2})$  expected complexity for  $\rho < \frac{1}{2L}$  with some drawbacks described in Appendix D. It is an open question to get a similar complexity improvement in our constrained setup with a wider range for  $\rho$ .

**Remark 4.3.** Pethick et al. (2023a) has complexity  $\widetilde{O}(\varepsilon^{-4})$ for a constrained problem with weak MVI. However, this work additionally assumed a stronger oracle model and Lipschitzness assumptions. In particular, denoting  $F(\cdot) = \mathbb{E}_{\xi \sim \Xi}[F_{\xi}(\cdot)]$  for an unknown  $\Xi$  that we can sample from, this work assumes mean-square (MS)-Lipschitzness:  $\mathbb{E}_{\xi \sim \Xi} \|F_{\xi}(x) - F_{\xi}(y)\|^2 \le L^2 \|x - y\|^2$ . This work also needs to query the operator for the same seed for two different points:  $F_{\xi}(x_k)$ ,  $F_{\xi}(x_{k-1})$ . These two assumptions define a different template. For nonconvex minimization, for example, lower bounds improve with these assumptions compared to our standard stochastic approximation setting in Assumption 4, see (Arjevani et al., 2023). Moreover, the additional assumption might not hold even for trivial problems:  $F_1(x) = x^2$ ,  $F_2(x) = -x^2$  where  $F = F_1 + F_2$ is clearly Lipschitz but not MS-Lipschitz.

#### 4.2. Weak MVI Case

We next modify Algorithm 3 for the stochastic case. The main observation from the analysis (see Lemma C.8) is that bounding the bias  $\|\mathbb{E}[\widetilde{J}_{\eta(F+G)}(x_k)] - J_{\eta(F+G)}(x_k)\|$  with square root of variance  $\mathbb{E}\|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\|^2$  by Jensen's inequality is loose and would give complexity  $\widetilde{O}(\varepsilon^{-6})$ , like (Pethick et al., 2023b, Cor. E.3(i)).

A natural candidate for a careful bias analysis is the multilevel Monte Carlo (MLMC) technique which helps control the bias-variance tradeoff (Giles, 2008; Blanchet & Glynn, 2015; Asi et al., 2021; Hu et al., 2021). The high level idea is that stochastic KM iteration, in our setting would give  $O(\varepsilon^{-4})$  complexity if we had unbiased samples of  $J_{\eta(F+G)}$  (see, e.g., (Bravo & Cominetti, 2024)). Obtaining such unbiased samples is highly non-trivial since  $J_{\eta(F+G)}$  is an inclusion problem in itself. Fortunately, MLMC is a way to get an estimator with bias  $O(\varepsilon')$  and variance  $\widetilde{O}(1)$  by making, in expectation,  $\widetilde{O}(1)$  calls to the oracle defined in Assumption 4. MLMC is used in (Asi et al., 2021) for the related proximal point algorithm.

Estimator 1 (MLMC). We set  $\widetilde{J}_{n(F+G)}$  as follows.

1. Given  $N_k \geq 1$ ,  $M_k \geq 1$ , set for  $m = 1, \dots, M_k$ ,

$$\begin{split} \widetilde{J}_{\eta(F+G)}^{(m)}(x_k) &= \begin{cases} y^0 + 2^I(y^I - y^{I-1}) \text{ if } I \leq N_k, \\ y^0, & \text{otherwise,} \end{cases} \\ \text{where } I \sim \text{Geom}(1/2) \\ \text{and } y^i &= \text{FBF}(x_k, 2^i, G, \text{Id} + \eta \widetilde{F}, 1 + \eta L) \ \forall i > 0. \end{split}$$

2. Given  $M_k$  independent draws of this estimator, we define  $\widetilde{J}_{\eta(F+G)}(x_k) = \frac{1}{M_k} \sum_{m=1}^{M_k} \widetilde{J}_{\eta(F+G)}^{(m)}(x_k)$ .

To show that the scheme is *implementable* we give the (non-optimized) values of  $M_k$ ,  $N_k$ . This is to ensure that they are

agnostic to unknown quantities  $\{||x_0 - x^*||^2, \sigma^2\}$ , unlike some MLMC methods (Chen & Luo, 2022).

Corollary 4.4. Let Assumptions 1, 3 and 4 hold. In Algorithm 3, set  $\eta < \frac{1}{L}$ ,  $\alpha_k \equiv \frac{\alpha}{\sqrt{k+2\log(k+3)}}$ , suppose that  $\rho < \eta$  and use Estimator 1 for computing  $\widetilde{J}_{\eta(F+G)}$  (see Algorithm 6) with  $N_k \equiv \lceil \frac{96(1-\eta L)^{-2}}{\min\{\frac{\alpha_k}{120\alpha(k+1)},\frac{1}{120}\}} \rceil$  and  $M_k \equiv \lceil \frac{672\times120(\log_2N_k)}{(1-\eta L)^2} \rceil$ . For any  $\varepsilon > 0$ , we have that  $\eta^{-1}\mathbb{E}\|(\mathrm{Id}-J_{\eta(F+G)})(x^{\mathrm{out}})\| \leq \varepsilon$ , with expected SFO complexity  $\widetilde{O}(\varepsilon^{-4})$  where  $x^{\mathrm{out}}$  is selected uniformly at random from  $\{x_0,\ldots,x_{K-1}\}$ .

Proof of this corollary appears in Appendix C.3.1. This result is an alternative to (Pethick et al., 2023a) that required additional assumptions as explained in Remark 4.3. In our setting under Assumption 4, the only  $O(\varepsilon^{-4})$  complexity was known in the special case of unconstrained problems  $(G \equiv 0)$ , due to (Diakonikolas et al., 2021) (see also (Choudhury et al., 2023)). Because of the use of MLMC, our complexity result is *expected* number of stochastic oracle calls and hence the results mentioned in this paragraph complement each other. See also Table 2.

MLMC is used in conditional/compositional stochastic minimization (Hu et al., 2021), distributionally robust optimization (Levy et al., 2020), and stochastic minimization with non-i.i.d. data (Dorfman & Levy, 2022). Our development of the KM iteration with MLMC can provide the potential to extend some of these results to stochastic min-max setting.

# 5. Conclusions

We conclude with some open questions. Even though our results for nonmonotone problems, either with cohypomonotonicity or weak MVI conditions, can go beyond the existing barrier for the  $\rho$  parameter, our algorithms rely on a double loop strategy, alternating between a Halpern or KM step and resolvent approximation. This strategy also results in an additional log factor in the final complexity bound. Two worthwhile directions in this context are: (i) developing a single loop algorithm with the extended  $\rho$  range (ii) obtaining first-order complexities without spurious log terms and extended range for  $\rho$ . It is also critical to find more problems in ML that satisfy these nonmonotonicity assumptions.

We next highlight a direction for stochastic problems. As mentioned before, the  $O(\varepsilon^{-4})$  first-order complexity seems to be the best-known for even constrained, convex-concave stochastic problems. However, for unconstrained problems, better complexities are known, see, e.g., (Chen & Luo, 2022; Cai et al., 2022a). The tools in these works may be combined with the ideas in our paper to develop better complexity results for constrained stochastic min-max problems with convex-concave or nonconvex-nonconcave functions.

#### Acknowledgements

This work was supported in part by the NSF grant 2023239, the NSF grant 2224213, the AFOSR award FA9550-21-1-0084, National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. 2019R1A5A1028324, 2022R1C1C1003940), and the Samsung Science & Technology Foundation grant (No. SSTF-BA2101-02).

A. Alacaoglu is thankful to Vidya Muthukumar and Panayotis Mertikopoulos for helpful discussions about the implications of the results in Section 4.1. We also thank a reviewer of our paper who helped us improve the presentation in Section 2.1.

This work was done while A. Alacaoglu was at the University of Wisconsin–Madison.

# **Impact Statement**

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

#### References

- Alacaoglu, A., Böhm, A., and Malitsky, Y. Beyond the golden ratio for variational inequality algorithms. *Journal of Machine Learning Research*, 24(172):1–33, 2023.
- Allen-Zhu, Z. How to make the gradients small stochastically: Even faster convex and nonconvex sgd. *Advances in Neural Information Processing Systems*, 31, 2018.
- Arjevani, Y., Carmon, Y., Duchi, J. C., Foster, D. J., Srebro, N., and Woodworth, B. Lower bounds for non-convex stochastic optimization. *Mathematical Programming*, 199 (1-2):165–214, 2023.
- Asi, H., Carmon, Y., Jambulapati, A., Jin, Y., and Sidford, A. Stochastic bias-reduced gradient methods. *Advances in Neural Information Processing Systems*, 34:10810–10822, 2021.
- Bartz, S., Dao, M. N., and Phan, H. M. Conical averagedness and convergence analysis of fixed point algorithms. *Journal of Global Optimization*, 82(2):351–373, 2022.
- Bauschke, H. H. and Combettes, P. L. Convex analysis and monotone operator theory in hilbert spaces. *CMS Books in Mathematics*, 2017.
- Bauschke, H. H., Moursi, W. M., and Wang, X. Generalized monotone operators and their averaged resolvents. *Mathematical Programming*, 189:55–74, 2021.

- Blanchet, J. H. and Glynn, P. W. Unbiased monte carlo for optimization and functions of expectations via multi-level randomization. In *2015 Winter Simulation Conference* (WSC), pp. 3656–3667. IEEE, 2015.
- Böhm, A. Solving nonconvex-nonconcave min-max problems exhibiting weak minty solutions. *Transactions on Machine Learning Research*, 2022.
- Böhm, A., Sedlmayer, M., Csetnek, E. R., and Bot, R. I. Two steps at a time—taking gan training in stride with tseng's method. *SIAM Journal on Mathematics of Data Science*, 4(2):750–771, 2022.
- Bravo, M. and Cominetti, R. Stochastic fixed-point iterations for nonexpansive maps: Convergence and error bounds. *SIAM Journal on Control and Optimization*, 62 (1):191–219, 2024.
- Bravo, M. and Contreras, J. P. Stochastic halpern iteration in normed spaces and applications to reinforcement learning. *arXiv:2403.12338*, 2024.
- Browder, F. E. and Petryshyn, W. V. Construction of fixed points of nonlinear mappings in hilbert space. *Journal of Mathematical Analysis and Applications*, 20(2):197–228, 1967
- Cai, X., Song, C., Guzmán, C., and Diakonikolas, J. Stochastic halpern iteration with variance reduction for stochastic monotone inclusions. *Advances in Neural Information Processing Systems*, 35:24766–24779, 2022a.
- Cai, X., Alacaoglu, A., and Diakonikolas, J. Variance reduced halpern iteration for finite-sum monotone inclusions. In *International Conference on Learning Representations*, 2023.
- Cai, Y. and Zheng, W. Accelerated single-call methods for constrained min-max optimization. In *The Eleventh International Conference on Learning Representations*, 2022.
- Cai, Y., Oikonomou, A., and Zheng, W. Accelerated algorithms for monotone inclusions and constrained nonconvex-nonconcave min-max optimization. *arXiv*:2206.05248, 2022b.
- Chen, L. and Luo, L. Near-optimal algorithms for making the gradient small in stochastic minimax optimization. *arXiv*:2208.05925, 2022.
- Choudhury, S., Gorbunov, E., and Loizou, N. Single-call stochastic extragradient methods for structured non-monotone variational inequalities: Improved analysis under weaker conditions. In *Advances in Neural Information Processing Systems*, 2023.

- Combettes, P. L. Quasi-fejérian analysis of some optimization algorithms. In *Studies in Computational Mathematics*, volume 8, pp. 115–152. Elsevier, 2001.
- Combettes, P. L. and Pennanen, T. Generalized mann iterates for constructing fixed points in hilbert spaces. *Journal of Mathematical Analysis and Applications*, 275(2): 521–536, 2002.
- Combettes, P. L. and Pennanen, T. Proximal methods for cohypomonotone operators. *SIAM journal on control and optimization*, 43(2):731–742, 2004.
- Dang, C. D. and Lan, G. On the convergence properties of non-euclidean extragradient methods for variational inequalities with generalized monotone operators. *Computational Optimization and applications*, 60:277–310, 2015.
- Dao, M. N. and Phan, H. M. Adaptive douglas—rachford splitting algorithm for the sum of two operators. *SIAM Journal on Optimization*, 29(4):2697–2724, 2019.
- Daskalakis, C., Foster, D. J., and Golowich, N. Independent policy gradient methods for competitive reinforcement learning. Advances in neural information processing systems, 33:5527–5540, 2020.
- Daskalakis, C., Skoulakis, S., and Zampetakis, M. The complexity of constrained min-max optimization. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pp. 1466–1478, 2021.
- Diakonikolas, J. Halpern iteration for near-optimal and parameter-free monotone inclusion and strong solutions to variational inequalities. In *Conference on Learning Theory*, pp. 1428–1451. PMLR, 2020.
- Diakonikolas, J., Daskalakis, C., and Jordan, M. I. Efficient methods for structured nonconvex-nonconcave min-max optimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 2746–2754. PMLR, 2021.
- Dorfman, R. and Levy, K. Y. Adapting to mixing time in stochastic optimization with markovian data. In *Interna*tional Conference on Machine Learning, pp. 5429–5446. PMLR, 2022.
- Facchinei, F. and Pang, J.-S. Finite-dimensional variational inequalities and complementarity problems. Springer, 2003.
- Fan, Y., Li, Y., and Chen, B. Weaker MVI condition: Extragradient methods with multi-step exploration. In *The Twelfth International Conference on Learning Representations*, 2024.
- Giles, M. B. Multilevel monte carlo path simulation. *Operations research*, 56(3):607–617, 2008.

- Giselsson, P. and Moursi, W. M. On compositions of special cases of lipschitz continuous operators. *Fixed Point Theory and Algorithms for Sciences and Engineering*, 2021 (1):1–38, 2021.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. Advances in neural information processing systems, 27, 2014.
- Gorbunov, E., Taylor, A., Horváth, S., and Gidel, G. Convergence of proximal point and extragradient-based methods beyond monotonicity: the case of negative comonotonicity. In *International Conference on Machine Learning*, pp. 11614–11641. PMLR, 2023.
- Grimmer, B., Lu, H., Worah, P., and Mirrokni, V. The landscape of the proximal point method for nonconvex–nonconcave minimax optimization. *Mathematical Programming*, 201(1-2):373–407, 2023.
- Groetsch, C. A note on segmenting mann iterates. *Journal of Mathematical Analysis and Applications*, 40(2):369–372, 1972.
- Hajizadeh, S., Lu, H., and Grimmer, B. On the linear convergence of extragradient methods for nonconvex–nonconcave minimax problems. *INFORMS Journal on Optimization*, 2023.
- Halpern, B. Fixed points of nonexpanding maps. *Bulletin of the American Mathematical Society*, 73(6):957–961, 1967.
- Hsieh, Y.-G., Iutzeler, F., Malick, J., and Mertikopoulos, P. On the convergence of single-call stochastic extragradient methods. *Advances in Neural Information Processing Systems*, 32, 2019.
- Hsieh, Y.-P., Mertikopoulos, P., and Cevher, V. The limits of min-max optimization algorithms: Convergence to spurious non-critical sets. In *International Conference on Machine Learning*, pp. 4337–4348. PMLR, 2021.
- Hu, Y., Chen, X., and He, N. On the bias-variance-cost tradeoff of stochastic optimization. *Advances in Neural Information Processing Systems*, 34:22119–22131, 2021.
- Kim, D. Accelerated proximal point method for maximally monotone operators. *Mathematical Programming*, 190 (1-2):57–87, 2021.
- Kohlenbach, U. On the proximal point algorithm and its halpern-type variant for generalized monotone operators in hilbert space. *Optimization Letters*, 16(2):611–621, 2022.

- Kotsalis, G., Lan, G., and Li, T. Simple and optimal methods for stochastic variational inequalities, i: operator extrapolation. *SIAM Journal on Optimization*, 32(3):2041–2073, 2022.
- Krasnosel'skii, M. A. Two remarks on the method of successive approximations. *Uspekhi matematicheskikh nauk*, 10(1):123–127, 1955.
- Lan, G. Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. *Mathematical programming*, 198(1):1059–1106, 2023.
- Lee, S. and Kim, D. Fast extra gradient methods for smooth structured nonconvex-nonconcave minimax problems. *Advances in Neural Information Processing Systems*, 34: 22588–22600, 2021.
- Lee, S. and Kim, D. Semi-anchored gradient methods for nonconvex-nonconcave minimax problems, 2024. URL https://openreview.net/forum?id=rmLTwKGiSP.
- Leuştean, L. and Pinto, P. Quantitative results on a halperntype proximal point algorithm. *Computational Optimization and Applications*, 79(1):101–125, 2021.
- Levy, D., Carmon, Y., Duchi, J. C., and Sidford, A. Large-scale methods for distributionally robust optimization. Advances in Neural Information Processing Systems, 33: 8847–8860, 2020.
- Liang, L., Toh, K.-C., and Zhu, J.-J. An inexact halpern iteration for with application to distributionally robust optimization. *arXiv:2402.06033*, 2024.
- Lieder, F. On the convergence rate of the halpern-iteration. *Optimization letters*, 15(2):405–418, 2021.
- Loizou, N., Berard, H., Gidel, G., Mitliagkas, I., and Lacoste-Julien, S. Stochastic gradient descent-ascent and consensus optimization for smooth games: Convergence analysis under expected co-coercivity. Advances in Neural Information Processing Systems, 34:19095–19108, 2021.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learn*ing Representations, 2018.
- Malitsky, Y. and Tam, M. K. A forward-backward splitting method for monotone inclusions without cocoercivity. *SIAM Journal on Optimization*, 30(2):1451–1472, 2020.
- Mann, W. R. Mean value methods in iteration. *Proceedings* of the American Mathematical Society, 4(3):506–510, 1953.

- Neumann, J. v. A model of general economic equilibrium. *The Review of Economic Studies*, 13(1):1–9, 1945.
- Pethick, T., Patrinos, P., Fercoq, O., Cevher, V., and Latafat, P. Escaping limit cycles: Global convergence for constrained nonconvex-nonconcave minimax problems. In *International Conference on Learning Representations*, 2022.
- Pethick, T., Fercoq, O., Latafat, P., Patrinos, P., and Cevher, V. Solving stochastic weak minty variational inequalities without increasing batch size. In *International Confer*ence on Learning Representations, 2023a.
- Pethick, T., Xie, W., and Cevher, V. Stable nonconvexnonconcave training via linear interpolation. In *Thirtyseventh Conference on Neural Information Processing Systems*, 2023b.
- Ryu, E. K. and Yin, W. *Large-scale convex optimization: algorithms & analyses via monotone operators.* Cambridge University Press, 2022.
- Sabach, S. and Shtern, S. A first order method for solving convex bilevel optimization problems. *SIAM Journal on Optimization*, 27(2):640–660, 2017.
- Tran-Dinh, Q. Sublinear convergence rates of extragradient-type methods: A survey on classical and recent developments. *arXiv:2303.17192*, 2023.
- Tran-Dinh, Q. and Luo, Y. Randomized block-coordinate optimistic gradient algorithms for root-finding problems. *arXiv*:2301.03113, 2023.
- Tseng, P. A modified forward-backward splitting method for maximal monotone mappings. *SIAM Journal on Control and Optimization*, 38(2):431–446, 2000.
- Wright, S. J. and Recht, B. *Optimization for data analysis*. Cambridge University Press, 2022.
- Yoon, T. and Ryu, E. K. Accelerated algorithms for smooth convex-concave minimax problems with o (1/k<sup>2</sup>) rate on squared gradient norm. In *International Conference on Machine Learning*, pp. 12098–12109. PMLR, 2021.
- Yoon, T. and Ryu, E. K. Accelerated minimax algorithms flock together. *arXiv:2205.11093*, 2022.

#### A. Proofs for Section 2

#### A.1. Preliminary Results

We start with the properties of the resolvent of a cohypomonotone operator and the properties of the subproblem for approximating this resolvent. These important points are also sketched in Section 2.1. We present this preliminary result here for the ease of reference throughout the proofs. Most of the conclusions follow from the results of (Bauschke et al., 2021). Note that  $\rho$ -cohypomonotone in our notation is  $-\rho$ -comonotone in the notation of (Bauschke et al., 2021). See also (Bauschke et al., 2021, Remark 2.5) for these two conventions.

**Fact A.1.** Let Assumptions 1 and 2 hold and let  $\eta > 0$ . Then, we have

- (i) The operator  $J_{\eta(F+G)}$  is single-valued and dom  $J_{\eta(F+G)} = \mathbb{R}^d$  when  $\rho < \eta$ .
- (ii) The operator  $J_{\eta(F+G)}$  is  $\frac{1}{2(1-\frac{\rho}{\eta})}$ -conically nonexpansive,  $\operatorname{Id} J_{\eta(F_G)}$  is  $\left(1-\frac{\rho}{\eta}\right)$ -cocoercive, and  $(1-\alpha)\operatorname{Id} + \alpha J_{\eta(F+G)}$  is firmly nonexpansive when  $\rho < \eta$ .
- (iii) For any  $\bar{x} \in \mathbb{R}^d$ , computing  $J_{\eta(F+G)}(\bar{x})$  is equivalent to solving the problem:

Find 
$$x \in \mathbb{R}^d$$
 such that  $0 \in (\mathrm{Id} + \eta(F+G))(x) - \bar{x}$ . (10)

*The problem* (10) *has a unique solution when*  $\rho < \eta$ .

- (iv) The operator  $\mathrm{Id} + \eta F$  is  $(1 + \eta L)$ -Lipschitz and  $(1 \eta L)$ -strongly monotone when  $\eta < \frac{1}{L}$ .
- *Proof.* (i) By Assumption 2 and the definition of cohypomonotonicity in (3), we have that  $\eta(F+G)$  is maximally  $\frac{\rho}{\eta}$ -cohypomonotone. Then for  $\frac{\rho}{\eta} < 1$ , (Bauschke et al., 2021, Corollary 2.14) gives the result.
- (ii) Since  $\eta(F+G)$  is maximally  $\frac{\rho}{\eta}$ -cohypomonotone, (Bauschke et al., 2021, Prop. 3.11(ii)) gives  $\frac{1}{2\left(1-\frac{\rho}{\eta}\right)}$ -conic nonexpansiveness. Cocoercivity of  $\mathrm{Id}-J_{\eta(F+G)}$  then follows from (Bauschke et al., 2021, Corollary 3.5(iii)). By the definition of conic nonexpansiveness, we have that  $T'=(1-2\alpha)\mathrm{Id}+2\alpha J_{\eta(F+G)}$  is nonexpansive. By definition, a firmly nonexpansive operator is one that can be written as  $\frac{1}{2}\mathrm{Id}+\frac{1}{2}N$  for a nonexpansive operator N (Bauschke & Combettes, 2017, Remark 4.34(iii)). Since  $(1-\alpha)\mathrm{Id}+\alpha J_{\eta(F+G)}=\frac{1}{2}\mathrm{Id}+\frac{1}{2}T'$  for the nonexpansive T' specified in this paragraph, we conclude.
- (iii) Let us denote  $\bar{x}^* = J_{\eta(F+G)}(\bar{x})$  and use the definition of a resolvent to obtain

$$\bar{x}^* = J_{\eta(F+G)}(\bar{x}) = (\mathrm{Id} + \eta(F+G))^{-1}(\bar{x}) \iff \bar{x}^* + \eta(F+G)(\bar{x}^*) \ni \bar{x},$$

where the existence of  $\bar{x}^*$  is guaranteed by (i). Rearranging the inclusion gives (10). Uniqueness of the solution is due to (i).

(iv) By Lipschitzness of F and Cauchy-Schwarz inequality, we have

$$\langle \eta F(x) - \eta F(y), x - y \rangle \ge -\eta \|F(x) - F(y)\| \|x - y\| \ge -\eta L \|x - y\|^2.$$

As a result, we have that  $\mathrm{Id} + \eta F$  is  $(1 - \eta L)$ -strongly monotone. We also have by triangle inequality that

$$\|(\mathrm{Id} + \eta F)(x) - (\mathrm{Id} + \eta F)y\| \le \|x - y\| + \eta \|F(x) - F(y)\| \le (1 + \eta L)\|x - y\|,$$

completing the proof.

#### A.2. Complexity of the Outer loop

Bounding the norm of the iterates.

**Lemma A.2.** Let Assumptions 1 and 2 hold. Suppose that the iterates  $(x_k)$  of Algorithm 1 satisfy  $||J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)|| \le \varepsilon_k$  for some  $\varepsilon_k > 0$  and  $\rho < \eta$ . Then, we have for  $k \ge 0$  that

$$||x_{k+1} - x^*|| \le ||x_0 - x^*|| + \left(1 - \frac{\rho}{\eta}\right) \frac{1}{k+2} \sum_{i=0}^k (i+1)\varepsilon_i.$$

*Proof.* Recall the following notation from Algorithm 1:

$$\alpha = 1 - \frac{\rho}{\eta} = \frac{\eta - \rho}{\eta}.$$

Then, by Fact A.1(ii), we know that  $J_{\eta(F+G)}$  is  $\frac{1}{2\alpha}$ -conically nonexpansive. This means that we can write  $J_{\eta(F+G)} = (1 - \frac{1}{2\alpha}) \operatorname{Id} + \frac{1}{2\alpha} N$  for a nonexpansive operator N.

Adding and subtracting  $\alpha(1-\beta_k)J_{\eta(F+G)}(x_k)$  in the definition of  $x_{k+1}$  in Algorithm 1, using conic nonexpansiveness of  $J_{\eta(F+G)}$ , and rearranging gives

$$\begin{split} x_{k+1} &= \beta_k x_0 + (1 - \beta_k) \left( (1 - \alpha) x_k + \alpha \widetilde{J}_{\eta(F+G)}(x_k) \right) \\ &= \beta_k x_0 + (1 - \beta_k) \left( (1 - \alpha) x_k + \alpha J_{\eta(F+G)}(x_k) \right) + \alpha (1 - \beta_k) \left( \widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k) \right) \\ &= \beta_k x_0 + \frac{1 - \beta_k}{2} x_k + \frac{1 - \beta_k}{2} N(x_k) + \alpha (1 - \beta_k) \left( \widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k) \right), \end{split}$$

where the last step is because  $J_{\eta(F+G)} = \frac{2\alpha-1}{2\alpha} \mathrm{Id} + \frac{1}{2\alpha} N$  for a nonexpansive operator N.

We now use triangle inequality, nonexpansiveness of N, the definition of  $\varepsilon_k$ , and the last equality to obtain

$$||x_{k+1} - x^*|| \le \beta_k ||x_0 - x^*|| + \frac{1 - \beta_k}{2} ||x_k - x^*|| + \frac{1 - \beta_k}{2} ||N(x_k) - x^*||$$

$$+ \alpha (1 - \beta_k) ||\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)||$$

$$\le \beta_k ||x_0 - x^*|| + (1 - \beta_k) ||x_k - x^*|| + \alpha (1 - \beta_k) \varepsilon_k,$$
(11)

where the inequality used that  $Nx^\star = x^\star$  since  $N = 2\alpha J_{\eta(F+G)} + (1-2\alpha)\mathrm{Id}$  and that  $J_{\eta(F+G)}(x^\star) = x^\star$  by the definition of  $x^\star$  in (2), and Fact A.1(i).

The result of the lemma now follows by induction after using the definition  $\beta_k = \frac{1}{k+2}$ . In particular, the assertion is true for k=0 by inspection. Assume the assertion holds for k=K-1, then (11) gives

$$||x_{K+1} - x^*|| \le \frac{1}{K+2} ||x_0 - x^*|| + \frac{K+1}{K+2} ||x_K - x^*|| + \frac{\alpha(K+1)}{K+2} \varepsilon_K$$

$$\le \frac{1}{K+2} ||x_0 - x^*|| + \frac{K+1}{K+2} \left( ||x_0 - x^*|| + \frac{\alpha}{K+1} \sum_{i=0}^{K-1} (i+1)\varepsilon_i \right) + \frac{\alpha(K+1)}{K+2} \varepsilon_K$$

$$= ||x_0 - x^*|| + \frac{\alpha}{K+2} \sum_{i=0}^{K} (i+1)\varepsilon_i,$$

which completes the induction. The statement follows after using  $\alpha=1-\frac{\rho}{\eta}$ .

#### **Iteration complexity**

**Lemma 2.5.** Let Assumptions 1 and 2 hold. Suppose that the iterates  $(x_k)$  of Algorithm 1 satisfy  $||J_{\eta(F+G)}(x_i) - \widetilde{J}_{\eta(F+G)}(x_i)|| \le \varepsilon_i$  for some  $\varepsilon_i > 0$  and  $\rho < \eta$ . Let  $R = \mathrm{Id} - J_{\eta(F+G)}$ . Then, we have for any  $K \ge 1$  that

$$\frac{\alpha K(K+1)}{4} \|R(x_K)\|^2 - \frac{K+1}{K\alpha} \|x^* - x_0\|^2 \le \sum_{k=0}^{K-1} \left( \frac{\alpha}{2} (k+1)(k+2)\varepsilon_k^2 + \alpha(k+1) \|R(x_k)\|\varepsilon_k \right),$$

where  $\alpha = 1 - \frac{\rho}{\eta}$ , as defined in Algorithm 1 and  $||x_k - x^\star|| \le ||x_0 - x^\star|| + \frac{\alpha}{k+1} \sum_{i=0}^{k-1} (i+1)\varepsilon_i$ .

*Proof of Lemma* 2.5. By Fact A.1(ii), we have that  $\mathrm{Id} - J_{\eta(F+G)}$  is  $\left(1 - \frac{\rho}{\eta}\right)$  cocoercive. Recall the definition of  $\alpha$  from Algorithm 1 and the notation for  $\mathrm{Id} - J_{\eta(F+G)}$  as:

$$\alpha = 1 - \frac{\rho}{\eta}$$
 and  $R = \operatorname{Id} - J_{\eta(F+G)}$ .

With these, we use  $\alpha$ -cocoercivity of R:

$$\langle R(x_{k+1}) - R(x_k), x_{k+1} - x_k \rangle > \alpha \|R(x_{k+1}) - R(x_k)\|^2. \tag{12}$$

By rearranging the update rule of  $x_{k+1}$  in Algorithm 1, we have for  $k \ge 0$  that

$$x_{k+1} = \beta_k x_0 + (1 - \beta_k) x_k - \alpha (1 - \beta_k) (\operatorname{Id} - \widetilde{J}_{\eta(F+G)})(x_k)$$
  
=  $\beta_k x_0 + (1 - \beta_k) x_k - \alpha (1 - \beta_k) R(x_k) + \alpha (1 - \beta_k) (\widetilde{J}_{\eta(F+G)} - J_{\eta(F+G)})(x_k),$  (13)

where we added and subtracted  $\alpha(1-\beta_k)J_{\eta(F+G)}(x_k)$  and used the definition  $R=\mathrm{Id}-J_{\eta(F+G)}$ .

We now use a step that is common in the rate analysis of Halpern-type methods, which can be seen for example in (Diakonikolas, 2020) or (Yoon & Ryu, 2021). In particular, from (13), we obtain two identical representations for  $x_{k+1} - x_k$ :

$$x_{k+1} - x_k = \beta_k(x_0 - x_k) - \alpha(1 - \beta_k)R(x_k) + \alpha(1 - \beta_k)(\widetilde{J}_{\eta(F+G)} - J_{\eta(F+G)})(x_k), \tag{14a}$$

$$x_{k+1} - x_k = \frac{\beta_k}{1 - \beta_k} (x_0 - x_{k+1}) - \alpha R(x_k) + \alpha (\widetilde{J}_{\eta(F+G)} - J_{\eta(F+G)})(x_k), \tag{14b}$$

where the second representation follows from subtracting  $\beta_k x_{k+1}$  from both sides of (13) and rearranging. With these at hand, we develop the left-hand side of (12). First, by using (14b), we have that

$$\langle R(x_{k+1}), x_{k+1} - x_k \rangle = \frac{\beta_k}{1 - \beta_k} \langle R(x_{k+1}), x_0 - x_{k+1} \rangle - \alpha \langle R(x_{k+1}), R(x_k) \rangle + \alpha \langle R(x_{k+1}), (\widetilde{J}_{\eta(F+G)} - J_{\eta(F+G)})(x_k) \rangle = \frac{\beta_k}{1 - \beta_k} \langle R(x_{k+1}), x_0 - x_{k+1} \rangle - \frac{\alpha}{2} \left( \|R(x_{k+1})\|^2 + \|R(x_k)\|^2 - \|R(x_{k+1}) - R(x_k)\|^2 \right) + \alpha \langle R(x_{k+1}), (\widetilde{J}_{\eta(F+G)} - J_{\eta(F+G)})(x_k) \rangle,$$
(15)

where the last step used the expansion  $||a - b||^2 = ||a||^2 - 2\langle a, b \rangle + ||b||^2$ .

Second, by using (14a), we have that

$$-\langle R(x_k), x_{k+1} - x_k \rangle = -\beta_k \langle R(x_k), x_0 - x_k \rangle + \alpha (1 - \beta_k) \|R(x_k)\|^2$$
$$-\alpha (1 - \beta_k) \langle R(x_k), (\widetilde{J}_{n(F+G)} - J_{n(F+G)})(x_k) \rangle. \tag{16}$$

After using (15) and (16) in (12) and rearranging, we obtain

$$\frac{\alpha}{2} \|R(x_{k+1})\|^{2} + \frac{\beta_{k}}{1 - \beta_{k}} \langle R(x_{k+1}), x_{k+1} - x_{0} \rangle 
\leq \frac{\alpha}{2} (1 - 2\beta_{k}) \|R(x_{k})\|^{2} + \beta_{k} \langle R(x_{k}), x_{k} - x_{0} \rangle 
+ \alpha \langle R(x_{k+1}) - (1 - \beta_{k}) R(x_{k}), (\widetilde{J}_{\eta(F+G)} - J_{\eta(F+G)})(x_{k}) \rangle - \frac{\alpha}{2} \|R(x_{k+1}) - R(x_{k})\|^{2}.$$
(17)

For the third term on the right-hand side of (17), we apply Cauchy-Schwarz, triangle and Young's inequalities along with the definition of  $\varepsilon_k$  to obtain

$$\alpha \langle R(x_{k+1}) - (1 - \beta_k) R(x_k), (\widetilde{J}_{\eta(F+G)} - J_{\eta(F+G)})(x_k) \rangle \leq \alpha \|R(x_{k+1}) - (1 - \beta_k) R(x_k)\| \varepsilon_k$$

$$\leq \alpha (\|R(x_{k+1}) - R(x_k)\| + \beta_k \|R(x_k)\|) \varepsilon_k$$

$$= \alpha \|R(x_{k+1}) - R(x_k)\| \varepsilon_k + \alpha \beta_k \|R(x_k)\| \varepsilon_k$$

$$\leq \frac{\alpha}{2} \|R(x_{k+1}) - R(x_k)\|^2 + \frac{\alpha}{2} \varepsilon_k^2 + \alpha \beta_k \|R(x_k)\| \varepsilon_k.$$
 (18)

This is the main point of departure from the existing analysis where this inequality is bounded by  $O(\|x_k - x^*\|\varepsilon_k)$ , cf. (Diakonikolas, 2020, display equation after (14)). We instead use the last term in (17) (which we obtained by using the firm nonexpansiveness of  $(1 - \alpha) \operatorname{Id} + \alpha J_{\eta(F+G)}$ ) to cancel the corresponding error term in (18). We use this last estimate in (17) and get

$$\frac{\alpha}{2} \|R(x_{k+1})\|^2 + \frac{\beta_k}{1 - \beta_k} \langle R(x_{k+1}), x_{k+1} - x_0 \rangle \le \frac{\alpha}{2} (1 - 2\beta_k) \|R(x_k)\|^2 + \beta_k \langle R(x_k), x_k - x_0 \rangle 
+ \frac{\alpha}{2} \varepsilon_k^2 + \alpha \beta_k \|R(x_k)\| \varepsilon_k.$$
(19)

Noting the identities

$$\beta_k = \frac{1}{k+2} \implies 1 - \beta_k = \frac{k+1}{k+2}, \quad \frac{\beta_k}{1-\beta_k} = \frac{1}{k+1}, \quad 1 - 2\beta_k = \frac{k}{k+2},$$

on (19) we obtain

$$\frac{\alpha}{2} \|R(x_{k+1})\|^2 + \frac{1}{k+1} \langle R(x_{k+1}), x_{k+1} - x_0 \rangle \le \frac{\alpha}{2} \frac{k}{k+2} \|R(x_k)\|^2 + \frac{1}{k+2} \langle R(x_k), x_k - x_0 \rangle + \frac{\alpha}{2} \varepsilon_k^2 + \frac{\alpha}{k+2} \|R(x_k)\| \varepsilon_k,$$

which holds for  $k \geq 0$ . Multiplying both sides by (k+1)(k+2) gives

$$\frac{\alpha(k+1)(k+2)}{2} \|R(x_{k+1})\|^2 + (k+2)\langle R(x_{k+1}), x_{k+1} - x_0 \rangle 
\leq \frac{\alpha k(k+1)}{2} \|R(x_k)\|^2 + (k+1)\langle R(x_k), x_k - x_0 \rangle 
+ \frac{\alpha}{2} (k+1)(k+2)\varepsilon_k^2 + \alpha(k+1) \|R(x_k)\|\varepsilon_k.$$

We sum the inequality for k = 0, 1, ..., K - 1 to get

$$\frac{\alpha K(K+1)}{2} \|R(x_K)\|^2 + (K+1)\langle R(x_K), x_K - x_0 \rangle 
\leq \sum_{k=0}^{K-1} \left( \frac{\alpha}{2} (k+1)(k+2)\varepsilon_k^2 + \alpha(k+1) \|R(x_k)\| \varepsilon_k \right).$$
(20)

By the standard estimation for the inner product on this left-hand side (using (i) monotonicity of R, which is implied by  $\alpha$ -cocoercivity of R with  $\alpha > 0$ ; (ii) definition of  $x^*$  as  $R(x^*) = (\operatorname{Id} - J_{\eta(F+G)})(x^*) = 0$  which uses Fact A.1(i); (iii) Young's inequality), we derive

$$(K+1)\langle R(x_K), x_K - x_0 \rangle = (K+1)\langle R(x_K), x^* - x_0 \rangle + (K+1)\langle R(x_K), x_K - x^* \rangle$$

$$\geq (K+1)\langle R(x_K), x^* - x_0 \rangle$$

$$\geq -\frac{\alpha K(K+1)}{4} \|R(x_K)\|^2 - \frac{K+1}{K\alpha} \|x^* - x_0\|^2.$$

We use this lower bound on (20) to conclude the first assertion. The second claim is essentially Lemma A.2.

# A.3. Complexity of the Inner Loop

**Theorem 2.7.** (See e.g., (Tseng, 2000, Theorem 3.4)) Let B be  $\mu$ -strongly monotone with  $\mu > 0$  and  $L_B$ -Lipschitz; A be maximally monotone, and  $z^* = (A+B)^{-1}(0) \neq \emptyset$ . For any  $\zeta > 0$ , running Algorithm 2 with  $\tau = \frac{1}{2L_B}$  and initial point  $z_0$  for  $N = \left\lceil \frac{4L_B}{\mu} \log \frac{\|z_0 - z^*\|}{\zeta} \right\rceil$  iterations give

$$||z_N - z^\star|| \leq \zeta,$$

where the number of calls to evaluations of B and resolvents of A is upper bounded by  $2\left\lceil \frac{4L_B}{\mu}\log\frac{\|z_0-z^\star\|}{\zeta}\right\rceil$ .

*Proof.* We only derive the number of iterations for ease of reference which follows trivially from (Tseng, 2000, Theorem 3.4). In particular, in the notation of (Tseng, 2000, Theorem 3.4(c)), we select  $\theta = \frac{1}{2}$ ,  $\alpha = \frac{1}{2L_B}$  and assume without loss of generality that  $\frac{\mu}{L_B} \leq \frac{3}{4}$  to obtain

$$||z_{t+1} - z^*||^2 \le \left(1 - \frac{\mu}{2L_B}\right) ||z_t - z^*||^2,$$

which after unrolling gives that

$$||z_N - z^*||^2 \le \left(1 - \frac{\mu}{2L_B}\right)^N ||z_0 - z^*||^2.$$

Standard manipulations give that after  $N = \left\lceil \frac{4L_B}{\mu} \log \frac{\|z_0 - z^\star\|}{\zeta} \right\rceil$  iterations, we have  $\|z_N - z^\star\|^2 \le \zeta^2$ .

#### A.4. Total complexity

**Theorem 2.1.** Let Assumptions 1 and 2 hold. Let  $\eta < \frac{1}{L}$  in Algorithm 1 and suppose  $\rho < \eta$ . For any k = 1, ..., K, we have that  $(x_k)$  from Algorithm 1 satisfies

$$\frac{1}{\eta^2} \|x_k - J_{\eta(F+G)}(x_k)\|^2 \le \frac{16 \|x_0 - x^*\|^2}{(\eta - \rho)^2 (k+1)^2}.$$

The number of first-order oracles used at iteration k of Algorithm 1 is upper-bounded by

$$\left\lceil \frac{4(1+\eta L)}{1-\eta L} \log(98\sqrt{k+2}\log(k+2)) \right\rceil.$$

*Proof of Theorem 2.1.* We recall the notations

$$\alpha = 1 - \frac{\rho}{\eta} \quad \text{and} \quad R = \operatorname{Id} - J_{\eta(F+G)}$$

and start from the result of Lemma 2.5 which states for  $K \geq 1$  that

$$\frac{\alpha K(K+1)}{4} \|R(x_K)\|^2 \le \frac{K+1}{K\alpha} \|x^* - x_0\|^2 + \sum_{k=0}^{K-1} \left(\frac{\alpha}{2} (k+1)(k+2)\varepsilon_k^2 + \alpha(k+1) \|R(x_k)\|\varepsilon_k\right).$$

Let us set

$$\varepsilon_k = \frac{\gamma \|R(x_k)\|}{\sqrt{k+2}\log(k+2)} \tag{21}$$

and note that we will not evaluate  $\varepsilon_k$  but we will prove that for a *computable* number of inner iterations  $N_k$ , this error criterion will be satisfied.

We substitute the definition of  $\varepsilon_k$  to the previous inequality and get

$$\frac{\alpha K(K+1)}{4} \|R(x_K)\|^2 \le \frac{K+1}{K\alpha} \|x_0 - x^*\|^2 + \sum_{k=0}^{K-1} \left( \frac{\alpha \gamma^2 (k+1) \|R(x_k)\|^2}{2 \log^2 (k+2)} + \frac{\alpha \gamma \sqrt{k+2} \|R(x_k)\|^2}{\log (k+2)} \right). \tag{22}$$

We now show by induction that

$$||R(x_k)|| \le \frac{4||x_0 - x^*||}{\alpha(k+1)} \quad \forall k \ge 1.$$
 (23)

Note that  $\alpha$ -cocoercivity of R and  $R(x^*) = 0$  gives  $||R(x_0)|| \le \frac{1}{\alpha} ||x_0 - x^*||$ . For k = 1, we have by  $\alpha$ -cocoercivity of R,  $R(x^*) = 0$  and Lemma A.2 that

$$||R(x_1)|| \le \frac{1}{\alpha} ||x_1 - x^*|| \le \frac{1}{\alpha} \left( ||x_0 - x^*|| + \frac{\gamma ||x_0 - x^*||}{2\sqrt{2} \log 2} \right) < \frac{2||x_0 - x^*||}{\alpha}, \tag{24}$$

for  $\gamma = \frac{1}{98}$ , which establishes the base case of induction. Now we assume (23) holds for all  $k \le K - 1$ . Then, we use (22) for  $K \ge 2$  (where we also use  $\frac{K+1}{K} \le 2$ ):

$$\frac{\alpha K(K+1)}{4} \|R(x_K)\|^2 \leq \frac{2}{\alpha} \|x_0 - x^*\|^2 + \sum_{k=0}^{K-1} \left( \frac{\alpha \gamma^2 (k+1) \|R(x_k)\|^2}{2 \log^2 (k+2)} + \frac{\alpha \gamma \sqrt{k+2} \|R(x_k)\|^2}{\log(k+2)} \right) \\
\leq \frac{2}{\alpha} \|x_0 - x^*\|^2 + \sum_{k=0}^{K-1} \left( \frac{8\gamma^2 \|x_0 - x^*\|^2}{\alpha(k+1) \log^2 (k+2)} + \frac{16\gamma \sqrt{k+2} \|x_0 - x^*\|^2}{\alpha(k+1)^2 \log(k+2)} \right).$$

Since we have that

$$\sum_{k=0}^{K-1} \frac{8}{(k+1)\log^2(k+2)} < 28 \text{ and } \sum_{k=0}^{K-1} \frac{16\sqrt{k+2}}{(k+1)^2\log(k+2)} < 49,$$

the value  $\gamma = \frac{1}{98}$  results in

$$\frac{\alpha K(K+1)}{4} \|R(x_K)\|^2 \le \frac{2.6}{\alpha} \|x_0 - x^*\|^2.$$

A direct implication of this inequality is that

$$||R(x_K)||^2 \le \frac{10.4}{\alpha^2 K(K+1)} ||x_0 - x^*||^2$$
  
$$\le \frac{15.6}{\alpha^2 (K+1)^2} ||x_0 - x^*||^2,$$

where we used  $\frac{1}{K(K+1)} \leq \frac{1.5}{(K+1)^2}$  which holds when  $K \geq 2$ . This completes the induction.

We next see that with  $N_k$  set as in Algorithm 1, we get the inexactness level specified by  $\varepsilon_k$  and the oracle complexity of each iteration is as claimed in the statement.

At iteration k, to apply the result in Theorem 2.7, we identify the following settings from Algorithm 1

$$A \equiv \eta G, \quad B(\cdot) \equiv (\mathrm{Id} + \eta F)(\cdot) - x_k, \quad z_0 \equiv x_k, \quad z^* \equiv J_{\eta(F+G)}(x_k), \quad \zeta \equiv \varepsilon_k$$
  
$$\implies z_0 - z^* = (\mathrm{Id} - J_{\eta(F+G)})(x_k) = R(x_k)$$

hence B is  $(1 + \eta L)$ -Lipschitz and  $(1 - \eta L)$ -strongly monotone due to Fact A.1(iv). Existence of  $z^*$  is guaranteed by Fact A.1(iii).

We now see that by the setting of

$$T = \left\lceil \frac{4(1+\eta L)}{1-\eta L} \log(98\sqrt{k+2}\log(k+2)) \right\rceil = \left\lceil \frac{4(1+\eta L)}{1-\eta L} \log \frac{\|R(x_k)\|}{\varepsilon_k} \right\rceil,$$

Theorem 2.7 tells us that

$$\|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\| \le \varepsilon_k,$$

for  $\varepsilon_k$  given in (21), as claimed.

Since each iteration of Algorithm 2 uses 2 evaluations of F and 1 resolvent for G, the first-order oracle complexity is  $2N_k$  and the result follows.

We now continue with the proof of Corollary 2.2 which follows trivially from Theorem 2.1.

*Proof of Corollary* 2.2. By Theorem 2.1, we have that after at most  $\left[\frac{4\|x_0-x^*\|}{(\eta-\rho)\varepsilon}\right]$  iterations, i.e., for a K such that

$$K \le \left\lceil \frac{4\|x_0 - x^*\|}{(\eta - \rho)\varepsilon} \right\rceil,\tag{25}$$

we are guaranteed to have

$$\eta^{-1} \| (\operatorname{Id} - J_{\eta(F+G)})(x_K) \| \le \varepsilon.$$

Total number of first-oracle calls during the run of the algorithm then be calculated as

$$\sum_{k=1}^{K} \left\lceil \frac{4(1+\eta L)}{1-\eta L} \log(98\sqrt{k+2}\log(k+2)) \right\rceil \leq K \cdot \left( \frac{4(1+\eta L)}{1-\eta L} \log(98\sqrt{K+2}\log(K+2)) + 1 \right).$$

We conclude after using (25).

#### B. Proofs for Section 3

#### **B.1. Preliminary results**

We now derive similar properties to Fact A.1 but with Assumption 3. These proofs are slightly more involved than Fact A.1 to accommodate the weaker assumption.

We start with the definition of conic quasi-nonexpansiveness that will be used in Fact B.2. Recall that an operator N is quasi-nonexpansive when  $||Nx - x^*|| \le ||x - x^*||$  where  $x^*$  is a fixed point of N.

**Definition B.1.**  $T: \mathbb{R}^d \to \mathbb{R}^d$  is  $\alpha$ -conically quasi-nonexpansive if there exists a quasi-nonexpansive operator  $N: \mathbb{R}^d \to \mathbb{R}^d$  such that  $T = (1 - \alpha) \mathrm{Id} + \alpha N$ .

This is a modification of conic nonexpansiveness in (Bauschke et al., 2021, Definition 3.1). In Appendix B.1.1, we show the conic quasi-nonexpansiveness (and related properties) of the resolvent of a star-cohypomonotone operator in view of Assumption 3, by invoking the corresponding arguments of (Bauschke et al., 2021) restricted to a point in the domain and a fixed point of the resolvent. Then we show *star*-cocoercivity of  $Id - J_{\eta(F+G)}$  which facilitates the analysis of KM iteration.

Fact B.2. Let Assumptions 1 and 3 hold. Then, we have

- (i) The operator  $J_{\eta(F+G)}$  is single-valued and dom  $J_{\eta(F+G)}=\mathbb{R}^d$  when  $\eta<\frac{1}{L}$ .
- (ii) The operator  $J_{\eta(F+G)}$  is  $\frac{1}{2\left(1-\frac{\rho}{\eta}\right)}$ -conically quasi-nonexpansive and  $\mathrm{Id}-J_{\eta(F_G)}$  is  $\left(1-\frac{\rho}{\eta}\right)$ -star-cocoercive when  $\rho<\eta$ .
- (iii) For any  $\bar{x} \in \mathbb{R}^d$ , computing  $J_{\eta(F+G)}(\bar{x})$  is equivalent to solving the problem

Find 
$$x \in \mathbb{R}^d$$
 such that  $0 \in (\mathrm{Id} + \eta(F+G))(x) - \bar{x}$ . (26)

The problem (26) has a unique solution when  $\eta < \frac{1}{L}$ .

(iv) The operator  $\mathrm{Id} + \eta F$  is  $(1+\eta L)$ -Lipschitz and  $(1-\eta L)$ -strongly monotone when  $\eta < \frac{1}{L}$ .

Proof.

(i) When F is L-Lipschitz, it is maximally L-hypomonotone (see e.g., (Giselsson & Moursi, 2021, Lemma 2.12)) and  $\eta F$  is maximally  $\eta L$ -hypomonotone since it is  $\eta L$ -Lipschitz.

By (Dao & Phan, 2019, Lemma 3.2(ii)), we know that  $\eta F + \mathrm{Id}$  is maximally  $(1 - \eta L)$ -(strongly) monotone. Then, using this and maximal monotonicity of G, we have by (Bauschke & Combettes, 2017, Corollary 25.5) that  $\mathrm{Id} + \eta (F + G)$  is maximally  $(1 - \eta L)$ -(strongly) monotone. Invoking (Dao & Phan, 2019, Lemma 3.2(ii)) again gives us that  $\eta (F + G)$  is maximally  $\eta L$ -hypomonotone.

We can then use (Bauschke et al., 2021, Lemma 2.8) to obtain that  $(\eta(F+G))^{-1}$  is maximally  $\eta L$ -cohypomonotone. This can be combined with (Bauschke et al., 2021, Corollary 2.14) to get the result when  $\eta L < 1$ .

(ii) Since F+G has a  $\rho$ -weak MVI solution under Assumption 3, we have that  $\eta(F+G)$  has  $\rho/\eta$ -weak MVI solution, i.e., by simple change of variables, we have for some  $\eta>0$ 

$$\begin{split} &\langle \eta u, x - x^* \rangle \geq \eta \rho \|u\|^2 \ \text{ where } \ u \in (F+G)(x) \\ &\iff \langle v, x - x^* \rangle \geq \frac{\rho}{\eta} \|v\|^2 \ \text{ where } \ v \in \eta(F+G)(x). \end{split}$$

Proposition B.5 then gives us that  $J_{\eta(F+G)}$  is  $\frac{1}{2\left(1-\frac{\rho}{\eta}\right)}$ -conically quasi-nonexpansive and then we have that  $\mathrm{Id}-J_{\eta(F+G)}$  is  $\left(1-\frac{\rho}{\eta}\right)$  star-cocoercive by Corollary B.4.

- (iii) The proof is the same as Fact A.1(iii) where the only difference is that now we ensure the existence of  $J_{\eta(F+G)}(\bar{x})$  with (i). Uniqueness also follows from this.
- (iv) The proof is the same as Fact A.1(iv).

#### B.1.1. PROPERTIES RELATED TO CONIC QUASI-NONEXPANSIVENESS

This section particularizes the notion and properties of the  $\alpha$ -conic nonexpansiveness in (Bauschke et al., 2021) to their *star* variants. The aim is to show that the properties extend to their *star* or *quasi*-variants when we use weak MVI condition instead of cohypomonotonicity. This sections implicitly assumes that  $J_A$  for operator  $A: \mathbb{R}^d \rightrightarrows \mathbb{R}^d$  is well-defined, sufficient conditions for which is shown in Fact B.2. We say that an operator N is quasi-nonexpansive when  $||Nx - x^*|| \leq ||x - x^*||$  where  $x^*$  is a fixed point of N.

**Lemma B.3.** (See (Bauschke et al., 2021, Lemma 3.4)) Consider  $T: \mathbb{R}^d \to \mathbb{R}^d$  and let  $T = (1 - \alpha) \mathrm{Id} + \alpha N$ . Then, N is quasi-nonexpansive if and only if we have, for all  $x \in \mathbb{R}^d$ ,

$$2\alpha \langle Tx - x^*, (\mathrm{Id} - T)x \rangle \ge (1 - 2\alpha) \|(\mathrm{Id} - T)x\|^2,$$

or equivalently

$$\left\| \left( 1 - \frac{1}{\alpha} \right) x + \frac{1}{\alpha} T x - x^{\star} \right\| \le \|x - x^{\star}\|. \tag{27}$$

*Proof.* Using  $\alpha^2 \|a\|^2 - \|(\alpha - 1)a + b\|^2 = 2\alpha \langle b, a - b \rangle - (1 - 2\alpha) \|a - b\|^2$  (see (Bauschke et al., 2021, Lemma 3.3)) with  $a = x - x^*$  and  $b = Tx - x^*$ , we have

$$0 \le 2\alpha \langle Tx - x^*, (\mathrm{Id} - T)x \rangle - (1 - 2\alpha) \| (\mathrm{Id} - T)x \|^2$$

$$= \alpha^2 \| x - x^* \|^2 - \| (\alpha - 1)(x - x^*) + Tx - x_* \|^2$$

$$= \alpha^2 \| x - x^* \|^2 - \| (\alpha - 1)(x - x^*) + (1 - \alpha)(x - x^*) + \alpha (Nx - x^*) \|^2$$

$$= \alpha^2 (\| x - x^* \|^2 - \| Nx - x^* \|^2),$$

which gives the assertion. Last claim follows by substituting  $N=\frac{1}{\alpha}T+\left(1-\frac{1}{\alpha}\right)\mathrm{Id}$  in the definition of quasi-nonexpansiveness for N.

**Corollary B.4.** (See (Bauschke et al., 2021, Corollary 3.5(iii)))  $T: \mathbb{R}^d \to \mathbb{R}^d$  is  $\alpha$ -conically quasi-nonexpansive if and only if  $\mathrm{Id} - T$  is  $\frac{1}{2\alpha}$ -star-cocoercive.

*Proof.* We use Lemma B.3:

$$\langle Tx - x^*, (\operatorname{Id} - T)x \rangle \ge \left(\frac{1}{2\alpha} - 1\right) \|(\operatorname{Id} - T)x\|^2 \quad \Leftrightarrow \quad \langle x - x^*, (\operatorname{Id} - T)x \rangle \ge \frac{1}{2\alpha} \|(\operatorname{Id} - T)x\|^2,$$

which is simply adding to both sides  $\|(\mathrm{Id} - T)(x)\|^2$ .

**Proposition B.5.** (See (Bauschke et al., 2021, Proposition 3.6(i))) Let  $A = T^{-1} - \operatorname{Id}$  and set  $N = \frac{1}{\alpha}T - \frac{1-\alpha}{\alpha}\operatorname{Id}$ , i.e.,  $T = J_A = (\operatorname{Id} + A)^{-1} = (1-\alpha)\operatorname{Id} + \alpha N$ . Then, T is  $\alpha$ -conically quasi-nonexpansive if and only if A is  $(1-\frac{1}{2\alpha})$ -star-cohypomonotone, i.e.,

$$\langle x - x^*, u \rangle \ge -\left(1 - \frac{1}{2\alpha}\right) \|u\|^2 \quad \forall (x, u) \in \operatorname{gra} A.$$

*Proof.* We see the two directions:

" $\Rightarrow$ " Let  $(x,u) \in \operatorname{gra} A$ . Then by definition of  $A = T^{-1} - \operatorname{Id}$  and manipulations, it follows that  $(x,u) = (T(x+u), (\operatorname{Id} - T)(x+u))$ . By Lemma B.3 invoked with  $x \leftarrow x + u$ , we have

$$2\alpha \langle T(x+u) - x^*, (\operatorname{Id} - T)(x+u) \rangle \ge (1 - 2\alpha) \| (\operatorname{Id} - T)(x+u) \|^2$$
  

$$\Leftrightarrow 2\alpha \langle x - x^*, u \rangle \ge (1 - 2\alpha) \| u \|^2,$$

where the last step substituted  $(x, u) = (T(x + u), (\mathrm{Id} - T)(x + u)).$ 

"\(\infty\)" Since  $(Tx, (\mathrm{Id} - T)x) \in \mathrm{gra} A$ , we have by star-cohypomonotonicity that  $\langle Tx - x^*, (\mathrm{Id} - T)x \rangle \geq \left(\frac{1}{2\alpha} - 1\right) \|(\mathrm{Id} - T)x\|^2$ . In view of Lemma B.3, we deduce conic quasi-nonexpansiveness.

# **B.2.** Complexity of the Outer Loop

**Bounding the norm of iterates.** Just like Appendix A, we start with the bound of the norms of the iterates.

**Lemma B.6.** Let Assumptions 1 and 3 hold. Suppose that the iterates  $(x_k)$  of Algorithm 3 satisfy  $||J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)|| \le \varepsilon_k$  for some  $\varepsilon_k > 0$  and  $\rho < \eta$ . Then, we have for  $k \ge 0$  that

$$||x_{k+1} - x^*|| \le ||x_k - x^*|| + \left(1 - \frac{\rho}{\eta}\right)\varepsilon_k.$$

*Proof.* From Fact B.2(ii), we know that  $J_{\eta(F+G)}$  is  $\frac{1}{2\left(1-\frac{\rho}{\eta}\right)}$ -conically quasi-nonexpansive. Then, by property (27) derived in Lemma B.3, since  $J_{\eta(F+G)}$  is also  $\frac{1}{1-\frac{\rho}{\eta}}$ -conically quasi-nonexpansive due to  $2\left(1-\frac{\rho}{\eta}\right)\geq 1-\frac{\rho}{\eta}$  (see also Corollary B.4), we have

$$\left\| \frac{\rho}{\eta} x_k + \left( 1 - \frac{\rho}{\eta} \right) J_{\eta(F+G)}(x_k) - x^* \right\| \le \|x_k - x^*\|. \tag{28}$$

By the definition of  $x_{k+1}$  in Algorithm 3, the definition of  $\varepsilon_k$  and triangle inequality, we have for  $k \geq 0$  that

$$||x_{k+1} - x^*|| \le \left\| \frac{\rho}{\eta} x_k + \left( 1 - \frac{\rho}{\eta} \right) J_{\eta(F+G)}(x_k) - x^* \right\| + \left( 1 - \frac{\rho}{\eta} \right) ||J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)||$$

$$\le \left\| \frac{\rho}{\eta} x_k + \left( 1 - \frac{\rho}{\eta} \right) J_{\eta(F+G)}(x_k) - x^* \right\| + \left( 1 - \frac{\rho}{\eta} \right) \varepsilon_k.$$

Combining with (28) gives the result.

**Iteration complexity.** Equipped with this result, we proceed to deriving the iteration complexity of the outer loop.

**Lemma 3.5.** Let Assumptions 1 and 3 hold. Suppose that the iterates  $(x_k)$  of Algorithm 3 satisfy  $||J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)|| \le \varepsilon_k$  for some  $\varepsilon_k > 0$  and  $\rho < \eta$ . Then, we have for  $K \ge 1$  that

$$\sum_{k=0}^{K-1} \| (\operatorname{Id} - J_{\eta(F+G)})(x_k) \|^2 \le \frac{2\eta^2}{(\eta - \rho)^2} \|x_0 - x^*\|^2 + 6 \sum_{k=0}^{K-1} \varepsilon_k^2 + \frac{4\eta}{\eta - \rho} \sum_{k=0}^{K-1} \|x_k - x^*\| \varepsilon_k, \tag{29}$$

where

$$||x_k - x^*|| \le ||x_{k-1} - x^*|| + \left(1 - \frac{\rho}{\eta}\right) \varepsilon_{k-1}.$$

*Proof.* From Fact B.2(ii), we have that  $\mathrm{Id}-J_{\eta(F+G)}$  is  $\left(1-\frac{\rho}{\eta}\right)$ -star cocoercive. Let us recall our running notations:

$$\alpha = 1 - \frac{\rho}{\eta}, \quad R = \operatorname{Id} - J_{\eta(F+G)}, \quad \widetilde{R} = \operatorname{Id} - \widetilde{J}_{\eta(F+G)}.$$

As a result, we have the following equivalent representation of  $x_{k+1}$  (see the definition in Algorithm 3):

$$x_{k+1} = x_k - \left(1 - \frac{\rho}{\eta}\right) \left(\text{Id} - \widetilde{J}_{\eta(F+G)}\right) (x_k)$$
$$= x_k - \alpha \widetilde{R}(x_k). \tag{30}$$

Then, by  $\alpha$ -star-cocoercivity of R, we have

$$\langle R(x_k), x_k - x^* \rangle \ge \alpha \|R(x_k)\|^2. \tag{31}$$

A simple decomposition gives

$$\langle R(x_k), x_k - x^* \rangle = \langle \widetilde{R}(x_k), x_k - x^* \rangle + \langle R(x_k) - \widetilde{R}(x_k), x_k - x^* \rangle. \tag{32}$$

We estimate the first term on the right-hand side of (32) as

$$\langle \widetilde{R}(x_{k}), x_{k} - x^{*} \rangle = \frac{1}{\alpha} \langle x_{k} - x_{k+1}, x_{k} - x^{*} \rangle$$

$$= \frac{1}{2\alpha} \left( \|x_{k} - x_{k+1}\|^{2} + \|x_{k} - x^{*}\|^{2} - \|x_{k+1} - x^{*}\|^{2} \right)$$

$$\leq \frac{1}{2\alpha} \left( \|x_{k} - x^{*}\|^{2} - \|x_{k+1} - x^{*}\|^{2} \right) + \frac{3\alpha}{4} \|R(x_{k})\|^{2} + \frac{3\alpha}{2} \|\widetilde{R}(x_{k}) - R(x_{k})\|^{2}$$

$$\leq \frac{1}{2\alpha} \left( \|x_{k} - x^{*}\|^{2} - \|x_{k+1} - x^{*}\|^{2} \right) + \frac{3\alpha}{4} \|R(x_{k})\|^{2} + \frac{3\alpha\varepsilon_{k}^{2}}{2}, \tag{33}$$

where we used the definition of  $x_{k+1}$  from (30) in the first step, standard expansion  $||a-b||^2 = ||a||^2 - 2\langle a,b\rangle + ||b||^2$  for the second step, the definition of  $x_{k+1}$  from (30) and Young's inequality in the third step, and the definitions of  $R_k$ ,  $\widetilde{R}_k$ ,  $\varepsilon_k$  in the last step.

For the second term on the right-hand side of (32), we have by Cauchy-Schwarz inequality and the definition of  $\widetilde{R}$  and  $\varepsilon_k$  that

$$\langle R(x_k) - \widetilde{R}(x_k), x_k - x^* \rangle \le ||R(x_k) - \widetilde{R}(x_k)|| ||x_k - x^*||$$

$$\le ||x_k - x^*|| \varepsilon_k.$$
(34)

We combine (33) and (34) in (32), plug in the result to (31) and rearrange to obtain

$$\frac{\alpha}{4} \|R(x_k)\|^2 \le \frac{1}{2\alpha} \left( \|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2 \right) + \frac{3\alpha \varepsilon_k^2}{2} + \|x_k - x^*\| \varepsilon_k.$$

The result follows by multiplying both sides by  $4/\alpha$ , summing for  $k=0,1,\ldots,K-1$ , and using the definition of  $\alpha$ . The bound on  $\|x_k-x^\star\|^2$  follows by Lemma B.6.

# **B.3.** Complexity of the Inner Loop

In a modular fashion, we will use precisely the same algorithm for the inner loop, i.e., the Forward-Backward-Forward (FBF) algorithm of (Tseng, 2000) like the Section A.3. Hence the complexity of the inner loop is the same as Theorem 2.7. As we see in the next section, the accuracy required by  $\widetilde{J}_{\eta(F+G)}$  is slightly different leading to the number of inner loop iterations  $N_k$  in Algorithm 3 to be slightly different than Algorithm 1.

# **B.4. Total Complexity**

**Theorem 3.1.** Let Assumptions 1 and 3 hold. Let  $\eta < \frac{1}{L}$  in Algorithm 3 and suppose that  $\rho < \eta$ . For any K > 1, we have that

$$\frac{1}{K} \sum_{k=0}^{K-1} \frac{1}{\eta^2} \|x_k - J_{\eta(F+G)}(x_k)\|^2 \le \frac{11 \|x_0 - x^\star\|^2}{(\eta - \rho)^2 K}.$$

The number of first-order oracles used at iteration k of Algorithm 3 is upper bounded by

$$\left[ \frac{4(1+\eta L)}{1-\eta L} \log(8(k+2)\log^2(k+2)) \right].$$

**Remark B.7.** It is straightforward to convert this to a last-iterate result if we additionally assume cohypomonotonicity as in (Pethick et al., 2023b), but we refrain from doing so since the main point of this section is to *relax* cohypomonotonicity.

*Proof of Theorem 3.1.* Recall the notations  $\alpha = 1 - \frac{\rho}{\eta}$  and  $R = \mathrm{Id} - J_{\eta(F+G)}$ . Let us set

$$\varepsilon_k = \frac{1}{8(k+1)\log^2(k+2)} \|x_k - J_{\eta(F+G)}(x_k)\|$$
(35)

and note, just as in the proof of Theorem 2.1, that we will not evaluate the value of  $\varepsilon_k$  but we will prove that for the number of iterations that FBF runs at each KM iteration in Algorithm 3, the error criterion required by  $\varepsilon_k$  is satisfied.

By using the definition of  $\varepsilon_k$  in Lemma B.6 gives

$$||x_{k+1} - x^*|| \le ||x_k - x^*|| + \frac{\alpha}{8(k+1)\log^2(k+2)} ||x_k - J_{\eta(F+G)}(x_k)||.$$
(36)

We note that  $\alpha=1-\frac{\rho}{\eta}\leq 1$  and since  $R=\mathrm{Id}-J_{\eta(F+G)}$  is  $\alpha$ -star cocoercive as shown in Fact B.2(ii), we have that  $\mathrm{Id}-J_{\eta(F+G)}$  is  $\alpha^{-1}$ -star Lipschitz and hence by  $(\mathrm{Id}-J_{\eta(F+G)})(x^\star)=0$  we have

$$\|(\operatorname{Id} - J_{\eta(F+G)})(x_k)\| = \|(\operatorname{Id} - J_{\eta(F+G)})(x_k) - (\operatorname{Id} - J_{\eta(F+G)})(x^*)\| \le \alpha^{-1} \|x_k - x^*\|.$$
(37)

Consequently, (36) becomes, after summing for  $k = 0, 1, \dots, K - 1$  that

$$||x_K - x^*|| \le ||x_0 - x^*|| + \sum_{i=0}^{K-1} \frac{1}{8(i+1)\log^2(i+2)} ||x_i - x^*||.$$

With this, we can show by induction that

$$||x_k - x^*|| \le 2||x_0 - x^*|| \quad \forall k > 0,$$
 (38)

because  $\sum_{i=0}^{\infty} \frac{1}{(i+1)\log^2(i+2)} < 4$ .

We use (38) in the result of Lemma 3.5 to obtain (also noting the definitions of  $\alpha$  and R)

$$\sum_{k=0}^{K-1} \|R(x_k)\|^2 \le \frac{2}{\alpha^2} \|x_0 - x^*\|^2 + \sum_{k=0}^{K-1} 6\varepsilon_k^2 + \frac{4}{\alpha} \sum_{k=0}^{K-1} 2\|x_0 - x^*\|\varepsilon_k.$$
 (39)

By using (38) and (37) in (35) we also know the following upper bound on  $\varepsilon_k$ :

$$\varepsilon_k \le \frac{\|x_0 - x^\star\|}{4\alpha(k+1)\log^2(k+2)}.$$

With this, (39) becomes

$$\sum_{k=0}^{K-1} \|R(x_k)\|^2 \le \frac{2}{\alpha^2} \|x_0 - x^*\|^2 + \sum_{k=0}^{K-1} \frac{3\|x_0 - x^*\|^2}{8\alpha^2 (k+1)^2 \log^4 (k+2)} + \sum_{k=0}^{K-1} \frac{2\|x_0 - x^*\|^2}{\alpha^2 (k+1) \log^2 (k+2)} < \frac{11}{\alpha^2} \|x_0 - x^*\|^2,$$
(40)

since  $\sum_{k=0}^{K-1} \frac{3}{8(k+1)^2 \log^4(k+2)} < 2$  and  $\sum_{k=0}^{K-1} \frac{2}{(k+1) \log^2(k+1)} < 7$ . This establishes the first part of the assertion.

We next see that, with  $N_k$  set as in Algorithm 3, we get the inexactness level specified by  $\varepsilon_k$  in (35) and we verify that the oracle complexity of each iteration is as claimed in the statement.

For the second part of the result, we proceed similar to the proof of Theorem 2.1. Namely, at iteration k, we apply the result in Theorem 2.7. For this, let us identify the following from the definitions in Algorithm 3

$$A \equiv \eta G, \quad B(\cdot) \equiv (\operatorname{Id} + \eta F)(\cdot) - x_k, \quad z_0 \equiv x_k, \quad z^* \equiv J_{\eta(F+G)}(x_k), \quad \zeta \equiv \varepsilon_k$$
  
$$\implies z_0 - z^* = (\operatorname{Id} - J_{\eta(F+G)})(x_k) = R(x_k).$$

As before, we have that B is  $(1 + \eta L)$ -Lipschitz and  $(1 - \eta L)$ -strongly monotone due to Fact B.2(iv). Existence of  $z^*$  is guaranteed by Fact B.2(iii) since  $\eta < \frac{1}{L}$ .

We now see that by the setting of  $N_k$  from Algorithm 3 and definition of  $\varepsilon_k$  in (35), we have

$$N_k = \left\lceil \frac{4(1+\eta L)}{1-\eta L} \log(8(k+1)\log^2(k+2)) \right\rceil = \left\lceil \frac{4(1+\eta L)}{1-\eta L} \log \frac{\|R(x_k)\|}{\varepsilon_k} \right\rceil.$$

With this value, Theorem 2.7 gives us

$$\|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\| \le \varepsilon_k$$

as claimed.

Since each iteration of Algorithm 2 uses 2 evaluations of F and 1 resolvent for G, the number of first-order oracle calls at iteration k is  $2N_k$  and the result follows.

We continue with the proof of Corollary 3.2 which follows trivially from Theorem 3.1.

*Proof of Corollary* 3.2. Based on Theorem 3.1, we have that after K iterations where

$$K \le \left\lceil \frac{11\|x_0 - x^*\|^2}{\eta^2 \alpha^2 \varepsilon^2} \right\rceil \tag{41}$$

we are guaranteed to obtain

$$\min_{0 \le k \le K - 1} \eta^{-1} ||R(x_k)|| \le \frac{1}{K} \sum_{k=0}^{K-1} \eta^{-1} ||R(x_k)|| \le \varepsilon.$$

Total number of first-oracle calls during the run of the algorithm then be calculated as

$$\sum_{k=1}^{K} \left\lceil \frac{4(1+\eta L)}{1-\eta L} \log(8(k+2) \log^2(k+2)) \right\rceil \leq K \cdot \left( \frac{4(1+\eta L)}{1-\eta L} \log(8(K+2) \log^2(K+2)) + 1 \right).$$

We conclude after using (41).

#### **B.5. Additional Results**

Let us re-emphasize the strategy in the previous proof: we set a target value for  $\varepsilon_k$  and then we prove that when we run the inner algorithm FBF for a certain *computable* number of iterations  $N_k$ , the criterion enforced on  $\widetilde{J}_{\eta(F+G)}$  by  $\varepsilon_k$  is satisfied. However, this number of inner iterations is *worst-case*. Another alternative, which could be more useful in practice is to set  $\varepsilon_k$  to a computable value and monitor the progress of the inner algorithm and break when  $\varepsilon_k$  is attained. One sidenote is that this is attainable in the deterministic case considered in this section, however it cannot be done in the stochastic case since the convergence guarantees are generally given in expectation.

This described strategy can be made rigorous with slight changes in the constants in our deterministic case. We now see this in the next proposition.

**Corollary B.8.** Let Assumptions 1 and 3 hold and let  $G = \partial \iota_C$  for a convex closed set C. Let  $\eta < \frac{1}{L}$  and  $\rho < \eta$  in Algorithm 1 with  $\|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\| \le \frac{c}{k \log^2(k+2)}$  for any c > 0 and use (Diakonikolas, 2020, Algorithm 4) to obtain such  $\widetilde{J}_{\eta(F+G)}(x_k)$  at iteration k. Then, we have that

$$\frac{1}{K} \sum_{k=0}^{K-1} \eta^{-1} \| (\operatorname{Id} - J_{\eta(F+G)})(x_k) \| \le \varepsilon,$$

with the number of calls to evaluation of F and resolvent of G is bounded by  $\tilde{O}\left(\frac{(1+\eta L)((1+c)\|x_0-x^\star\|^2+c^2)}{\varepsilon^2(\eta-\rho)^2(1-\eta L)}\right)$ .

**Remark B.9.** Note that (Diakonikolas, 2020, Algorithm 4) has a built-in stopping criterion to terminate the algorithm when the required accuracy is achieved. The value for  $\varepsilon_k$  defined in this corollary is computable since it only depends on k and a user-defined constant c. This is an alternative to FBF we used in the main text where we use a computable number of iterations to run the inner algorithm rather than using a stopping criterion as (Diakonikolas, 2020, Algorithm 4). On the one hand, in practice, a stopping criterion can be more desirable since the worst-case number of iterations can be pessimistic. On the other hand, the strategy of using a stopping criterion is inherently more complicated in the stochastic case whereas using a worst-case computable number of inner iteration is still easily implementable. This is why we considered the latter setting throughout the paper. However, this corollary is still included for the former strategy.

Proof of Corollary B.8. We obtain the result by modifying the proof of Theorem 3.1. We set

$$\varepsilon_k = \frac{c}{(k+1)\log^2(k+2)},$$

for any c > 0.

By using this on Lemma B.6 and summing the result for  $k = 0, 1, \dots, K-1$  we obtain

$$||x_k - x^*|| \le ||x_0 - x^*|| + \alpha \sum_{k=0}^{K-1} \frac{c}{(k+1)\log^2(k+2)}$$
  
$$\le ||x_0 - x^*|| + 4\alpha c,$$

since  $\sum_{k=0}^{K-1} \frac{1}{(k+1)\log^2(k+2)} < 4$ . We use this bound on the result of Lemma 3.5 to obtain

$$\sum_{k=0}^{K-1} \|R(x_k)\|^2 \le \frac{2}{\alpha^2} \|x_0 - x^\star\|^2 + \sum_{k=0}^{K-1} \frac{6c^2}{(k+1)^2 \log^4(k+2)} + \frac{4}{\alpha} \sum_{k=0}^{K-1} \frac{c(\|x_0 - x^\star\| + 4\alpha c)}{(k+1) \log^2(k+2)}$$
$$\le \left(\frac{2}{\alpha^2} + \frac{16c}{\alpha}\right) \|x_0 - x^\star\|^2 + 30c^2 + 64c^2,$$

which gives the result after dividing by  $\eta^2$  and noting that (Diakonikolas, 2020, Lemma 17) gives complexity  $\widetilde{O}\left(\frac{1+\eta L}{1-\eta L}\right)$  for obtaining such a  $\widetilde{J}_{\eta(F+G)}(x_k)$  with (Diakonikolas, 2020, Algorithm 4).

We continue with the result mentioned in Remark 3.4.

**Corollary B.10.** Let Assumptions 1 and 3 hold. Let  $\eta < \frac{1}{L}$  and  $\rho < \eta$ .

(i) Let  $G \equiv 0$  and consider Algorithm 3. Then we have that  $\min_{0 \le k \le K-1} ||F(x_k)|| \le 2\varepsilon$  with the first-order oracle calls bounded by

$$\widetilde{O}\left(\frac{(1+\eta L)\|x_0 - x^{\star}\|^2}{\varepsilon^2 (\eta - \rho)^2 (1 - \eta L)}\right).$$

(ii) Let  $G \equiv \partial \iota_C$  for a convex closed set  $C \subseteq \mathbb{R}^d$ . Given  $\varepsilon > 0$ , consider Algorithm 3 with the update  $\widetilde{J}_{\eta(F+G)}(x_k)$  replaced with (Diakonikolas, 2020, Algorithm 4) with error criterion  $\|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\| \le \frac{\eta \varepsilon^2}{(k+1)\log^2(k+3)}$ . Then, for  $x^{out} = \arg\min_{x \in \{x_0, \dots, x_{k-1}\}} \|x - \widetilde{J}_{\eta(F+G)}(x)\|$ , we have that  $\eta^{-1}\|(\operatorname{Id} - J_{\eta(F+G)})(x^{out})\| \le 2\varepsilon + 3\varepsilon^2$  with the first-order oracle calls bounded by

$$\widetilde{O}\left(\frac{(1+\eta L)\|x_0 - x^\star\|^2}{\varepsilon^2 (\eta - \rho)^2 (1 - \eta L)}\right). \tag{42}$$

See also Remark 2.3 for details on how we can use this result to further obtain a guarantee like  $\operatorname{dist}(0,(F+G)(x^{\operatorname{out}})) \leq \varepsilon$ .

*Proof of Corollary B.10.* (i) In this case, we start from the final steps of the proof of Theorem 3.1 (see (40)) which, after using  $R = \text{Id} - J_{\eta(F+G)}$ , gives us that

$$\frac{1}{K} \sum_{k=0}^{K-1} \eta^{-2} \|x_k - J_{\eta F}(x_k)\|^2 \le \varepsilon^2, \tag{43}$$

with the prescribed complexity bound given in Corollary 3.2. Let us define  $\bar{x}_k = J_{\eta F}(x_k)$ .

On the one hand, we use the definition of resolvent to obtain

$$\bar{x}_k = J_{\eta F}(x_k) \iff \bar{x}_k + \eta F(\bar{x}_k) = x_k \iff x_k - \bar{x}_k = \eta F(\bar{x}_k),$$

which, in view of (43), means that we have

$$\frac{1}{K} \sum_{k=0}^{K-1} ||F(\bar{x}_k)||^2 \le \varepsilon^2.$$
 (44)

On the other hand, we know by Young's inequality and Lipschitzness of F that

$$\frac{1}{K} \sum_{k=0}^{K-1} \|F(x_k)\|^2 \le \frac{1}{K} \sum_{k=0}^{K-1} 2 \|F(\bar{x}_k)\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} 2 \|F(x_k) - F(\bar{x}_k)\|^2 
\le \frac{1}{K} \sum_{k=0}^{K-1} 2 \|F(\bar{x}_k)\|^2 + \frac{1}{K} \sum_{k=0}^{K-1} 2 L^2 \|x_k - \bar{x}_k\|^2 
\le (2 + 2\eta^2 L^2) \varepsilon^2 
< 4\varepsilon^2.$$

where we used (43) and (44).

(ii) A slight modification of the proof of Corollary B.8 by using  $\varepsilon_k = \frac{\eta \varepsilon^2}{(k+1) \log^2(k+3)}$  gives us that

$$\frac{1}{K} \sum_{k=0}^{K-1} \eta^{-2} \| (\operatorname{Id} - J_{\eta(F+G)})(x_k) \|^2 \le \varepsilon^2$$
 (45)

with the complexity bound (42). This is because (Diakonikolas, 2020, Lemma 17) showed that (Diakonikolas, 2020, Algorithm 4) outputs a  $\widetilde{J}_{\eta(F+G)}(x_k)$  satisfying the requirement set by  $\varepsilon_k = \frac{\eta \varepsilon^2}{(k+1)\log^2(k+3)}$ , with the same worst-case complexity as Theorem 2.7. The difference is that (Diakonikolas, 2020, Algorithm 4) has a computable stopping criterion (instead of the maximum number of iterations Algorithm 2 takes) where we can check if  $\varepsilon_k = \frac{\eta \varepsilon^2}{(k+1)\log^2(k+3)}$  accuracy is achieved and break the loop.

Since we have the pointwise bound  $||J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)||^2 \le \eta^2 \varepsilon^4$ , we derive from (45) that

$$\frac{1}{K} \sum_{k=0}^{K-1} \eta^{-2} \| (\operatorname{Id} - \widetilde{J}_{\eta(F+G)})(x_k) \|^2 \le 2(\varepsilon^2 + \varepsilon^4).$$

Hence, for  $x^{\text{out}}$  defined in the statement, we get

$$\eta^{-2} \| (\operatorname{Id} - \widetilde{J}_{\eta(F+G)})(x^{\operatorname{out}}) \|^2 \le 2(\varepsilon^2 + \varepsilon^4).$$
(46)

Then, by using the pointwise bound  $||J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)||^2 \le \eta^2 \varepsilon^4$  for all k, we know that

$$\eta^{-1} \| (\operatorname{Id} - J_{\eta(F+G)})(x^{\operatorname{out}}) \| \leq \eta^{-1} \| (\operatorname{Id} - \widetilde{J}_{\eta(F+G)})(x^{\operatorname{out}}) \| + \eta^{-1} \| (J_{\eta(F+G)} - \widetilde{J}_{\eta(F+G)})(x^{\operatorname{out}}) \|$$

$$\leq \varepsilon^{2} + \sqrt{2(\varepsilon^{4} + \varepsilon^{2})} < 2\varepsilon + 3\varepsilon^{2},$$

which uses (46) and the implication of the error criterion  $||J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)|| \le \eta \varepsilon^2$ , completing the proof.

Assumption	Reference	Limit of $\rho$	Constraints	Oracle <sup>†</sup>	Complexity
weak MVI	(Diakonikolas et al., 2021)	$\frac{1}{4\sqrt{2}L}$	×	Single	$O(\varepsilon^{-4})$
	(Choudhury et al., 2023)	$\frac{1}{2L}$	×	Single	$O(\varepsilon^{-4})$
	(Böhm, 2022)	$\frac{1}{2L}$	×	Single	$O(\varepsilon^{-6})$
	(Pethick et al., 2023a)	$\frac{1}{2L}$	$\checkmark$	Multiple	$\widetilde{O}(arepsilon^{-4})$
	Theorem C.11	$\frac{1}{L}$	$\checkmark$	Single	$\widetilde{O}(\varepsilon^{-4})$
cohypomonotone	(Pethick et al., 2023b)	$\frac{1}{2L}$	✓	Single	$\widetilde{O}(\varepsilon^{-6})  (\mathrm{best})^{\ddagger}$
	(Pethick et al., 2023b)	$\frac{1}{2L}$	$\checkmark$	Single	$\widetilde{O}(arepsilon^{-16})$ (last)
	(Chen & Luo, 2022)*	$\frac{1}{2L}$	×	Single	$\widetilde{O}(\varepsilon^{-2})$
	Corollary C.5*	$\frac{1}{L}$	$\checkmark$	Single	$\widetilde{O}(\varepsilon^{-4})$

Table 2. Comparison of first order algorithms for stochastic problems. Complexity refers to the number of oracle calls to get the fixed point residual  $\mathbb{E}\|(\mathrm{Id}-J_{\eta(F+G)})(x^{\mathrm{out}})\| \leq \varepsilon$ . See also Remark 2.3. <sup>†</sup>Oracle access refers to the number of operator evaluations algorithm makes with one random seed given  $F(x) = \mathbb{E}_{\xi \sim \Xi}[F_{\xi}(x)]$ . For example, "Single" refers to algorithms that only access one sample per seed, i.e., only  $F_{\xi_t}(x_t)$ , "Multiple" is for algorithms that access multiple samples per seed, i.e.,  $F_{\xi_t}(x_t)$  and  $F_{\xi_t}(x_{t-1})$ . Algorithms with "Multiple" access also make the additional assumption that  $\mathbb{E}_{\xi \sim \Xi} \|F_{\xi}(x) - F_{\xi}(y)\|^2 \leq L^2 \|x - y\|^2$  which is stronger than mere Lipschitzness of F. <sup>‡</sup>(best) refers to best iterate in view of Remark 3.3; (last) refers to a last iterate convergence rate. \*These works have complexity as expected number of oracle calls due to the use of MLMC estimator. See also Appendix D.1 for derivations of the complexities when they are not written explicitly in the existing works.

#### C. Proofs for Section 4

**Notation.** We use the following definitions for conditional expectations: For expectation conditioned on the filtration generated by the randomness of  $x_k, \ldots, x_1$ , we use  $\mathbb{E}_k[\cdot]$  while analyzing Algorithm 4 and Algorithm 6. In the notation of Algorithm 5, we similarly use  $\mathbb{E}_{t+1/2}[\cdot]$  for the expectation conditioned on the filtration generated by the randomness of  $z_{t+1/2}, z_t, \ldots, z_1, z_{1/2}$ . Unif denotes the uniform distribution and Geom denotes the geometric distribution.

Table 2 summarizes the existing works for stochastic min-max problems satisfying cohypomonotonicity or weak MVI conditions.

#### C.1. Analysis of the inner loop for stochastic problems

The main change for algorithms in the stochastic case is computing the resolvent approximation  $\widetilde{J}_{\eta(F+G)}(x_k)$ . We now need to invoke FBF with unbiased oracles for F, see for example (9). For ease of reference, we specify the algorithm below. Note that Algorithm 4 is precisely Algorithm 1 when (9) is used for estimating the resolvent and Algorithm 5 is precisely Algorithm 2 when unbiased oracle  $\widetilde{B}$  is inputted rather than full operator B. Algorithm 5 is a stochastic version of FBF, which is analyzed in the monotone case by (Böhm et al., 2022).

# Algorithm 4 Stochastic Inexact Halpern iteration for problems with cohypomonotonicity

Input: Parameters 
$$\beta_k = \frac{1}{k+2}, \eta, L, \rho, \alpha = 1 - \frac{\rho}{\eta}, K > 0$$
, initial iterate  $x_0 \in \mathbb{R}^d$ , subroutine FBF given in Algorithm 5 for  $k = 0, 1, 2, \ldots, K - 1$  do 
$$\widetilde{J}_{\eta(F+G)}(x_k) = \text{FBF}\left(x_k, N_k, \eta G, \text{Id} + \eta \widetilde{F}, 1 + \eta L\right), \text{ where } N_k = \left\lceil \frac{1734(k+2)^3 \log^2(k+2)}{(1-\eta L)^2} \right\rceil$$
 
$$x_{k+1} = \beta_k x_0 + (1-\beta_k)((1-\alpha)x_k + \alpha \widetilde{J}_{\eta(F+G)}(x_k))$$
 end for

More particularly, we solve the following stochastic strongly monotone inclusion problem:

Find 
$$x^* \in \mathbb{R}^d$$
 such that  $0 \in (A+B)(x^*)$ , where  $B = \mathbb{E}_{\xi \sim \Xi}[B_{\xi}]$ .

Similar results to next theorem appeared in (Hsieh et al., 2019; Kotsalis et al., 2022). We provide a proof for being complete and precise since we could not find a particular reference for stochastic FBF with strong monotonicity and explicit constants.

**Algorithm 5** FBF $(z_0, T, A, \widetilde{B}_{\rm in}, L_B)$  from (Tseng, 2000) – Stochastic

Input: Parameter 
$$\tau_t = \frac{2}{(t+1)\mu+6L_B}$$
, initial iterate  $z_0 \in \mathbb{R}^d$ ,  $\widetilde{B}(\cdot) = \widetilde{B}_{\rm in}(\cdot) - z_0$  for  $t = 0, 1, 2, \dots, N-1$  do  $z_{t+1/2} = J_{\tau_t A}(z_t - \tau_t \widetilde{B}(z_t))$ 

$$z_{t+1} = z_{t+1/2} + \tau_t \widetilde{B}(z_t) - \tau_t \widetilde{B}(z_{t+1/2})$$

end for

It is also worth noting that we do not focus on optimizing the non-dominant terms. A tight bound for all the terms can be found in (Kotsalis et al., 2022) who analyzed a different algorithm.

**Theorem C.1.** Let  $z^* = (A+B)^{-1}(0) \neq \emptyset$ , the operator B be  $L_B$ -Lipschitz and  $\mu$ -strongly monotone with  $\mu > 0$ , A be maximally monotone. Let  $\widetilde{B} \colon \mathbb{R}^d \to \mathbb{R}^d$  satisfy  $\mathbb{E}[\widetilde{B}(x)] = B(x)$  and  $\mathbb{E}\|\widetilde{B}(x) - B(x)\|^2 \leq \sigma^2$ . Then, we have that the last iterate of Algorithm 5 after running for T iterations, when initialized with  $z_0$ , and step size  $\tau_t = \frac{2}{(t+1)\mu + 6L_B}$  satisfies the bound

$$\mathbb{E}||z_T - z^*||^2 \le \frac{6L_B/\mu ||z_0 - z^*||^2 + 48\sigma^2/\mu^2}{T + 6L_B/\mu}.$$

Each iteration of the algorithm uses two evaluations of  $\widetilde{B}$  and one resolvent of A

*Proof.* Note that the definition of  $z_{t+1/2}$  implies  $\tau_t A(z_{t+1/2}) \ni z_t - z_{t+1/2} - \tau_t B_{\xi_t}(z_t)$ . The definition of  $z^*$  implies  $\tau_t A(z^*) \ni -\tau_t B(z^*)$  By using this with monotonicity of A, we get

$$\langle z_{t+1/2} - z_t + \tau_t \widetilde{B}(z_t) - \tau_t B(z^*), z^* - z_{t+1/2} \rangle \ge 0.$$

By the definition of  $z_{t+1}$ , we then have

$$\langle z_{t+1} - z_t + \tau_t \widetilde{B}(z_{t+1/2}) - \tau_t B(z^*), z^* - z_{t+1/2} \rangle \ge 0.$$
 (47)

By taking expectation conditioned on  $z_{t+1/2}$  and also using strong monotonicity of B, we also have

$$\mathbb{E}_{t+1/2}\langle \tau_{t}\widetilde{B}(z_{t+1/2}) - \tau_{t}B(z^{*}), z_{t+1/2} - z^{*}\rangle = \langle \tau_{t}B(z_{t+1/2}) - \tau_{t}B(z^{*}), z_{t+1/2} - z^{*}\rangle$$

$$\geq \mu \tau_{t} \|z^{*} - z_{t+1/2}\|^{2}$$

$$\geq \frac{\mu \tau_{t}}{2} \|z^{*} - z_{t+1}\|^{2} - \mu \tau_{t} \|z_{t+1} - z_{t+1/2}\|^{2}$$

$$\geq \frac{\mu \tau_{t}}{2} \|z^{*} - z_{t+1}\|^{2} - \frac{1}{3} \|z_{t+1} - z_{t+1/2}\|^{2}, \tag{48}$$

where the third step is by Young's inequality and last step is by the definition of  $\tau_t$ , i.e.,  $\tau_t \mu = \frac{2\mu}{(t+1)\mu + 6L_B} \le \frac{2\mu}{6L_B} \le \frac{1}{3}$  since  $\mu \le L_B$ .

We have, by the elementary identities  $\langle a,b \rangle = \frac{1}{2} \left( \|a\|^2 + \|b\|^2 - \|a-b\|^2 \right) = \frac{1}{2} \left( -\|a\|^2 - \|b\|^2 + \|a+b\|^2 \right)$ , that

$$\langle z_{t+1} - z_t, z^* - z_{t+1/2} \rangle = \langle z_{t+1} - z_t, z^* - z_{t+1} \rangle + \langle z_{t+1} - z_t, z_{t+1} - z_{t+1/2} \rangle$$

$$= \frac{1}{2} \left( \|z_t - z^*\|^2 - \|z_{t+1} - z^*\|^2 - \|z_t - z_{t+1/2}\|^2 + \|z_{t+1} - z_{t+1/2}\|^2 \right). \tag{49}$$

Using (48) and (49) on (47) after taking total expectation, using tower rule and dividing both sides by  $\tau_t$  gives

$$\left(\frac{1}{2\tau_t} + \frac{\mu}{2}\right) \mathbb{E}\|z^* - z_{t+1}\|^2 \le \frac{1}{2\tau_t} \mathbb{E}\|z^* - z_t\|^2 + \frac{5}{6\tau_t} \mathbb{E}\|z_{t+1} - z_{t+1/2}\|^2 - \frac{1}{2\tau_t} \mathbb{E}\|z_t - z_{t+1/2}\|^2.$$
 (50)

Definition of  $z_{t+1}$  in Algorithm 5 gives

$$\begin{split} \frac{5}{6}\|z_{t+1} - z_{t+1/2}\| &= \frac{5\tau_t^2}{6}\|\widetilde{B}(z_t) - \widetilde{B}(z_{t+1/2})\|^2 \\ &\leq \frac{5\tau_t^2}{2}\left(\|\widetilde{B}(z_t) - B(z_t)\|^2 + \|B(z_t) - B(z_{t+1/2})\|^2 + \|B(z_{t+1/2}) - \widetilde{B}(z_{t+1/2})\|^2\right) \\ &\leq 5\tau_t^2\sigma^2 + \frac{5\tau_t^2L_B^2}{2}\|z_t - z_{t+1/2}\|^2, \end{split}$$

where the last line is by the variance bound assumed on  $\widetilde{B}$  and Lipschitzness of B.

With this, we get in place of (50) that

$$\left(\frac{1}{2\tau_t} + \frac{\mu}{2}\right) \mathbb{E}\|z^* - z_{t+1}\|^2 \le \frac{1}{2\tau_t} \mathbb{E}\|z^* - z_t\|^2 + \frac{1}{\tau_t} \left(\frac{5\tau_t^2 L_B^2}{2} - \frac{1}{2}\right) \mathbb{E}\|z_t - z_{t+1/2}\|^2 + 5\tau_t \sigma^2.$$
 (51)

The definition of  $\tau_t = \frac{2}{(t+1)\mu + 6L_B}$  has two consequences:

$$\frac{1}{2\tau_t} + \frac{\mu}{2} = \frac{6L_B + (t+3)\mu}{4}$$

and

$$\tau_t = \frac{2}{(t+1)\mu + 6L_B} \le \frac{1}{3L_B} \Longrightarrow \tau_t^2 \le \frac{1}{5L_B^2} \iff 5\tau_t^2 L_B^2 \le 1.$$

This last estimate shows that the second term on the right-hand side of (51) is nonpositive.

Then, we obtain, after multiplying both sides of (51) by  $\left(\frac{1}{2\tau_t} + \frac{\mu}{2}\right)^{-1} = \frac{4}{6L_B + (t+3)\mu}$  that

$$\mathbb{E}\|z^{\star} - z_{t+1}\|^{2} \le \left(\frac{(t+1)\mu + 6L_{B}}{(t+3)\mu + 6L_{B}}\right)\|z^{\star} - z_{t}\|^{2} + \frac{40\sigma^{2}}{(6L_{B} + (t+1)\mu)(6L_{B} + (t+3)\mu)}.$$
 (52)

We next show by induction that

$$\mathbb{E}||z^* - z_t||^2 \le \frac{6L_B/\mu ||z_0 - z^*||^2 + 48\sigma^2/\mu^2}{t + 6L_B/\mu} \quad \forall t \ge 0.$$

For brevity, let us denote  $\kappa = 6L_B/\mu$ .

The base case t=0 holds by inspection. Next we assume the assertion holds for t=T and consider (52) to deduce

$$\mathbb{E}\|z^{\star} - z_{T+1}\|^{2} \le \frac{T + 1 + \kappa}{T + 3 + \kappa} \frac{\kappa \|z_{0} - z^{\star}\|^{2} + 48\sigma^{2}/\mu^{2}}{T + \kappa} + \frac{40\sigma^{2}/\mu^{2}}{(T + 1 + \kappa)(T + 3 + \kappa)} \\ = \left(\frac{(T + 1 + \kappa)}{(T + 3 + \kappa)(T + \kappa)} + \frac{1}{1.2(T + 1 + \kappa)(T + 3 + \kappa)}\right) \left(\kappa \|z_{0} - z^{\star}\|^{2} + 48\sigma^{2}/\mu^{2}\right).$$

As a result, the inductive step will be implied by

$$\left(\frac{(T+1+\kappa)^2}{(T+1+\kappa)(T+3+\kappa)(T+\kappa)} + \frac{(T+\kappa)}{1.2(T+1+\kappa)(T+3+\kappa)(T+\kappa)}\right) \le \frac{1}{T+1+\kappa},$$

which, after letting  $\nu = T + \kappa$ , is equivalent to

$$\left(\frac{1.2(\nu+1)^2}{(\nu+3)\nu} + \frac{\nu}{(\nu+3)\nu}\right) \le 1.2 \iff 1.2(\nu+1)^2 \le 1.2\nu^2 + 2.6\nu \iff 1.2 \le 0.2\nu \iff 6 \le \nu.$$

This holds because  $\nu = T + \kappa = T + 6L_B/\mu \ge 6$  since  $L_B/\mu > 1$ . This completes the induction.

#### C.2. Stochastic Problem with Cohypomonotonicity

We have a stochastic version of Lemma A.2 proof of which is almost equivalent.

**Lemma C.2.** Let Assumptions 1 and 2 hold. For the sequence  $(x_k)$  generated by Algorithm 4 with  $\mathbb{E}_k ||J_{\eta(F+G)}(x_k) - \widetilde{J}_{\eta(F+G)}(x_k)||^2 \le \varepsilon_k^2$ , we have for  $k \ge 0$  that

$$\mathbb{E}||x_{k+1} - x^*|| \le ||x_0 - x^*|| + \frac{\alpha}{k+2} \sum_{i=0}^k (i+1) \mathbb{E}[\varepsilon_i].$$

*Proof.* The proof follows the same steps as Lemma A.2 after taking expectation on (11) and using Jensen's inequality since

$$\mathbb{E}_k\left[\|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\|\right] \le \sqrt{\mathbb{E}_k\left[\|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\|^2\right]} \le \varepsilon_k.$$

Hence the result follows by tower rule and the same induction as the proof of Lemma A.2.

**Lemma C.3.** Let Assumptions 1 and 2 hold. Consider Algorithm 4 with  $\mathbb{E}_k \|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\|^2 \le \varepsilon_k^2$ . Then, we have for any  $\gamma > 0$  and  $K \ge 1$  that

$$\frac{\alpha K(K+1)}{4} \mathbb{E} \|R(x_K)\|^2 \le \frac{K+1}{K\alpha} \|x^* - x_0\|^2 + \sum_{k=0}^{K-1} \left( \frac{\alpha(\gamma+1)}{2\gamma} (k+1)(k+2) \mathbb{E}[\varepsilon_k^2] + \frac{\gamma \alpha \mathbb{E} \|R(x_k)\|^2}{2} \right).$$

*Proof.* We follow the proof of Lemma 2.5 until (17) and then we take expectation to obtain

$$\frac{\alpha}{2} \mathbb{E} \|R(x_{k+1})\|^{2} + \frac{\beta_{k}}{1 - \beta_{k}} \mathbb{E} \langle R(x_{k+1}), x_{k+1} - x_{0} \rangle 
\leq \frac{\alpha}{2} (1 - 2\beta_{k}) \mathbb{E} \|R(x_{k})\|^{2} + \beta_{k} \mathbb{E} \langle R(x_{k}), x_{k} - x_{0} \rangle 
+ \alpha \mathbb{E} \langle R(x_{k+1}) - (1 - \beta_{k}) R(x_{k}), (\widetilde{J}_{\eta(F+G)} - J_{\eta(F+G)})(x_{k}) \rangle - \frac{\alpha}{2} \mathbb{E} \|R(x_{k+1}) - R(x_{k})\|^{2}.$$
(53)

We then consider (18) after taking expectation and using Cauchy-Schwarz, triangle and Young's inequalities to obtain

$$\alpha \mathbb{E} \langle R(x_{k+1}) - (1 - \beta_k) R(x_k), (\widetilde{J}_{\eta(F+G)} - J_{\eta(F+G)})(x_k) \rangle 
\leq \alpha \mathbb{E} \left[ (\|R(x_{k+1}) - R(x_k)\| + \beta_k \|R(x_k)\|) \|\widetilde{J}_{\eta(F+G)}(x_k) + J_{\eta(F+G)}(x_k)\| \right] 
\leq \frac{\alpha}{2} \mathbb{E} \|R(x_{k+1}) - R(x_k)\|^2 + \frac{\gamma \alpha \beta_k^2}{2} \mathbb{E} \|R(x_k)\|^2 + \frac{\alpha}{2} \left(1 + \frac{1}{\gamma}\right) \mathbb{E} \|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\|^2 
\leq \frac{\alpha}{2} \mathbb{E} \|R(x_{k+1}) - R(x_k)\|^2 + \frac{\gamma \alpha \beta_k^2}{2} \mathbb{E} \|R(x_k)\|^2 + \frac{\alpha}{2} \left(1 + \frac{1}{\gamma}\right) \mathbb{E}[\varepsilon_k^2], \tag{54}$$

where the last step also used the tower rule along with the definition of  $\varepsilon_k$ .

We then use the same arguments as those after (19) to get the result.

#### C.2.1. Proof for Corollary 4.1

Corollary 4.1 is essentially the summary of the results proven below.

**Theorem C.4.** Let Assumptions 1, 2, and 4 hold. Let  $\eta < \frac{1}{L}$  in Algorithm 4 and  $\rho < \eta$ . For any  $k \ge 1$ , we have that

$$\frac{1}{\eta^2} \mathbb{E} \|x_k - J_{\eta(F+G)}(x_k)\|^2 \le \frac{36(\|x_0 - x^*\|^2 + \sigma^2)}{(\eta - \rho)^2 k^2}.$$

The number of first-order oracles used at iteration k of Algorithm 4 is upper bounded by

$$2\left[\frac{1734(k+2)^3\log^2(k+2)}{(1-\eta L)^2}\right]. \tag{55}$$

**Corollary C.5.** Let Assumptions 1 and 2 hold. Let  $\eta < \frac{1}{L}$  in Algorithm 4 and  $\rho < \eta$ . For any  $\varepsilon > 0$ , we have that  $\mathbb{E}\left[\eta^{-1}\|x_k - J_{\eta(F+G)}(x_k)\|\right] \leq \varepsilon$  with stochastic first-order oracle complexity

$$\widetilde{O}\left(\frac{\|x_0 - x^*\|^4 + \sigma^4}{(\eta - \rho)^4 (1 - \eta L)^2 \varepsilon^4}\right)$$

*Proof.* This corollary immediately follows from Theorem C.4 by combining the number of outer iterations and the number of stochastic first-order oracle calls for each outer iteration.

**Remark C.6.** The complexity in the previous corollary has the same dependence on  $||x_0 - x^*||$ ,  $\sigma$  as (Pethick et al., 2023a; Bravo & Cominetti, 2024). As we see in the next remark, the dependence on  $(\eta - \rho)$  can be improved by using the knowledge of the target accuracy  $\varepsilon$  and the variance upper bound  $\sigma^2$  as done in (Diakonikolas et al., 2021; Lee & Kim, 2021; Chen & Luo, 2022).

**Remark C.7.** By using parameters depending on target accuracy  $\varepsilon$  and noise variance  $\sigma^2$ , we can improve the complexity to

$$\widetilde{O}\left(\frac{\|x_0 - x^*\|^2 \sigma^2}{(\eta - \rho)^2 (1 - \eta L)^2 \varepsilon^4}\right)$$

*Proof of Theorem C.4.* Let us set

$$\varepsilon_k^2 = \frac{\gamma^2(\alpha^2 ||R(x_k)||^2 + 8\sigma^2)}{\alpha^2 (k+2)^3 \log^2(k+2)}$$
(56)

and plug this in to the result of Lemma C.3 to obtain

$$\frac{\alpha K(K+1)}{4} \mathbb{E} \|R(x_K)\|^2 
\leq \frac{K+1}{K\alpha} \|x_0 - x^*\|^2 + \sum_{k=0}^{K-1} \mathbb{E} \left( \frac{(\gamma^2 + \gamma)(\alpha^2 \|R(x_k)\|^2 + 8\sigma^2)}{2\alpha(k+2)\log^2(k+2)} + \frac{\gamma\alpha \|R(x_k)\|^2}{2} \right) 
= \frac{K+1}{K\alpha} \|x_0 - x^*\|^2 + \sum_{k=0}^{K-1} \left( \frac{4(\gamma^2 + \gamma)\sigma^2}{\alpha(k+2)\log^2(k+2)} + \frac{\alpha(\gamma^2 + \gamma)\mathbb{E} \|R(x_k)\|^2}{2(k+2)\log^2(k+2)} + \frac{\gamma\alpha\mathbb{E} \|R(x_k)\|^2}{2} \right) 
< \frac{K+1}{K\alpha} \|x_0 - x^*\|^2 + \frac{12(\gamma^2 + \gamma)\sigma^2}{\alpha} + \sum_{k=0}^{K-1} \left( \frac{\alpha(\gamma^2 + \gamma)\mathbb{E} \|R(x_k)\|^2}{2(k+2)\log^2(k+2)} + \frac{\gamma\alpha\mathbb{E} \|R(x_k)\|^2}{2} \right), \tag{57}$$

since  $\sum_{k=0}^{K-1} \frac{1}{(k+2)\log^2(k+2)} < 3$ .

We now show by induction that

$$\mathbb{E}||R(x_k)||^2 \le \frac{36(||x_0 - x^*||^2 + \sigma^2)}{\alpha^2(k+1)^2}.$$

The base case for the induction with K = 0, 1 hold the same way as the proof of Theorem 2.1 where the only change is we use Lemma C.2 and the definition of  $\varepsilon_k$  in (56), see also (24).

Let us consider (57) for  $K \ge 2$  and assume that the assertion holds for  $k \le K - 1$ . We then have that

$$\frac{\alpha K(K+1)}{4} \mathbb{E} \|R(x_K)\|^2 \\ \leq \frac{2}{\alpha} \|x_0 - x^*\|^2 + \frac{12(\gamma^2 + \gamma)\sigma^2}{\alpha} + \sum_{k=0}^{K-1} \left( \frac{18(\gamma^2 + \gamma)(\|x_0 - x^*\|^2 + \sigma^2)}{\alpha(k+2)(k+1)^2 \log^2(k+2)} + \frac{18\gamma(\|x_0 - x^*\|^2 + \sigma^2)}{\alpha(k+1)^2} \right),$$

where we also used  $\frac{K+1}{K} \leq 2$ .

By using  $\sum_{k=0}^{\infty} \frac{18}{(k+1)^2} < 30$  and  $\sum_{k=0}^{\infty} \frac{18}{(k+2)(k+1)^2 \log^2(k+2)} < 21$  and  $\gamma = \frac{1}{17}$ , we have that

$$\frac{\alpha K(K+1)}{4} \mathbb{E} \|R(x_K)\|^2 \le \frac{6(\|x_0 - x^*\|^2 + \sigma^2)}{2}.$$

We use  $\frac{1}{K(K+1)} \leq \frac{1.5}{(K+1)^2}$  which holds for  $K \geq 2$  to complete the induction

To see the number of first-order oracles, we use the result for stochastic FBF in Theorem C.1. For our subproblem at iteration k, this result implies

$$\mathbb{E}_{k}\left[\|\widetilde{J}_{\eta(F+G)}(x_{k}) - J_{\eta(F+G)}(x_{k})\|^{2}\right] \leq \frac{6\left(\frac{1+\eta L}{1-\eta L}\|x_{k} - J_{\eta(F+G)}(x_{k})\|^{2} + \frac{8\sigma^{2}}{(1-\eta L)^{2}}\right)}{N_{k}}$$

$$\leq \frac{\frac{6}{(1-\eta L)^{2}}\left(\|x_{k} - J_{\eta(F+G)}(x_{k})\|^{2} + 8\sigma^{2}\right)}{N_{k}}.$$

Recall that (56), with  $\gamma = \frac{1}{17}$  and  $R = \mathrm{Id} - J_{\eta(F+G)}$ , requires

$$\mathbb{E}_{k}\left[\|\widetilde{J}_{\eta(F+G)}(x_{k}) - J_{\eta(F+G)}(x_{k})\|^{2}\right] \leq \frac{(\alpha^{2}\|(\mathrm{Id} - J_{\eta(F+G)})(x_{k})\|^{2} + 8\sigma^{2})}{289\alpha^{2}(k+2)^{3}\log^{2}(k+2)}$$

Noting that  $\frac{1}{\alpha^2} > 1$ , a sufficient condition to attain this requirement is

$$\frac{\frac{6}{(1-\eta L)^2} \left( \|x_k - J_{\eta(F+G)}(x_k)\|^2 + 8\sigma^2 \right)}{N_k} \le \frac{\|x_k - J_{\eta(F+G)}(x_k)\|^2 + 8\sigma^2}{289(k+2)^3 \log^2(k+2)},$$

verifying the required number of iterations  $N_k$  as given in Algorithm 4 to be sufficient for the inexactness criterion. Since each iteration of FBF takes 2 stochastic operator evaluations  $\widetilde{F}$  and one resolvent of G, we have the result.

#### C.3. Stochastic Problem with weak MVI condition

As motivated in Section 4.2, we use the multilevel Monte Carlo (MLMC) estimator (Giles, 2008; Blanchet & Glynn, 2015; Asi et al., 2021; Hu et al., 2021). In Section 4.2, we only sketched the main changes in Algorithm 3 because of space limitations. We start by explicitly writing down the algorithm with MLMC estimator.

### Algorithm 6 Inexact KM iteration for problems with weak MVI

**Input:** Parameters  $\eta, L, \rho, \alpha = 1 - \frac{\rho}{\eta}$ ,  $\alpha_k = \frac{\alpha}{\sqrt{k+2\log(k+3)}} K > 0$ , initial iterate  $x_0 \in \mathbb{R}^d$ , subroutine MLMC-FBF given in Algorithm 7

$$\begin{aligned} &\text{for } k=0,1,2,\ldots,K-1 \text{ do} \\ &N_k = \lceil \frac{96(1-\eta L)^{-2}}{\min\{\frac{2}{120\alpha(k+1)},\frac{1}{120}\}} \rceil \text{ and } M_k = \lceil \frac{672\times120(\log_2N_k)}{(1-\eta L)^2} \rceil \\ &\widetilde{J}_{\eta(F+G)}^{(m)}(x_k) = \text{MLMC-FBF}\left(x_k,N_k,\eta G,\operatorname{Id} + \eta \widetilde{F},1+\eta L\right) \text{ independently for each } m=1,\ldots,M_k \\ &\widetilde{J}_{\eta(F+G)}(x_k) = \frac{1}{M_k} \sum_{i=1}^{M_k} \widetilde{J}_{\eta(F+G)}^{(i)}(x_k) \\ &x_{k+1} = (1-\alpha_k)x_k + \alpha_k \widetilde{J}_{\eta(F+G)}(x_k) \end{aligned}$$
 end for

# **Algorithm 7** MLMC-FBF $(z_0, N, A, B, L_B)$

**Input:** Initial iterate  $z_0 \in \mathbb{R}^d$ , subsolver FBF from Algorithm 5

Define  $y^i = \text{FBF}(z_0, 2^i, \widetilde{B}, A, L_B)$  for any  $i \geq 0$ . Draw  $I \sim \text{Geom}(1/2)$ 

**Output:**  $y^{\text{out}} = y^0 + 2^I(y^I - y^{I-1})$  if  $2^I \le N$ , otherwise  $y^{\text{out}} = y^0$ .

We start by modifying the proof of Lemma 3.5 for the stochastic problem, which is the most important for getting the final complexity.

**Lemma C.8.** Let Assumptions 1 and 3 hold. Suppose that the iterates generated by Algorithm 6 satisfy  $\mathbb{E}_k \|\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)\|^2 \le \varepsilon_{k,v}^2$  and  $\|\mathbb{E}_k[\widetilde{J}_{\eta(F+G)}(x_k)] - J_{\eta(F+G)}(x_k)\| \le \varepsilon_{k,b}$ . Then, we have that

$$\frac{\alpha}{4} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \| (\text{Id} - J_{\eta(F+G)})(x_k) \|^2 \le \frac{1}{2} \|x_0 - x^*\|^2 + \frac{3}{2} \sum_{k=0}^{K-1} \alpha_k^2 \mathbb{E} [\varepsilon_{k,v}^2] + \sum_{k=0}^{K-1} \alpha_k \mathbb{E} [\|x_k - x^*\| \varepsilon_{k,b}].$$

*Proof.* We proceed mostly as the proof of Lemma 3.5 apart from minor changes due to the stochastic setting such as iteration-dependent step sizes.

From Fact B.2(ii), we have that  $\mathrm{Id} - J_{\eta(F+G)}$  is  $\left(1 - \frac{\rho}{\eta}\right)$ -star cocoercive. Recall our running notations:

$$\alpha = 1 - \frac{\rho}{\eta}, \quad R = \operatorname{Id} - J_{\eta(F+G)}, \quad \widetilde{R} = \operatorname{Id} - \widetilde{J}_{\eta(F+G)}.$$

As a result, we have the following equivalent representation of  $x_{k+1}$  (see the definition in Algorithm 6):

$$x_{k+1} = x_k - \alpha_k \widetilde{R}(x_k). \tag{58}$$

By  $\alpha$ -star-cocoercivity of  $R = \mathrm{Id} - J_{\eta(F+G)}$ , we have

$$\langle R(x_k), x_k - x^* \rangle \ge \alpha \|R(x_k)\|^2. \tag{59}$$

By a simple decomposition, we write

$$\langle R(x_k), x_k - x^* \rangle = \langle \widetilde{R}(x_k), x_k - x^* \rangle + \langle R(x_k) - \widetilde{R}(x_k), x_k - x^* \rangle. \tag{60}$$

For the expectation of the first term on the right-hand side of (60), we derive that (cf. (33))

$$\mathbb{E}\langle \widetilde{R}(x_{k}), x_{k} - x^{*} \rangle = \frac{1}{\alpha_{k}} \mathbb{E}\langle x_{k} - x_{k+1}, x_{k} - x^{*} \rangle 
= \frac{1}{2\alpha_{k}} \mathbb{E} \left( \|x_{k} - x_{k+1}\|^{2} + \|x_{k} - x^{*}\|^{2} - \|x_{k+1} - x^{*}\|^{2} \right) 
\leq \frac{1}{2\alpha_{k}} \mathbb{E} \left( \|x_{k} - x^{*}\|^{2} - \|x_{k+1} - x^{*}\|^{2} \right) + \frac{3\alpha_{k}}{4} \mathbb{E} \|R(x_{k})\|^{2} + \frac{3\alpha_{k}}{2} \mathbb{E} \|\widetilde{R}(x_{k}) - R(x_{k})\|^{2} 
\leq \frac{1}{2\alpha_{k}} \mathbb{E} \left( \|x_{k} - x^{*}\|^{2} - \|x_{k+1} - x^{*}\|^{2} \right) + \frac{3\alpha}{4} \mathbb{E} \|R(x_{k})\|^{2} + \frac{3\alpha_{k} \mathbb{E}[\varepsilon_{k,v}^{2}]}{2}, \tag{61}$$

where we used the definition of  $x_{k+1}$  from (58) in the first step, standard expansion  $||a-b||^2 = ||a||^2 - 2\langle a,b\rangle + ||b||^2$  for the second step, the definition of  $x_{k+1}$  from (58) and Young's inequality in the third step, the definition of  $\varepsilon_{k,v}$  with tower rule and  $\alpha_k \leq \alpha$  in the last step.

For the second term on the right-hand side of (60), we have, by Cauchy-Schwarz inequality and the definition of  $\tilde{R}$  and  $\varepsilon_{k,b}$ , that

$$\mathbb{E}\langle R(x_k) - \widetilde{R}(x_k), x_k - x^* \rangle = \mathbb{E}[\mathbb{E}_k \langle R(x_k) - \widetilde{R}(x_k), x_k - x^* \rangle]$$

$$= \mathbb{E}\langle \mathbb{E}_k [R(x_k) - \widetilde{R}(x_k)], x_k - x^* \rangle$$

$$\leq \mathbb{E}\left[ \|R(x_k) - \mathbb{E}_k [\widetilde{R}(x_k)] \| \|x_k - x^* \|\right]$$

$$\leq \mathbb{E}\left[ \|x_k - x^* \| \varepsilon_{k,b} \right], \tag{62}$$

where the first step is by tower rule and the second step is by  $x_k - x^*$  being measurable under the conditioning of  $\mathbb{E}_k$ .

We combine (61) and (62) in (60), plug in the result to (59) and rearrange to obtain

$$\frac{\alpha}{4} \mathbb{E} \|R(x_k)\|^2 \le \frac{1}{2\alpha_k} \mathbb{E} \left( \|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2 \right) + \frac{3\alpha_k \mathbb{E}[\varepsilon_{k,v}^2]}{2} + \mathbb{E}[\|x_k - x^*\|\varepsilon_{k,b}].$$

We conclude after multiplying both sides by  $\alpha_k$  and summing for  $k = 0, 1, \dots, K - 1$ .

The next lemma considers the bias and variance of the MLMC estimator and follows the same arguments as (Asi et al., 2021, Proposition 1). The only change is that we use the algorithm FBF (see Algorithm 5 for a stochastic version) as the subsolver and we consider a strongly monotone inclusion problem rather than minimization. These do not alter the estimations significantly as can be seen in the proof.

**Lemma C.9.** Under the same setting as Theorem C.1 and  $N \ge 2$ , for the output of Algorithm 7, we have that

$$\|\mathbb{E}[y^{\text{out}}] - z^*\|^2 \le \frac{12L_B/\mu \|z_0 - z^*\|^2 + 96\sigma^2/\mu^2}{N}$$
$$\mathbb{E}\|y^{\text{out}} - z^*\|^2 \le 14(6L_B/\mu \|z_0 - z^*\|^2 + 48\sigma^2/\mu^2)\log_2 N$$

where the expected number of calls to  $\widetilde{F}$  is  $O(\log_2 N)$ .

*Proof.* We argue as (Asi et al., 2021, Property 1). The only difference is that we call Theorem C.1 which is our main solver for the strongly monotone problem.

Let us denote  $i_N = \max\{i \ge 0 \colon 2^i \le N\}$ . For a given event E, consider also the following notation for the characteristic function:  $\mathbf{1}_E = 1$  if E is true and  $\mathbf{1}_E = 0$  if E is false.

Then, we have by the definition of  $y^{\text{out}}$  in Algorithm 7 that

$$\mathbb{E}[y^{\text{out}}] = \mathbb{E}[y^{0}] + \mathbb{E}[\mathbf{1}_{\{2^{I} \leq N\}} \cdot 2^{I}(y^{I} - y^{I-1})]$$

$$= \mathbb{E}[y^{0}] + \sum_{i=1}^{i_{N}} \Pr(I = i) 2^{i} \mathbb{E}[y^{i} - y^{i-1}]$$

$$= \mathbb{E}[y^{0}] + \mathbb{E}[y^{i_{N}} - y^{0}]$$

$$= \mathbb{E}[y^{i_{N}}]. \tag{63}$$

By the definition of  $i_N$ , we have that  $2^{i_N} \geq \frac{N}{2}$  and hence, by Jensen's inequality and Theorem C.1, we have

$$\|\mathbb{E}[y^{i_N}] - z^*\|^2 \le \mathbb{E}\|y^{i_N} - z^*\|^2$$

$$\le \frac{12L_B\|z_0 - z^*\|^2 + 96\sigma^2/\mu}{N\mu},$$

which is the claimed bound on the bias due to (63).

We continue with estimating the variance of  $y^{\text{out}}$ . First, Young's inequality gives that

$$\mathbb{E}\|y^{\text{out}} - z^{\star}\|^{2} \le 2\mathbb{E}\|y^{\text{out}} - y^{0}\|^{2} + 2\mathbb{E}\|y^{0} - z^{\star}\|^{2}.$$
(64)

We estimate the first term on the right-hand side:

$$\mathbb{E}\|y^{\text{out}} - y^{0}\|^{2} = \sum_{i=1}^{i_{N}} \Pr(I = i)\mathbb{E}\|2^{i}(y^{i} - y^{i-1})\|^{2}$$

$$= \sum_{i=1}^{i_{N}} 2^{i}\mathbb{E}\|y^{i} - y^{i-1}\|^{2}$$

$$\leq \sum_{i=1}^{i_{N}} 2^{i+1} \left(\mathbb{E}\|y^{i} - z^{\star}\|^{2} + \mathbb{E}\|y^{i-1} - z^{\star}\|^{2}\right), \tag{65}$$

where the last step is by Young's inequality.

By the definitions of  $y^i, y^{i-1}$  and Theorem C.1, we have that

$$\mathbb{E}||y^{i} - z^{\star}||^{2} \le \frac{6L_{B}||z_{0} - z^{\star}||^{2} + 48\sigma^{2}/\mu}{2^{i}\mu},$$

$$\mathbb{E}||y^{i-1} - z^{\star}||^{2} \le \frac{6L_{B}||z_{0} - z^{\star}||^{2} + 48\sigma^{2}/\mu}{2^{i-1}\mu}.$$

This gives, in view of (65), that

$$\mathbb{E}\|y^{\text{out}} - y^0\|^2 \le \frac{6(6L_B\|z_0 - z^*\|^2 + 48\sigma^2/\mu)}{\mu} i_N.$$

The second term on the right-hand side of (64) is estimated the same way by using Theorem C.1:

$$\mathbb{E}||y^0 - z^*||^2 \le \frac{6L_B||z_0 - z^*||^2 + 48\sigma^2/\mu}{\mu}.$$

Combining the last two estimates in (64) gives the claimed bound on the variance after using  $i_N \leq \log_2 N$ .

The expected number of calls to  $\widetilde{B}$  is calculated as

$$2 + 2\sum_{i=1}^{i_N} P(I=i)(2^i + 2^{i-1}) = O(1 + i_N) = O(1 + \log_2 N),$$

since each iteration of stochastic FBF uses 2 unbiased samples of F. This completes the proof.

In fact, Algorithm 6 computes independent draws of MLMC-FBF and averages them to get a better control on the variance as (Asi et al., 2021, Theorem 1).

**Corollary C.10.** Let  $\widetilde{J}_{\eta(F+G)}(x_k)$  be defined as in Algorithm 6 and consider the setting of Theorem C.1. Then, for any  $b_k$ , v, we have the bias and variance bounds given as

$$\|\mathbb{E}_{k}[\widetilde{J}_{\eta(F+G)}(x_{k})] - J_{\eta(F+G)}(x_{k})\|^{2} \le b_{k}^{2}(\|(\mathrm{Id} + J_{\eta(F+G)})(x_{k})\|^{2} + \sigma^{2}),$$
  
$$\mathbb{E}_{k}\|\widetilde{J}_{\eta(F+G)}(x_{k}) - J_{\eta(F+G)}(x_{k})\|^{2} \le v^{2}(\|(\mathrm{Id} + J_{\eta(F+G)})(x_{k})\|^{2} + \sigma^{2}),$$

where

$$N_k = \left\lceil \frac{\max\{12L_B/\mu, 96/\mu^2\}}{\min\{b_k^2, \frac{v^2}{2}\}} \right\rceil, \quad \text{and} \quad M_k = \left\lceil \frac{2\log_2 N_k \max\{84L_B/\mu, 672/\mu^2\}}{v^2} \right\rceil.$$

Each iteration makes in expectation  $O(\log N_k \cdot M_k)$  calls to stochastic first-order oracle.

*Proof.* This proof follows the arguments in (Asi et al., 2021, Theorem 1). The difference is that we set the values of  $N_k$ ,  $M_k$  independent of  $||R(x_k)||^2$  and  $\sigma^2$ , to make  $N_k$ ,  $M_k$  computable, which results in these terms appearing in the bias and variance upper bounds.

We first note that  $\mathbb{E}_k[\widetilde{J}_{\eta(F+G)}(x_k)] = \mathbb{E}_k[\widetilde{J}_{\eta(F+G)}^{(1)}(x_k)]$ . We next have by direct expansion that

$$\mathbb{E}_{k} \| \widetilde{J}_{\eta(F+G)}(x_{k}) - J_{\eta(F+G)}(x_{k}) \|^{2} = \frac{1}{M_{k}} \mathbb{E}_{k} \| \widetilde{J}_{\eta(F+G)}^{(1)}(x_{k}) - J_{\eta(F+G)}(x_{k}) \|^{2} \\
+ \left( 1 - \frac{1}{M_{k}} \right) \| \mathbb{E}_{k} [\widetilde{J}_{\eta(F+G)}^{(1)}(x_{k})] - J_{\eta(F+G)}(x_{k}) \|^{2},$$

since  $\widetilde{J}_{\eta(F+G)}^{(i)}$  are independent draws of the same estimator.

By applying the identity  $\mathbb{E}\|X\|^2 = \mathbb{E}\|X - \mathbb{E}X\|^2 + \|\mathbb{E}X\|^2$  with  $X = \widetilde{J}_{\eta(F+G)}^{(1)}(x_k) - J_{\eta(F+G)}(x_k)$ , we obtain

$$\mathbb{E}_{k} \| \widetilde{J}_{\eta(F+G)}(x_{k}) - J_{\eta(F+G)}(x_{k}) \|^{2} = \frac{1}{M_{k}} \mathbb{E}_{k} \| \widetilde{J}_{\eta(F+G)}^{(1)}(x_{k}) - \mathbb{E}_{k} [\widetilde{J}_{\eta(F+G)}^{(1)}(x_{k})] \|^{2} 
+ \| \mathbb{E}_{k} [\widetilde{J}_{\eta(F+G)}^{(1)}(x_{k})] - J_{\eta(F+G)}(x_{k}) \|^{2}.$$
(66)

On the one hand, the fact  $\mathbb{E}||X - \mathbb{E}X||^2 \le \mathbb{E}||X||^2$  gives

$$\mathbb{E}_{k} \| \widetilde{J}_{\eta(F+G)}^{(1)}(x_{k}) - \mathbb{E}_{k} [\widetilde{J}_{\eta(F+G)}^{(1)}(x_{k})] \|^{2} \le \mathbb{E}_{k} \| \widetilde{J}_{\eta(F+G)}^{(1)}(x_{k}) - J_{\eta(F+G)}(x_{k}) \|^{2}.$$

$$(67)$$

On the other hand, the bounds in Lemma C.9 gives, after substituting  $z_0 = x_k$  and  $z^* = J_{\eta(F+G)}(x_k)$  that

$$\|\mathbb{E}_{k}[\widetilde{J}_{\eta(F+G)}^{(1)}(x_{k})] - J_{\eta(F+G)}(x_{k})\|^{2} \le \frac{12L_{B}/\mu\|(\mathrm{Id} + J_{\eta(F+G)})(x_{k})\|^{2} + 96\sigma^{2}/\mu^{2}}{N_{k}},\tag{68a}$$

$$\mathbb{E}_{k} \| \widetilde{J}_{\eta(F+G)}^{(1)}(x_{k}) - J_{\eta(F+G)}(x_{k}) \|^{2} \le \left( 84L_{B}/\mu \| (\operatorname{Id} + J_{\eta(F+G)})(x_{k}) \|^{2} + 672\sigma^{2}/\mu^{2} \right) \log_{2} N_{k}. \tag{68b}$$

Using  $\mathbb{E}_k[\widetilde{J}_{\eta(F+G)}(x_k)] = \mathbb{E}_k[\widetilde{J}_{\eta(F+G)}^{(1)}(x_k)]$  gives the bias bound after using the definition of  $N_k$  and (68a)

Plugging in (68b) and (67) in (66) gives the variance bound after substituting the values of  $N_k$  and  $M_k$ .

#### C.3.1. Proof for Corollary 4.4

Corollary 4.4 is essentially the summary of the results proven below.

Let us remark the recent work (Bravo & Cominetti, 2024, Corollary 5.4) that studied stochastic KM iteration for nonexpansive operators on normed spaces. This work assumes access to an unbiased oracle of the nonexpansive operator at hand and get the complexity  $\widetilde{O}(\varepsilon^{-4})$ . As mentioned in Section 4.2, this corresponds to requiring unbiased samples of  $J_{\eta(F+G)}$  in our setting, which is difficult due to the definition of the resolvent. We get the same complexity up to logarithmic factors without access to unbiased samples of  $J_{\eta(F+G)}$ , which we go around by using the MLMC technique. We also do not require nonexpansiveness from  $J_{\eta(F+G)}$  and work with conic quasi-nonexpansiveness.

**Theorem C.11.** Let Assumptions 1, 3, and 4 hold. Consider Algorithm 6 with  $\eta < \frac{1}{L}$  and  $\rho < \eta$ . Then, we have for  $K \ge 1$  that

$$\mathbb{E}_{x^{\text{out}} \sim \text{Unif}\{x_0, \dots, x_{K-1}\}} [\mathbb{E} \| (\text{Id} - J_{\eta(F+G)})(x^{\text{out}}) \|^2] \le \frac{64(\|x_0 - x^*\|^2 + \alpha^2 \sigma^2) \log(K+3)}{\alpha^2 \sqrt{K}},$$

where  $\alpha = 1 - \frac{\rho}{\eta}$ . Each iteration makes, in expectation,  $O(\log^2(k+2))$  calls to stochastic oracle  $\widetilde{B}$  and resolvent of A. Hence to obtain  $\mathbb{E}\|(\mathrm{Id} - J_{\eta(F+G)})(x^{\mathrm{out}})\| \leq \varepsilon$ , we have the expected stochastic first-order complexity  $\widetilde{O}(\varepsilon^{-4})$ .

The main reason for the length of the following proof is the lack of boundedness of  $(x_k)$ . In particular, proving this theorem is rather straightforward when we assume a bounded domain. We have to handle the complications without this assumption. There are also additional difficulties that arise because we are making sure that the inputs to MLMC-FBF will not involve unknown quantities such as  $||x_0 - x^*||$  or  $\sigma$  to run the algorithm. These are, for example, used in (Chen & Luo, 2022) for setting the parameters. Because of this reason, the bounds for  $\varepsilon_{k,v}$  and  $\varepsilon_{k,b}$  involve  $||(\mathrm{Id} + J_{\eta(F+G)})(x_k)||$  and  $\sigma^2$ .

The main reason for the difficulty here is  $\|(\mathrm{Id} + J_{\eta(F+G)})(x_k)\|^2$ , since we do not have a uniform bound on this quantity, unlike  $\sigma^2$  and this term appears in many summands. We will carry these terms coming from the MLMC bounds to get a recursion involving the sum of  $\|(\mathrm{Id} + J_{\eta(F+G)})(x_k)\|^2$  for different ranges on both sides. We then go around the issue of lacking of a bound on  $(x_k)$  by using an inductive argument on  $\sum_{k=0}^K \|(\mathrm{Id} + J_{\eta(F+G)})(x_k)\|^2$ .

*Proof of Theorem C.11.* Recall our running notations:

$$\alpha = 1 - \frac{\rho}{\eta}, \quad R = \operatorname{Id} - J_{\eta(F+G)}, \quad \widetilde{R} = \operatorname{Id} - \widetilde{J}_{\eta(F+G)}.$$

We start by following the proof of Lemma B.6. By  $\alpha$ -star-cocoercivity of  $\mathrm{Id}-J_{\eta(F+G)}$  and  $\alpha \geq \alpha_k$  (which gives that  $J_{\eta(F+G)}$  is  $\frac{1}{\alpha_k}$ -star-conic nonexpansive), we can use property (27) derived in Lemma B.3 to obtain

$$\|(1 - \alpha_k)x_k + \alpha_k J_{\eta(F+G)}(x_k) - x^*\| \le \|x_k - x^*\|$$

and by the definition of  $x_{k+1}$ , we get

$$||x_{k+1} - x^*|| \le ||(1 - \alpha_k)x_k + \alpha_k J_{\eta(F+G)}(x_k) - x^*|| + \alpha_k ||\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)||$$

$$\le ||x_k - x^*|| + \alpha_k ||\widetilde{J}_{\eta(F+G)}(x_k) - J_{\eta(F+G)}(x_k)||.$$
(69)

Summing the inequality for  $0, \ldots, k-1$  gives

$$||x_{k} - x^{*}|| \leq ||x_{0} - x^{*}|| + \sum_{i=0}^{k-1} \alpha_{i}||\widetilde{J}_{\eta(F+G)}(x_{i}) - J_{\eta(F+G)}(x_{i})||$$

$$\implies \mathbb{E}||x_{k} - x^{*}||^{2} \leq 2\mathbb{E}||x_{0} - x^{*}||^{2} + 2k \sum_{i=0}^{k-1} \alpha_{i}^{2} \mathbb{E}||\widetilde{J}_{\eta(F+G)}(x_{i}) - J_{\eta(F+G)}(x_{i})||^{2},$$
(70)

where we first squared both sides, used Young's inequality and then took expectation.

We continue by restating the result of Lemma C.8 after applying Young's inequality on the last term to obtain

$$\frac{\alpha}{4} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \| (\text{Id} - J_{\eta(F+G)})(x_k) \|^2 \le \frac{1}{2} \|x_0 - x^*\|^2 + \frac{3}{2} \sum_{k=0}^{K-1} \alpha_k^2 \mathbb{E} [\varepsilon_{k,v}^2] 
+ \sum_{k=0}^{K-1} \left( \frac{\alpha_k^2}{2\alpha^2 (k+1)} \mathbb{E} \|x_k - x^*\|^2 + \frac{(k+1)\alpha^2}{2} \mathbb{E} [\varepsilon_{k,b}^2] \right).$$
(71)

We now estimate the second and third terms on the right-hand side. By using Corollary C.10 and the definition of  $R(x_k)$ ,  $\alpha_k = \frac{\alpha}{\sqrt{k+2}\log(k+3)} \le \frac{\alpha}{\sqrt{2}\log 3}$  and using  $v^2 = \frac{1}{60} \le \frac{\sqrt{2}\log 3}{24}$ , we obtain

$$\frac{3}{2} \sum_{k=0}^{K-1} \alpha_k^2 \mathbb{E}[\varepsilon_{k,v}^2] \le \frac{3}{2} \sum_{k=0}^{K-1} \alpha_k^2 v^2 \left( \mathbb{E} \| R(x_k) \|^2 + \sigma^2 \right) \\
\le \frac{\alpha \alpha_{K-1}}{16} \left( \mathbb{E} \| R(x_{K-1}) \|^2 + \sigma^2 \right) + \frac{3}{2} \sum_{k=0}^{K-2} \alpha_k^2 v^2 \left( \mathbb{E} \| R(x_k) \|^2 + \sigma^2 \right).$$
(72)

We continue with the first part of the third term on the right-hand side of (71) and bound it using (70):

$$\sum_{k=0}^{K-1} \frac{\alpha_k^2}{2\alpha^2(k+1)} \mathbb{E} \|x_k - x^*\|^2 \le \frac{1}{2} \|x_0 - x^*\|^2 + \sum_{k=1}^{K-1} \frac{\alpha_k^2}{2\alpha^2(k+1)} \mathbb{E} \|x_k - x^*\|^2 
\le \frac{1}{2} \|x_0 - x^*\|^2 + \sum_{k=1}^{K-1} \frac{\alpha_k^2}{2\alpha^2(k+1)} \left( 2\|x_0 - x^*\|^2 + 2k \sum_{i=0}^{k-1} \alpha_i^2 \mathbb{E}[\varepsilon_{i,v}^2] \right), \quad (73)$$

where the last line identified  $\varepsilon_{i,v}^2$  in view of Lemma C.8.

We focus on the last term here to get

$$\begin{split} \sum_{k=1}^{K-1} \frac{\alpha_k^2}{2\alpha^2(k+1)} \cdot 2k \sum_{i=0}^{k-1} \alpha_i^2 \mathbb{E}[\varepsilon_{i,v}^2] &= \frac{1}{\alpha^2} \sum_{i=0}^{K-2} \sum_{k=i+1}^{K-1} \frac{k}{k+1} \alpha_k^2 \alpha_i^2 \mathbb{E}[\varepsilon_{i,v}^2] \\ &\leq \frac{1}{\alpha^2} \left( \sum_{k=0}^{K-1} \alpha_k^2 \right) \sum_{i=0}^{K-2} \alpha_i^2 \mathbb{E}[\varepsilon_{i,v}^2] \\ &\leq \frac{1}{\alpha^2} \left( \sum_{k=0}^{K-1} \alpha_k^2 \right) \sum_{i=0}^{K-2} \alpha_i^2 v^2 \left( \mathbb{E} \|R(x_i)\|^2 + \sigma^2 \right), \end{split}$$

where the last step used Corollary C.10.

Plugging in back to (73) gives

$$\sum_{k=0}^{K-1} \frac{\alpha_k^2}{2\alpha^2(k+1)} \mathbb{E} \|x_k - x^*\|^2 \le \left(\frac{1}{2} + \sum_{k=1}^{K-1} \frac{\alpha_k^2}{\alpha^2(k+1)}\right) \|x_0 - x^*\|^2 + \left(\sum_{k=0}^{K-1} \frac{\alpha_k^2}{\alpha^2}\right) \sum_{i=0}^{K-2} \alpha_i^2 v^2 \left(\mathbb{E} \|R(x_i)\|^2 + \sigma^2\right).$$
(74)

By using  $\alpha_k = \frac{\alpha}{\sqrt{k+2}\log(k+3)}$ , we have

$$\sum_{k=0}^{K-1} \frac{\alpha_k^2}{\alpha^2} < 3, \quad \sum_{k=0}^{K-1} \frac{\alpha_k^2}{\alpha^2(k+1)} < 0.25, \quad \alpha_k \le \alpha \, \forall k \ge 0,$$

which helps estimate (74) as

$$\sum_{k=0}^{K-1} \frac{\alpha_k^2}{2\alpha^2(k+1)} \mathbb{E} \|x_k - x^\star\|^2 \le \frac{3}{4} \|x_0 - x^\star\|^2 + 3 \sum_{i=0}^{K-2} \alpha_i^2 v^2 \left( \mathbb{E} \|R(x_i)\|^2 + \sigma^2 \right). \tag{75}$$

We finally estimate the second part of the third term on the right-hand side of (71) by using Corollary C.10:

$$\alpha^2 \sum_{k=0}^{K-1} \frac{k+1}{2} \mathbb{E}[\varepsilon_{k,b}^2] \le \alpha^2 \sum_{k=0}^{K-1} (k+1) b_k^2 \mathbb{E}[\|R(x_k)\|^2 + \sigma^2].$$

We use the setting  $b_k^2=\frac{\alpha_k}{120\alpha(k+1)}$  and  $b_{K-1}^2<\frac{\alpha_{K-1}}{16\alpha K}$  to obtain

$$\alpha^{2} \sum_{k=0}^{K-1} \frac{k+1}{2} \mathbb{E}[\varepsilon_{k,b}^{2}] \leq \frac{\alpha \alpha_{K-1}}{16} (\mathbb{E} \|R(x_{K-1})\|^{2} + \sigma^{2}) + \alpha^{2} \sum_{k=0}^{K-2} (k+1) b_{k}^{2} \mathbb{E}[\|R(x_{k})\|^{2} + \sigma^{2}]. \tag{76}$$

We collect (72), (75), and (76) in (71) to get

$$\frac{\alpha}{8} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|R(x_k)\|^2 \le \frac{5}{4} \|x_0 - x^*\|^2 + \frac{\alpha \alpha_{K-1}}{8} \sigma^2 + \frac{9}{2} \sum_{k=0}^{K-2} \alpha_k^2 v^2 \left( \mathbb{E} \|R(x_k)\|^2 + \sigma^2 \right) + \alpha^2 \sum_{k=0}^{K-2} (k+1) b_k^2 \mathbb{E} [\|R(x_k)\|^2 + \sigma^2].$$
(77)

We now show by induction that

$$\alpha \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|R(x_k)\|^2 \le C \left( \|x_0 - x^*\|^2 + \alpha^2 \sigma^2 \right) \quad \forall K \ge 1, \tag{78}$$

for some C to be determined. With  $\alpha < 1$ , (77) becomes

$$\frac{\alpha}{8} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|R(x_k)\|^2 \le 1.25 \|x_0 - x^*\|^2 + \alpha^2 \sigma^2 + 4.5 \sum_{k=0}^{K-2} v^2 (\alpha \alpha_k \mathbb{E} \|R(x_k)\|^2 + \alpha^2 \sigma^2) 
+ \alpha \sum_{k=0}^{K-2} \frac{(k+1)b_k^2}{\alpha_k} \mathbb{E} [\alpha \alpha_k \|R(x_k)\|^2 + \alpha^2 \sigma^2].$$
(79)

Let us set

$$C = 32$$
,  $b_k^2 = \frac{\alpha_k}{120\alpha(k+1)}$ ,  $v^2 = \frac{1}{60}$ 

and use the inductive assumption  $\alpha \sum_{k=0}^{K-2} \alpha_k \mathbb{E} \|R(x_k)\|^2 \leq 32(\|x_0 - x^\star\|^2 + \alpha^2 \sigma^2)$  in (79) to obtain

$$\frac{\alpha}{8} \sum_{k=0}^{K-1} \alpha_k \mathbb{E} ||R(x_k)||^2 \le 4(||x_0 - x^*||^2 + \alpha^2 \sigma^2),$$

which verifies  $\alpha \sum_{k=0}^{K-1} \alpha_k \mathbb{E} \|R(x_k)\|^2 \le 32(\|x_0 - x^\star\|^2 + \alpha^2 \sigma^2)$ .

For the base case, we use  $\alpha_0 = \frac{\alpha}{\sqrt{2}\log 3} < 1$  and  $\alpha^{-1}$ -star-Lipschitzness of  $R = \mathrm{Id} - J_{\eta(F+G)}$  to get  $\alpha\alpha_0\|R(x_0)\|^2 \le \|x_0 - x^\star\|^2$ . This establishes the base case and completes the induction.

By using  $\alpha_k \geq \alpha_K = \frac{\alpha}{\sqrt{K+2\log(K+3)}}$  in (78) with C=32 and multiplying both sides by  $\frac{1}{K\alpha_K}$  and using  $\frac{\sqrt{K+2}}{K} \leq \frac{2}{\sqrt{K}}$  which is true for  $K \geq 1$ , we get the claimed rate result. Finally, in view of Corollary C.10, and definitions of  $b_k, v$ , each iteration makes expected number of calls  $O(\log^2(k+1))$ . By using this expected cost of each iteration, we also get the final expected stochastic first-order complexity result.

#### D. Additional Remarks on Related Work

There exist a line of works that attempted to construct local estimation of Lipschitz constants to offer an improved range for  $\rho$  depending on the curvature (Pethick et al., 2022; Alacaoglu et al., 2023). However, these results cannot bring global improvements in the worst-case range of  $\rho$  where the limit for  $\rho$  is still  $\frac{1}{2L}$ . This is because it is easy to construct examples where the local Lipschitz constants are the same as the global Lipschitz constant.

The work (Hajizadeh et al., 2023) gets linear rate of convergence for interaction dominant problems which is shown to be closely related to cohypomonotonity, see Example 1. One important difference is that cohypomonotonicity is equivalent to  $\alpha$  interaction dominance with  $\alpha \geq 0$  whereas (Hajizadeh et al., 2023) requires  $\alpha > 0$  for linear convergence. This is an important difference because (i) we know that cohypomonotonicity relaxes monotonicity and (ii) we know that even monotonicity is not sufficient for linear convergence. For monotone problems  $O(\varepsilon^{-1})$  is the optimal first-order oracle complexity (see, e.g., (Yoon & Ryu, 2021)) and hence it is also optimal with cohypomonotonicity.

In the literature for fixed point iterations, several works considered inexact Halpern or KM iterations without characterizing explicit first-order complexity results, see for example (Leuştean & Pinto, 2021; Bartz et al., 2022; Kohlenbach, 2022; Combettes & Pennanen, 2004). In particular, Bartz et al. (2022) used conic nonexpansiveness to analyze KM iteration. The dependence of the range of  $\rho$  on L arises when we start characterizing the first-order complexity. This is the reason these works have not been included in comparisons in Table 1.

For the stochastic cohypomonotone problems, the best complexity result to our knowledge is due to (Chen & Luo, 2022). This paper can obtain the optimal complexity  $\widetilde{O}(\varepsilon^{-2})$  with cohypomonotone stochastic problems with a 6-loop algorithm using many carefully designed regularization techniques, extending the work of Allen-Zhu (2018) that focused on minimization. Some disadvantages of this approach compared to ours: (i) the bound for cohypomonotonicity is  $\rho \leq \frac{1}{2L}$ ; (ii) the algorithm needs estimates of variance upper bound  $\sigma^2$  and, more importantly,  $\|x^0 - x^\star\|^2$ ; (iii) the result is only given for unconstrained problems, which also makes it difficult to assume a bounded domain since there is no guarantee a priori for the iterates to stay bounded for an unconstrained problem. Given that the 6-loop algorithm and analysis of (Chen & Luo, 2022) is rather complicated, it is not clear to us if their arguments generalize to constraints or if the other drawbacks can be alleviated.

The work (Tran-Dinh & Luo, 2023) focused on problems with  $\rho$ -weakly MVI solutions for  $\rho < \frac{1}{8L}$  and derived  $O(\varepsilon^{-2})$  for a randomized coordinate algorithm. Due to randomization, the complexity result in this work holds for the expectation of the optimality measure. Because of the coordinatewise updates, the problem focused in this work is deterministic, similar to the setup in Section 3. Bravo & Contreras (2024) studied stochastic inexact Halpern iteration in normed spaces and obtained complexity  $O(\varepsilon^{-5})$  for finding fixed-points, by using an oracle providing unbiased samples of a nonexpansive map.

#### D.1. Clarifications about Table 2

Since the complexity results have not been written explicitly in some of the references, we provide details on how we computed the complexities that we report for the existing works.

(Choudhury et al., 2023): We use Theorem 4.5 in this corresponding paper to see that squared operator norm is upper bounded by  $O(K^{-1})$ . To make the operator norm smaller than  $\varepsilon$ , the order of K is  $\varepsilon^{-2}$ . The batch-size has order K and hence the total number of oracle calls is  $O(K^2) = O(\varepsilon^{-4})$ .

(Böhm et al., 2022): We use Theorem 3.3 in this corresponding paper. The paper stated that to make the squared operator norm smaller than  $\varepsilon$ , number of iterations is  $O(\varepsilon^{-2})$  and the batch size is  $O(\varepsilon^{-3})$ . This gives complexity  $O(\varepsilon^{-3})$  for making the *squared* operator norm smaller than  $\varepsilon$ . Hence, to make the operator norm smaller than  $\varepsilon$ , the complexity is  $O(\varepsilon^{-6})$ .

(Pethick et al., 2023b): (i) For "best rate" result, we use Corollary E.3(i) in this corresponding paper. The dominant term in the bound for the squared residual is  $O(K^{-1})$ . Hence to make the norm of the residual smaller than  $\varepsilon$  (equivalently, the squared norm smaller than  $\varepsilon^{-2}$ ), one needs K to be of the order  $\varepsilon^{-2}$ . Then, the squared variance is assumed to decrease at the order of  $k^2$  which requires the batch size at iteration k to be  $k^2$ . Then the complexity is upper bounded by  $\sum_{k=1}^K \tau k^2 = \widetilde{O}(K^3) = \widetilde{O}(\varepsilon^{-6})$ , (ii) for the "last iterate", we use the Corollary E.3(ii) given in the paper to see that the dominant term in the bound of the squared residual is  $O\left(\frac{1}{\sqrt{K}}\right)$ . To make the squared residual smaller than  $\varepsilon^2$ , this means K is of the order  $\varepsilon^{-4}$ . The squared variance is assumed to decrease at the rate  $k^3$  which requires a batch size of  $k^3$  at iteration k. Then, with the same calculation as before, the complexity of stochastic first-order oracles to make the residual less than  $\varepsilon$  is  $\sum_{k=1}^K \tau k^3 = \widetilde{O}(K^4) = \widetilde{O}(\varepsilon^{-16})$ .