

# Optimal Extragradient-Based Algorithms for Stochastic Variational Inequalities with Separable Structure

Huizhuo Yuan<sup>◊</sup>   Chris Junchi Li<sup>†</sup>   Gauthier Gidel<sup>‡</sup>  
 Michael I. Jordan<sup>†,□</sup>   Quanquan Gu<sup>◊</sup>   Simon S. Du<sup>★</sup>

Department of Computer Sciences, University of California, Los Angeles<sup>◊</sup>

Department of Electrical Engineering and Computer Sciences, University of California, Berkeley<sup>†</sup>  
 DIRO, Université de Montréal and Mila<sup>‡</sup>

Department of Statistics, University of California, Berkeley<sup>□</sup>

Paul G. Allen School of Computer Science and Engineering, University of Washington<sup>★</sup>

September 19, 2024

## Abstract

We consider the problem of solving stochastic monotone variational inequalities (VIs) with a separable structure when only stochastic oracle queries are available. Built upon standard extragradient for variational inequalities, we propose a novel algorithm called stochastic *accelerated gradient-extragradient* (AG-EG) for strongly monotone VIs with a separable structure. Our approach combines extragradient and Nesterov’s acceleration strengths, enabling faster convergence. By showing that its iterates remain in a bounded domain and applying scheduled restarting, we prove that AG-EG has an optimal convergence rate for strongly monotone VIs. Furthermore, when specializing to the particular case of bilinearly coupled strongly-convex-strongly-concave saddle-point problems as well as bilinear games, our algorithm achieves fine-grained convergence rates that match the respective lower bounds, with the effect of stochasticity being characterized by additive statistical error term that is optimal up to a constant prefactor. Our results shed light on designing accelerated extragradient-based algorithms, which can be applied to a wide range of VI problems.

## 1 Introduction

We consider the variational inequality (VI) problem which plays a central role in a wide range of optimization and control applications, including but not limited to convex minimization, saddle point problems, and games [FP03, Nem04, NJLS09, JNT11, JLZ23]. Extensive treatments on both theory and applications of variational inequality trace back to as early as [KS80]. In particular, our formulation of general VI problem aims to find a solution  $\mathbf{z}^* \in \mathcal{Z}$  that satisfies:

$$\langle \mathcal{W}(\mathbf{z}^*), \mathbf{z}^* - \mathbf{z} \rangle \leq 0, \quad \forall \mathbf{z} \in \mathcal{Z} \quad (1)$$

where  $\mathcal{Z}$  is a finite-dimensional closed and convex feasible set and  $\mathcal{W}(\cdot)$  is a monotone operator of the following form:

$$\mathcal{W}(\mathbf{z}) = \nabla \mathcal{F}(\mathbf{z}) + \mathcal{H}(\mathbf{z}) + J'(\mathbf{z}) \equiv \mathbb{E}_{\xi}[\nabla \tilde{\mathcal{F}}(\mathbf{z}; \xi)] + \mathbb{E}_{\zeta}[\tilde{\mathcal{H}}(\mathbf{z}; \zeta)] + J'(\mathbf{z}) \quad (2)$$

where  $\mathcal{F}$  is a continuously differentiable function which is  $L$ -smooth and  $\mu$ -strongly convex,  $\mathcal{H}$  is an  $M$ -Lipschitz monotone operator,  $J' \in \partial J$  is the subgradient of a simple and convex function, and  $\xi$  and  $\zeta$  are drawn from distributions  $\mathcal{D}_{\xi}$  and  $\mathcal{D}_{\zeta}$ , respectively. This formulation captures

the separable structure of the problem where  $\nabla\mathcal{F}$  usually models the cooperations,  $\mathcal{H}$  models the competing forces in a system and  $J$  models the nonsmooth factor, respectively. In addition, we consider the stochastic setting where we can only access  $\nabla\mathcal{F}$  and  $\mathcal{H}$  through their unbiased estimators  $\nabla\tilde{\mathcal{F}}(\mathbf{z};\xi)$  and  $\tilde{\mathcal{H}}(\mathbf{z};\zeta)$  respectively.

A notable instance of the VI problem (1) with separable structure (2) is the widely studied *bilinearly coupled strongly-convex-strongly-concave saddle-point problem*:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \max_{\mathbf{y} \in \mathbb{R}^m} \mathcal{F}(\mathbf{x}, \mathbf{y}) = F(\mathbf{x}) + H(\mathbf{x}, \mathbf{y}) - G(\mathbf{y}) \equiv \mathbb{E}_{\xi} [f(\mathbf{x}; \xi)] + \mathbb{E}_{\zeta} [h(\mathbf{x}, \mathbf{y}; \zeta)] - \mathbb{E}_{\xi} [g(\mathbf{y}; \xi)], \quad (3)$$

where  $H(\mathbf{x}, \mathbf{y}) \equiv \mathbf{x}^{\top} \mathbf{B} \mathbf{y} - \mathbf{x}^{\top} \mathbf{u}_x + \mathbf{u}_y^{\top} \mathbf{y}$  is the bilinear coupling function with the coupling matrix  $\mathbf{B} \in \mathbb{R}^{n \times m}$ . Note that (3) is a special instance of (1) when taking  $\mathcal{F}(\mathbf{z}) = F(\mathbf{x}) + G(\mathbf{y})$  with  $\mathbf{z} = [\mathbf{x}; \mathbf{y}]$ ,  $\mathcal{H}(\mathbf{z}) = [\nabla_{\mathbf{x}} H(\mathbf{x}, \mathbf{y}); -\nabla_{\mathbf{y}} H(\mathbf{x}, \mathbf{y})]$  and  $J = 0$ . In addition to a wide range of applications in economics, problems of form (3) are becoming increasingly important in machine learning. For instance, (3) appears in reinforcement learning, differentiable games, regularized empirical risk minimization, and robust optimization formulations. It can also be seen as a local approximation of the nonconvex-nonconcave minimax games—e.g., the generative adversarial network (GAN) [GPAM<sup>+</sup>20]—around a local Nash equilibrium [MNG17, NK17].

In this paper, we aim to improve the efficiency of solving (1) by utilizing the structural information of the monotone operator in (2). More specifically, we consider the case when  $\mathcal{F}$  is strongly monotone, or zero. Although optimal convergence results have been achieved for the monotone VI problem (1) [CLO17] as well as the special case of convex-concave saddle point problem with bilinear coupling (3) [CLO14], it remains open how to design an optimal algorithm for the strongly monotone VI problem. Notably, for the special case (3) when  $F$  and/or  $G$  are strongly convex, several concurrent works have independently obtained the optimal convergence rates [KGR22, THO22, JST22, MRGK22, LYGJ22]. On the other hand, when both  $F$  and  $G$  are zero, optimal convergence results have been obtained by [LYL<sup>+</sup>22] and the accelerated-gradient optimistic gradient approach [LYGJ22]. We defer a more complete list of related work to the appendix.

## 1.1 Main Contributions

We start with the strongly monotone VI problem in an unbounded feasible set, extending the scope of recent work such as [JST22] and going beyond earlier studies that focus on non-strongly monotone VIs in a bounded feasible set [JNT11, CLO17].<sup>1</sup> We propose a class of algorithms named stochastic *accelerated gradient-extragradient* (AG-EG), which combines Nesterov’s acceleration with the extragradient method. By employing either a strongly-convexity shifting technique or a scheduled restarting scheme, our algorithm achieves the optimal convergence rates that match the lower bounds for the general strongly monotone VI problem (1), the special SC-SC saddle point problem (3), and the bilinear games, in both deterministic and stochastic settings, thus providing a unified optimal solution. In sharp contrast to the accelerated mirror-prox (AMP) algorithm proposed by [CLO17, JLZ23], our analysis does not rely on the boundedness of the feasible set  $\mathcal{Z}$ , which makes our algorithm projection-free. We also extend our algorithm to VIs with bounded feasible set and/or non-differential convex regularization through proximal mapping. In particular, we summarize the contributions as follows:

---

<sup>1</sup>VIs in an unbounded feasible set is more difficult to solve because existing algorithms and analyses highly rely on the boundedness of the feasible set.

- (1) We present a direct approach of stochastic AG-EG, dubbed as *Stochastic AG-EG-Direct*, for separable strongly monotone VIs, which achieves an iteration complexity of  $\mathcal{O}\left(\sqrt{\frac{L}{\mu}} + \frac{M}{\mu} + \frac{\sigma^2}{\mu^2 \varepsilon^2}\right) \log\left(\frac{L}{\mu} \frac{1}{\varepsilon}\right)$  for finding an  $\varepsilon$ -optimal point, where  $\sigma^2$  is the uniform variance bound on the stochastic operator. When  $\sigma = 0$  such a query complexity admits the sharpest possible near-unity coefficient [§2.3, Theorem 1]. The deterministic part matches the complexity lower bound in [ZHZ22], while the stochastic part matches the optimal statistical error.
- (2) We also present a stochastic AG-EG algorithm equipped with scheduled restarting, where the iteration complexity lower bound due to [ZHZ22] is matched as  $\mathcal{O}\left(\left(\sqrt{\frac{L}{\mu}} + \frac{M}{\mu}\right) \log\left(\frac{1}{\varepsilon}\right) + \frac{\sigma^2}{\mu^2 \varepsilon^2}\right)$  for finding an  $\varepsilon$ -optimal point. Such a query complexity improves by a logarithmic factor upon the previous one attained by Stochastic AG-EG-Direct.

When specializing the VI problem to bilinearly coupled SC-SC and bilinear game saddle-point problem, our results have the following implications:

**Strongly-convex-strongly-concave (SC-SC) Saddle Point Problem.** For the class of SC-SC saddle point problem, the stochastic AG-EG descent-ascent Algorithm 1, equipped with scaling reduction, achieves an iteration complexity of

$$\mathcal{O}\left(\left(\sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}\right) \log\left(\frac{1}{\varepsilon}\right) + \frac{\sigma^2}{\mu_F^2 \varepsilon^2}\right) \quad (4)$$

where  $F : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $L_F$ -smooth and  $\mu_F$ -strongly convex,  $G : \mathbb{R}^m \rightarrow \mathbb{R}$  is  $L_G$ -smooth and  $\mu_G$ -strongly convex, and  $\sigma^2$  is the weighted, uniform variance bound on the stochastic gradients. When the optimization problem is deterministic, the complexity upper bound matches the lower bound [ZHZ22][§3.1, Corollary 1].

**Bilinear Games.** For bilinear games ( $\nabla f(\mathbf{x}; \xi) = \mathbf{0}$  and  $\nabla g(\mathbf{y}; \xi) = \mathbf{0}$  almost surely), Algorithm 1, equipped with scheduled restarting achieves an

$$\mathcal{O}\left(\sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}} \log\left(\frac{\sqrt[4]{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top) \lambda_{\max}(\mathbf{B}^\top \mathbf{B})}}{\sigma_{\text{Bil}}}\right) + \frac{\sigma_{\text{Bil}}^2}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top) \varepsilon^2}\right) \quad (5)$$

iteration complexity, where  $\sigma_{\text{Bil}}^2$  is the variance of the stochastic gradient on the bilinear coupling term. When there is no randomness, this complexity result reduces to  $\mathcal{O}\left(\sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}} \log\left(\frac{1}{\varepsilon}\right)\right)$  for the bilinear games, matching the lower bound of [IAGM20] [§3.2, Corollary 4].<sup>2</sup>

**Organization.** The rest of this paper is organized as follows. Section 2 proposes for strongly monotone variational inequalities (VIs) the Accelerated Gradient-Extragradient Descent-Ascent algorithm that achieves an accelerated convergence rate, and extends to VIs with bounded domains with proximal operator. Section 3 discusses two specific instances of saddle-point problem, where our proposed AG-EG algorithm achieves upper bound results that matches the corresponding lower bounds. Finally, Section 4 summarizes the contributions of this work and suggests future directions.

<sup>2</sup>For the function class of bilinear games, we assume that  $n = m$  where  $\mathbf{B}$  is a nonsingular square matrix, so that  $\lambda_{\min}(\mathbf{B}\mathbf{B}^\top) > 0$  and the complexity makes sense. See §3.2 for more on this.

**Notation.** Let  $\lambda_{\max}(\mathbf{M})$  (resp.  $\lambda_{\min}(\mathbf{M})$ ) be the largest (resp. smallest) eigenvalue of a real symmetric matrix  $\mathbf{M}$ . Let  $a \vee b \equiv \max(a, b)$  (resp.  $a \wedge b \equiv \min(a, b)$ ) denote the maximum (resp. minimum) value of two reals  $a, b$ . For two nonnegative real sequences  $(a_n)$  and  $(b_n)$ , we write  $a_n = \mathcal{O}(b_n)$  or  $a_n \lesssim b_n$  (resp.  $a_n = \Omega(b_n)$  or  $a_n \gtrsim b_n$ ) to denote  $a_n \leq Cb_n$  (resp.  $a_n \geq Cb_n$ ) for all  $n \geq 1$  for a positive, numerical constant  $C$ , and let  $a_n \asymp b_n$  if both  $a_n \lesssim b_n$  and  $a_n \gtrsim b_n$  hold. We also let  $a_n = \tilde{\mathcal{O}}(b_n)$  denote  $a_n \leq Cb_n$  where  $C$  hides a polylogarithmic factor in problem-dependent constants. We let  $[\mathbf{x}; \mathbf{y}] \in \mathbb{R}^{n+m}$  concatenate two vectors  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^m$ . Finally for two real symmetric matrices  $\mathbf{A}$  and  $\mathbf{B}$ , we denote  $\mathbf{A} \preceq \mathbf{B}$  (resp.  $\mathbf{A} \succeq \mathbf{B}$ ) when  $\mathbf{v}^\top (\mathbf{A} - \mathbf{B}) \mathbf{v} \leq 0$  (resp.  $\mathbf{v}^\top (\mathbf{A} - \mathbf{B}) \mathbf{v} \geq 0$ ) holds for all vectors  $\mathbf{v}$ .

## 2 Accelerated Gradient-Extragradient Descent-Ascent Algorithm

In this section, we focus on accelerating the extragradient algorithm for the strongly monotone VI problem in (1) with separable structure (2). Our algorithm design draws inspiration from the work of [CLO17] on the stochastic Accelerated MirrorProx (AMP) algorithm for non-strongly monotone VIs. The AMP algorithm applies Nesterov-type acceleration on top of the mirror-prox method [Kor76, Nem04] and attains the optimal iteration complexity of  $\mathcal{O}\left(\sqrt{\frac{L}{\varepsilon}} + \frac{M}{\varepsilon}\right)$ . However, the big-O notation hides the diameter of the feasible set, thus the AMP algorithm and analysis can only deal with VIs with bounded domain. Our algorithm can not only achieve the optimal convergence rates for the strongly monotone VI problem with separable structure, but also get rid of the dependency on the diameter of the feasible set. Therefore, our algorithm can deal with VIs with unbounded domain.

Throughout §2, we maintain conceptual simplicity by presenting all our algorithm designs in the deterministic setting, while presenting the convergence results in the more general stochastic setting. These results can easily reduce to the deterministic setting when the stochastic noise vanishes.

### 2.1 Setting and Assumptions

In this section, we formally introduce our assumptions. We first state the smoothness and monotonicity assumptions that we impose on  $\mathcal{F}$  and  $\mathcal{H}$ .

**Assumption 1** (Monotonicity, strong convexity and smoothness). *We assume that function  $\mathcal{F}(\cdot)$  is continuous differentiable with  $L$ -Lipschitz continuous gradient and is  $\mu$ -strongly convex. That is, for any  $\mathbf{z}, \mathbf{z}' \in \mathcal{Z}$ ,*

$$\frac{\mu}{2} \|\mathbf{z} - \mathbf{z}'\|^2 \leq \mathcal{F}(\mathbf{z}) - \mathcal{F}(\mathbf{z}') - \nabla \mathcal{F}(\mathbf{z}')^\top (\mathbf{z} - \mathbf{z}') \leq \frac{L}{2} \|\mathbf{z} - \mathbf{z}'\|^2$$

Furthermore, operator  $\mathcal{H}(\cdot)$  is monotone and  $M$ -Lipschitz in the sense that for any  $\mathbf{z}, \mathbf{z}' \in \mathcal{Z}$ ,

$$\langle \mathcal{H}(\mathbf{z}) - \mathcal{H}(\mathbf{z}'), \mathbf{z} - \mathbf{z}' \rangle \geq 0, \quad \|\mathcal{H}(\mathbf{z}) - \mathcal{H}(\mathbf{z}')\| \leq M \|\mathbf{z} - \mathbf{z}'\|$$

Second, we impose assumptions on the noise variance.

**Assumption 2** (Unbiased gradients and variance bounds). *We assume that  $\mathbf{z} \in \mathcal{Z}$ , samples  $\xi \sim \mathcal{D}_\xi$  and  $\zeta \sim \mathcal{D}_\zeta$  are drawn from given distributions such that the following conditions hold:  $\mathbb{E}_\xi[\nabla \tilde{\mathcal{F}}(\mathbf{z}; \xi)] = \nabla \mathcal{F}(\mathbf{z})$ ,  $\mathbb{E}_\zeta[\tilde{\mathcal{H}}(\mathbf{z}; \zeta)] = \mathcal{H}(\mathbf{z})$ , and*

$$\mathbb{E}_\xi \left[ \|\nabla \tilde{\mathcal{F}}(\mathbf{z}; \xi) - \nabla \mathcal{F}(\mathbf{z})\|^2 \right] \leq \sigma_{\text{Str}}^2, \quad \mathbb{E}_\zeta \left[ \|\tilde{\mathcal{H}}(\mathbf{z}; \zeta) - \mathcal{H}(\mathbf{z})\|^2 \right] \leq \sigma_{\text{Bil}}^2 \quad (6)$$

For all results in this work, we suppose that Assumptions 1 and 2 hold with appropriate parameter settings. Given a desired accuracy  $\varepsilon > 0$ , our goal is to find an  $\varepsilon$ -optimal minimax point defined as:

**Definition 1** ( $\varepsilon$ -Optimal point). *A point  $\mathbf{z} \in \mathcal{Z}$  is called an  $\varepsilon$ -optimal point for the VI problem in (1) if  $\|\mathbf{z} - \mathbf{z}^*\| \leq \varepsilon$ .*

## 2.2 The ExtraGradient (EG) Algorithm

We first consider the cases where  $\mathcal{Z}$  is the entire space  $\mathbf{R}^n$  and the objective is smooth ( $J = 0$ ). The extragradient (EG) algorithm, introduced by [Kor76], is designed to address cyclic behavior in saddle-point problems by introducing an extrapolated point for gradient evaluation. In the context of VI problems (1), let  $\mathbf{z}_t$  represents the  $t$ -th iterate of the EG algorithm. The update rule of EG is as follows:

$$\mathbf{z}_{t+1} = \mathbf{z}_t - \eta \mathcal{W}(\mathbf{z}_t - \eta \mathcal{W}(\mathbf{z}_t)) \quad (7)$$

where  $\eta > 0$  is the step size. For  $L$ -smooth and  $\mu$ -strongly monotone operator  $\mathcal{W}$ , [Tse95, MOP20, GBV<sup>+</sup>19] have shown that the EG algorithm achieves an iteration complexity of  $\mathcal{O}(\kappa \log(1/\varepsilon))$ , where  $\kappa = L/\mu$  denotes the condition number of the problem.

## 2.3 Accelerating the ExtraGradient Algorithm, Direct Approach

The convergence rate of the EG algorithm is far from optimal for the strongly monotone VI problem in (1) with separable structure (2). Firstly, the update rule in (7) takes  $\mathcal{W}$  as a whole without utilizing the separable structure. This oversight prevents us from exploiting the properties of  $\nabla \mathcal{F}$ . Secondly, in the case of bilinear games, the established lower bound for EG is  $\Omega(\sqrt{\kappa} \log(1/\varepsilon))$  rather than  $\Omega(\kappa \log(1/\varepsilon))$ . This discrepancy highlights the potential for accelerating the EG algorithm in various directions. We first rewrite the EG update rule in (7) as follows:

$$\begin{aligned} \mathbf{z}_{t-\frac{1}{2}} &= \mathbf{z}_{t-1} - \eta \mathcal{W}(\mathbf{z}_{t-1}) = \mathbf{z}_{t-1} - \eta(\mathcal{H}(\mathbf{z}_{t-1}) + \nabla \mathcal{F}(\mathbf{z}_{t-1})) \\ \mathbf{z}_t &= \mathbf{z}_{t-1} - \eta \mathcal{W}(\mathbf{z}_{t-\frac{1}{2}}) = \mathbf{z}_{t-1} - \eta(\mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}) + \nabla \mathcal{F}(\mathbf{z}_{t-\frac{1}{2}})) \end{aligned} \quad (8)$$

To accelerate the process based on  $\nabla \mathcal{F}$ , we consider the Nesterov's second acceleration scheme on minimizing a single convex function  $\mathcal{F}$  [Tse08, LZ18, LLF20]:

$$\mathbf{z}_{t-1}^{\text{md}} = (1 - \alpha_t) \mathbf{z}_{t-1}^{\text{ag}} + \alpha_t \mathbf{z}_{t-1}, \quad \mathbf{z}_t = \mathbf{z}_{t-1} - \frac{\eta}{\alpha_t} \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \quad \mathbf{z}_t^{\text{ag}} = (1 - \alpha_t) \mathbf{z}_{t-1}^{\text{ag}} + \alpha_t \mathbf{z}_t \quad (9)$$

where  $\alpha_t$  is extrapolation stepsize in standard three-line Nesterov's scheme. Here we adopt the notations  $\mathbf{z}^{\text{md}}$  and  $\mathbf{z}^{\text{ag}}$  to indicate the middle (or extrapolated) point and the aggregated point [CLO17], respectively. Next, to achieve acceleration, we replace the gradient of  $\nabla \mathcal{F}$  evaluated at both  $\mathbf{z}_{t-1}$  and  $\mathbf{z}_{t-\frac{1}{2}}$  in (8) by the gradient evaluated at  $\mathbf{z}_{t-1}^{\text{md}}$  in (9). Furthermore, we shift the index of  $\mathbf{z}^{\text{ag}}$  by  $\frac{1}{2}$  to indicate the use of  $\mathbf{z}_{t-\frac{1}{2}}$  instead of  $\mathbf{z}_t$  in the  $\mathbf{z}^{\text{ag}}$  update in (9). In addition, we take into account the  $\mu$ -strongly convexity of  $\mathcal{F}$  and shift the gradient of the strongly-convex part  $\frac{\mu}{2} \|\mathbf{z} - \mathbf{z}_0\|^2$  from  $\nabla \mathcal{F}(\mathbf{z})$  to  $\mathcal{H}(\mathbf{z})$  as

$$\mathcal{W}(\mathbf{z}) = (\nabla \mathcal{F}(\mathbf{z}) - \mu(\mathbf{z} - \mathbf{z}_0)) + (\mathcal{H}(\mathbf{z}) + \mu(\mathbf{z} - \mathbf{z}_0))$$

we obtain the following update rule for a direct approach of accelerated EG algorithm (different stepsize schemes on  $\eta_t$  are required for different algorithmic designs):

$$\begin{cases} \mathbf{z}_{t-1}^{\text{md}} = (1 - \alpha_t) \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t \mathbf{z}_{t-1}, \\ \mathbf{z}_{t-\frac{1}{2}} = \mathbf{z}_{t-1} - \eta_t \left( \mathcal{H}(\mathbf{z}_{t-1}) + \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) - \mu(\mathbf{z}_{t-1}^{\text{md}} - \mathbf{z}_{t-1}) \right), \\ \mathbf{z}_t = \mathbf{z}_{t-1} - \eta_t \left( \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}) + \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) - \mu(\mathbf{z}_{t-1}^{\text{md}} - \mathbf{z}_{t-\frac{1}{2}}) \right), \\ \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} = (1 - \alpha_t) \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t \mathbf{z}_{t-\frac{1}{2}} \end{cases} \quad (10)$$

We call the algorithm in (10) as *accelerated gradient-extragradient, direct approach* (AG-EG-Direct), and postpone its full description to Algorithm 2 in §C.1. The final output of the direct approach is  $\mathbf{z}_{\mathcal{T}}$  after  $\mathcal{T}$  iterates. The following theorem states the convergence rate and iteration complexity of AG-EG (direct approach).

**Theorem 1** (Convergence of stochastic AG-EG, direct approach). *Suppose Assumptions 1 and 2 hold. Fix any  $r \in (0, 1)$ ,  $\beta \in (0, \infty)$ , let  $\kappa_\beta = \frac{L}{\mu} + \frac{(1+\beta)M^2}{\mu^2}$  and set the stepsize upper bound  $\bar{\alpha} \equiv \frac{r}{1+\sqrt{1+r\kappa_\beta}}$ . For any sequence of stepsizes  $\alpha_t \in (0, \bar{\alpha}]$  and  $\eta_t = \frac{\alpha_t}{\mu}$ , the iterates of stochastic AG-EG (direct approach) satisfy that for all  $t = 1, \dots, \mathcal{T}$*

$$\mathbb{E} \|\mathbf{z}_t - \mathbf{z}^*\|^2 \leq \|\mathbf{z}_0 - \mathbf{z}^*\|^2 \left( \frac{L}{\mu} + 1 \right) \prod_{s=1}^t (1 - \alpha_s) + \frac{3\sigma^2}{\mu^2} \sum_{s=1}^t \alpha_s^2 \prod_{\tau=s+1}^t (1 - \alpha_\tau) \quad (11)$$

where we define  $\sigma = \frac{1}{\sqrt{3}} \sqrt{\frac{1}{1-r} \sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta}) \sigma_{\text{Bil}}^2}$  such that  $\sigma^2$  is the weighted, uniform variance bound on the stochastic gradients.

Throughout all sections afterward, we use the same definition  $\sigma$  as in Theorem 1. The proof of Theorem 1 is provided in §D.4. We further note that one possible stepsize choice is to let  $\alpha_t \equiv \alpha$ , and Eq. (11) reduces to

$$\mathbb{E} \|\mathbf{z}_t - \mathbf{z}^*\|^2 \leq \|\mathbf{z}_0 - \mathbf{z}^*\|^2 \left( \frac{L}{\mu} + 1 \right) e^{-\alpha t} + \frac{3\sigma^2}{\mu^2} \alpha$$

For any given  $\mathcal{T} \geq 1$ , by choosing the optimal  $\alpha = \frac{1}{\mathcal{T}} \left( 1 + \log \left( \frac{\mu^2 \mathcal{T}}{3\sigma^2} \left( \frac{L}{\mu} + 1 \right) \|\mathbf{z}_0 - \mathbf{z}^*\|^2 \right) \right) \wedge \bar{\alpha}$ , Eq. (11) implies

$$\mathbb{E} \|\mathbf{z}_{\mathcal{T}} - \mathbf{z}^*\|^2 \leq \|\mathbf{z}_0 - \mathbf{z}^*\|^2 \left( \frac{L}{\mu} + 1 \right) e^{-\bar{\alpha} \mathcal{T}} + \frac{3\sigma^2}{\mu^2 \mathcal{T}} \left( 1 + \log \left( \frac{\mu^2 \mathcal{T}}{3\sigma^2} \left( \frac{L}{\mu} + 1 \right) \|\mathbf{z}_0 - \mathbf{z}^*\|^2 \right) \right)$$

Prescribing the desired accuracy  $\varepsilon > 0$ , the iteration complexity to output an  $\varepsilon$ -optimal minimax point is<sup>3</sup>

$$\mathcal{O} \left( \left( \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} + \frac{\sigma^2}{\mu^2 \varepsilon^2} \right) \log \left( \left( \frac{L}{\mu} + 1 \right) \|\mathbf{z}_0 - \mathbf{z}^*\|^2 / \varepsilon^2 \right) \right)$$

We conjecture that the logarithmic factor in the optimal statistical rate  $\frac{\sigma^2}{\mu^2 \varepsilon^2}$  is removable using a proper diminishing stepsize, a possibility that we reserve for future study. In the setting of

<sup>3</sup>Throughout this work, we focus on the iteration complexity whereas the required number of queries to the stochastic gradient oracle is three times the iteration complexity.

nonstochastic variational inequality, setting  $\sigma = 0$  and  $r \rightarrow 1^-$ ,  $\beta \rightarrow 0^+$  in Theorem 1, we obtain the optimal iteration complexity bound as follows:

$$\sim \left(1 + \sqrt{1 + \frac{L}{\mu} + \frac{M^2}{\mu^2}}\right) \log \left(\left(\frac{L}{\mu} + 1\right)/\varepsilon^2\right) \quad (12)$$

**Remark.** Our complexity bounds *fundamentally differs* from the previous analysis [CLO17, JLZ23] for separable smooth (strongly) monotone VIs. The convergence results in previous studies are dependent on the diameter of the domain, whereas our convergence rate is independent of the domain parameters and eliminates the need for projection onto a bounded domain. Moreover, our contributions go beyond those of [CLO17] by extending the analysis to the strongly monotone case. In comparison with [JLZ23], we design an algorithm when  $\nabla \mathcal{F}$  is strongly monotone and resolve the open problem of extending the analysis to the stochastic case. Additionally, our complexity bound in (12) indicates a near-unity coefficient on the condition-number exponent, improving the corresponding coefficient in [JLZ23, Theorem 15] by an asymptotic factor of 4.

The direct approach, which reduces to EG when  $\nabla \mathcal{F} = 0$  and  $\mu = 0$ , falls short of attaining optimality within the specific regime of bilinear games. In the next subsection, we will introduce a new algorithm that can overcome this limitation.

## 2.4 Accelerating the ExtraGradient Algorithm with Scheduled Restarting

In this subsection, we solve problem (1) by further accelerating the stochastic EG algorithm. Rather than directly relying on the strongly-monotonicity of  $\nabla \mathcal{F}$ , the inner updates of our new algorithm is identical to the updates in (10) with  $\mu = 0$ . Due to the domain-independent nature of our analysis, we can apply the scheduled restarting technique [OC15, Rd17, RG22] to the outer loop, accelerating the algorithm from sublinear convergence to linear convergence. In addition, the output of our algorithm is the aggregated point  $\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}}$  after  $\mathcal{T}$  iterates. We present the full algorithm in Algorithm 1.

We first present the convergence rate of a single-epoch (i.e., the inner loop) of Algorithm 1 in Theorem 2. To accommodate more flexibility in the choice of parameters, we introduce three constants  $r, \beta$ , and  $C$  in the theorem statement.

**Theorem 2** (Convergence of stochastic AG-EG, one epoch). *Suppose Assumptions 1 and 2 hold. For any fixed epoch length  $\mathcal{T} \geq 1$ , any constant  $r \in (0, 1)$ ,  $\beta \in (0, \infty)$ ,  $C \in (0, \infty)$ , choose stepsizes  $\alpha_t = \frac{2}{t+1}$  and  $\eta_t$  such that*

$$\frac{t}{\eta_t} = \frac{2}{r}L \vee \mathcal{B} + \sqrt{\frac{1+\beta}{r}}Mt \quad (13)$$

where  $\mathcal{B} = \frac{\sigma\sqrt{\mathcal{T}(\mathcal{T}+1)}}{C\sqrt{\mathbb{E}\|\mathbf{z}_0 - \mathbf{z}^*\|^2}}$ . The output  $\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}}$  of a single-epoch of Algorithm 1 satisfies

$$\mathbb{E} \left\| \mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} - \mathbf{z}^* \right\|^2 \leq \frac{2}{\mu(\mathcal{T}+1)} \left( \frac{2L}{r\mathcal{T}} + \mathcal{A} \sqrt{\frac{1+\beta}{r}} M \right) \mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2 + \frac{2(\frac{1}{C} + C)\sigma}{\mu\sqrt{\mathcal{T}}} \sqrt{\mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2} \quad (14)$$

where the prefactor  $\mathcal{A} \equiv 1 + C^2\mathcal{B}\eta_1 \leq 1 + C^2$  reduces to 1 when  $\sigma = 0$ .

The proof of Theorem 2 is provided in §D.3. We make a few remarks on Theorem 2 as follows:

---

**Algorithm 1** Stochastic Accelerated Gradient-Extra Gradient (AG-EG) Descent-Ascent Algorithm, with Scheduled Restarting

---

**Require:** Initialization  $\mathbf{z}_0^{[0]}$ , total number of epochs  $\mathcal{S} \geq 1$ , total number of per-epoch iterates  $(\mathcal{T}_s : s = 1, \dots, \mathcal{S})$ , stepsizes  $(\alpha_t, \eta_t : t = 1, 2, \dots)$

**for**  $s = 1, 2, \dots, \mathcal{S}$  **do**

    Set  $\mathbf{z}_{-\frac{1}{2}}^{\text{ag}} \leftarrow \mathbf{z}_0^{[s-1]}, \mathbf{z}_0 \leftarrow \mathbf{z}_0^{[s-1]}, \mathbf{z}_0^{\text{md}} \leftarrow \mathbf{z}_0^{[s-1]}$

**for**  $t = 1, 2, \dots, \mathcal{T}_s$  **do**

        Draw samples  $\xi_{t-\frac{1}{2}} \sim \mathcal{D}_\xi$  from oracle, and also  $\zeta_{t-\frac{1}{2}}, \zeta_t \sim \mathcal{D}_\zeta$  independently from oracle

$\mathbf{z}_{t-\frac{1}{2}} \leftarrow \mathbf{z}_{t-1} - \eta_t \left( \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}}) + \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right)$

$\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \leftarrow (1 - \alpha_t) \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t \mathbf{z}_{t-\frac{1}{2}}$

$\mathbf{z}_t \leftarrow \mathbf{z}_{t-1} - \eta_t \left( \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) + \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right)$

$\mathbf{z}_t^{\text{md}} \leftarrow (1 - \alpha_{t+1}) \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} + \alpha_{t+1} \mathbf{z}_t$

**end for**

    Set  $\mathbf{z}_0^{[s]} \leftarrow \mathbf{z}_{\mathcal{T}_s-\frac{1}{2}}^{\text{ag}}$  {//Warm-start using the output of the previous epoch}

**end for**

**Output:**  $\mathbf{z}_0^{[\mathcal{S}]}$

---

**Remark.** In the setting of deterministic optimization, by taking  $\sigma = 0, r \rightarrow 1^-, \beta \rightarrow 0^+$  in our analysis, with stepsize choice  $\eta_t = \frac{t}{2L+Mt}$ , we obtain that

$$\|\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} - \mathbf{z}^*\|^2 \leq \frac{2}{\mu(\mathcal{T}+1)} \left( \frac{2L}{\mathcal{T}} + M \right) \|\mathbf{z}_0 - \mathbf{z}^*\|^2 \quad (15)$$

Under this circumstance, the algorithm is independent of  $\mathcal{B}$  and requires no knowledge of  $\|\mathbf{z}_0 - \mathbf{z}^*\|^2$ . In the face of stochasticity, we choose  $C = 1$  when the initial distance to the optimal point is known. Alternatively, when only an over-estimate  $\Gamma_0$  of  $\sqrt{\mathbb{E}\|\mathbf{z}_0 - \mathbf{z}^*\|^2}$  is provided, we can set (large enough)  $C = \frac{\Gamma_0}{\sqrt{\mathbb{E}\|\mathbf{z}_0 - \mathbf{z}^*\|^2}} \geq 1$  to obtain

$$\mathbb{E}\|\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} - \mathbf{z}^*\|^2 \leq \frac{2}{\mu(\mathcal{T}+1)} \left( \frac{2L}{r\mathcal{T}} + 2\sqrt{\frac{1+\beta}{r}}M \right) \Gamma_0^2 + \frac{4\sigma}{\mu\sqrt{\mathcal{T}}} \Gamma_0 \quad (16)$$

**Remark.** When the constants are not a concern, the coarse-grained choices of  $r = \frac{1}{2}$  and  $\beta = 1$  would suffice. Nevertheless, to optimize the constants, the tradeoff between the deviation of  $r$  from 1 and  $\beta$  from 0 is crucial, as it determines a balance between the stochastic gradient noise variance and the convergence rate coefficients.

To prepare for our multi-epoch result with the help of scheduled restarting, we perform an induction based on (16) as follows: suppose  $\mathbb{E}\|\mathbf{z}_0^{[s-1]} - \mathbf{z}^*\|^2 \leq \Gamma_0^2 e^{1-s}$  hold, by taking  $r = \frac{1}{2}$  and  $\beta = 1$ , we have

$$\mathbb{E}\|\mathbf{z}_0^{[s]} - \mathbf{z}^*\|^2 \lesssim \frac{L}{\mu\mathcal{T}_s^2} \Gamma_0^2 e^{1-s} + \frac{M}{\mu\mathcal{T}_s} \Gamma_0^2 e^{1-s} + \frac{\sigma}{\mu\sqrt{\mathcal{T}_s}} \Gamma_0 e^{\frac{1-s}{2}}$$

Setting the right-hand-side of the above inequality to satisfy  $\leq \Gamma_0^2 e^{-s}$ , and solving for  $\mathcal{T}_s$ , we need the epoch length satisfies  $\mathcal{T}_s \asymp \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} + \frac{\sigma^2}{\mu^2 \Gamma_0^2 e^{1-s}}$ . Thus, we can obtain the total iteration



complexity as

$$\sum_{s=1}^S \left[ \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} + \frac{\sigma^2}{\mu^2 \Gamma_0^2 e^{1-s}} \right] = \left( \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} \right) S + \frac{\sigma^2}{\mu^2 \Gamma_0^2} \cdot \frac{e^S - 1}{e - 1}$$

where  $S \equiv \left\lceil \log \frac{\Gamma_0^2}{\varepsilon^2} \right\rceil$ . This yields the following multi-epoch iteration complexity bound:

**Corollary 1** (Iteration complexity of stochastic AG-EG with scheduled restarting). *Under the same condition of Theorem 2, the stochastic AG-EG with scheduled restarting in Algorithm 1 with epoch length  $\mathcal{T}_s \asymp \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} + \frac{\sigma^2}{\mu^2 \Gamma_0^2 e^{1-s}}$  has a total iteration complexity of*

$$\mathcal{O} \left( \left( \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} \right) \log \left( \frac{1}{\varepsilon} \right) + \frac{\sigma^2}{\mu^2 \varepsilon^2} \right) \quad (17)$$

In the deterministic setting, the hard instance constructed by [ZHZ22] can be seen as a special case of a monotone VI (1), which implies a lower bound of  $\Omega \left( \left( \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} \right) \log \left( \frac{1}{\varepsilon} \right) \right)$ . This demonstrates the optimality of the complexity provided by Corollary 1. It is worth noting that while both complexity bounds in Corollary 1 and Theorem 1 match the lower bound in [ZHZ22] for strongly monotone VIs with separable structure, the direct approach in §2.3 reduces to the *last-iterate* independent-sample stochastic extragradient (SEG) algorithm in *bilinear games*. Consequently, the deterministic part ( $\sigma = 0$ ) fails to match the lower bound in [IAGM20]. In the stochastic case with noise variance bounded away from zero, the direct approach in §2.3 can exhibit *non-convergence behavior* [HIMM20, §3]. The AG-EG algorithm in §2.4 resolves this issue by restarting the *average-iterate* SEG, matching the lower bound results (see §3.2 for more details). In addition, the complexity bound in (17) also eliminates the log prefactor of the statistical error term  $\frac{\sigma^2}{\mu^2 \varepsilon^2}$  compared to Theorem 1.

## 2.5 Extension of AG-EG to Proximal Algorithms

In previous subsections, we focus on the case where the feasible set  $\mathcal{Z}$  represents the entire space and the non-differentiable convex function  $J$  is dropped. In this subsection, we extend the AG-EG algorithm and its analysis to the more general setting that has a bounded feasible set (via Euclidean projection onto the feasible set) as well as a non-differentiable convex regularization term (i.e., via proximal operator). In particular, these settings are useful in various applications, such as the variational inequality on the Lorentz cone where projection onto  $\mathcal{Z} = \{(\mathbf{x}, t) \in \mathbb{R}^{(n+1)} : \|\mathbf{x}\| \leq t\}$  is required [CLO17], and the two-player game that involves projection onto the probability simplex, among others. To deal with bounded feasible set  $\mathcal{Z}$ , we adopt a variant of the EG algorithm, where we project the extrapolated point and the main iterates back onto the feasible set  $\mathcal{Z}$  of  $\mathcal{W}$ .

$$\begin{aligned} \mathbf{z}_{t-\frac{1}{2}} &= P_{\mathcal{Z}} [\mathbf{z}_{t-1} - \eta \mathcal{W}(\mathbf{z}_{t-1})] = \operatorname{argmin}_{\mathbf{z} \in \mathcal{Z}} \langle \mathbf{z} - \mathbf{z}_{t-1}, \eta \mathcal{W}(\mathbf{z}_{t-1}) \rangle + \frac{1}{2} \|\mathbf{z} - \mathbf{z}_{t-1}\|^2 \\ \mathbf{z}_t &= P_{\mathcal{Z}} \left[ \mathbf{z}_{t-1} - \eta \mathcal{W}(\mathbf{z}_{t-\frac{1}{2}}) \right] = \operatorname{argmin}_{\mathbf{z} \in \mathcal{Z}} \left\langle \mathbf{z} - \mathbf{z}_{t-1}, \eta \mathcal{W}(\mathbf{z}_{t-\frac{1}{2}}) \right\rangle + \frac{1}{2} \|\mathbf{z} - \mathbf{z}_{t-1}\|^2 \end{aligned} \quad (18)$$

where  $P_{\mathcal{Z}}(\mathbf{z}) = \operatorname{argmin}_{\mathbf{z}' \in \mathcal{Z}} \|\mathbf{z} - \mathbf{z}'\|^2$  is the Euclidean projection operator. To handle non-differentiable simple convex function  $J$ , we can replace projection operator in (18) by the following proximal

mapping defined via a Bregman divergence  $\mathcal{D}(\cdot, \cdot)$

$$\text{prox}_{\mathcal{Z}}^J(\mathbf{v}) \equiv \underset{\mathbf{u} \in \mathcal{Z}}{\text{argmin}} \langle \mathbf{v}, \mathbf{u} - \mathbf{z} \rangle + \mathcal{D}(\mathbf{z}, \mathbf{u}) + J(\mathbf{u}) \quad (19)$$

In fact, (18) can be seen as a special case of (19) when choosing the Bregman divergence  $\mathcal{D}(\mathbf{z}, \mathbf{u}) = 1/2\|\mathbf{z} - \mathbf{u}\|^2$  and  $J(\mathbf{u})$  as the set indicator function of the feasible set  $\mathcal{Z}$ . Therefore, by substituting the prox-mapping (19) into the AG-EG updates introduced in §2.4, we obtain the more general proximal AG-EG algorithm in Algorithm 3 (See in §C.2), which reduces to Algorithm 1 when  $J = 0$ ,  $\mathcal{D}(\mathbf{z}, \mathbf{u}) = 1/2\|\mathbf{z} - \mathbf{u}\|^2$  and  $\mathcal{Z} = \mathbf{R}^n$ . Similar to Corollary 1, we have the following iteration complexity result, whose proof is deferred to §D:

**Corollary 2** (Iteration complexity of stochastic proximal AG-EG with scheduled restarting). *Under the same condition of Theorem 2, the stochastic proximal AG-EG with scheduled restarting in Algorithm 3 with epoch length  $\mathcal{T}_s \asymp \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} + \frac{\sigma^2}{\mu^2 \Gamma_0^2 e^{1-s}}$  has a total iteration complexity of*

$$\mathcal{O} \left( \left( \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} \right) \log \left( \frac{1}{\varepsilon} \right) + \frac{\sigma^2}{\mu^2 \varepsilon^2} \right)$$

For the deterministic case, proximal AG-EG with scheduled restarting has a total iteration complexity of  $\mathcal{O} \left( \left( \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} \right) \log \left( \frac{1}{\varepsilon} \right) \right)$  to output an  $\varepsilon$ -optimal point of (1).

### 3 Implications for Specific Instances

In this section, we discuss the implications of our AG-EG algorithm and its convergence rates when applying to two instances of saddle-point problem.

#### 3.1 Strongly-Convex-Strongly-Concave Saddle Point Problem

For the stochastic bilinearly-coupled SC-SC saddle point problem (3), we note that the smoothness and strong convexity parameters  $L_F$ ,  $L_G$ ,  $\mu_F$ , and  $\mu_G$  of  $F$  and  $G$  may differ. To accommodate these variations in curvature information, we employ a scaling reduction technique. This technique enables us to convert the SCSC with equal strong convexity parameters for  $F$  and  $G$  by reparametrizing the objective function. The same argument is also applicable to the direct approach.

In lieu to (3), we consider

$$\min_{\hat{\mathbf{x}}} \max_{\hat{\mathbf{y}}} \widehat{\mathcal{F}}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = F(\hat{\mathbf{x}}) + \widehat{H}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) - \widehat{G}(\hat{\mathbf{y}})$$

where  $\widehat{\mathcal{F}}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \mathcal{F}(\mathbf{x}, \mathbf{y})$  with the symbolic reparametrization  $\hat{\mathbf{x}} = \mathbf{x}$ ,  $\hat{\mathbf{y}} = \sqrt{\frac{\mu_G}{\mu_F}} \mathbf{y}$ ,  $\widehat{H}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = H(\mathbf{x}, \mathbf{y})$ ,  $\widehat{G}(\hat{\mathbf{y}}) = G(\mathbf{y})$  and also their derivatives  $\nabla_{\hat{\mathbf{y}}} \widehat{H}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \sqrt{\frac{\mu_F}{\mu_G}} \nabla_{\mathbf{y}} H(\mathbf{x}, \mathbf{y})$ ,  $\nabla \widehat{G}(\hat{\mathbf{y}}) = \sqrt{\frac{\mu_F}{\mu_G}} \nabla G(\mathbf{y})$  (the stochastic oracles  $\widehat{h}, \widehat{g}$  follows the same rule). It is straightforward to verify that  $\widehat{\mathcal{F}}(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  is  $\mu$ -strongly-convex- $\mu$ -strongly-concave. The essence of our update rules can be summarized by the rescaled updates on  $\mathbf{y}$ , as illustrated below:

$$\begin{aligned} \hat{\mathbf{y}}_t &= \hat{\mathbf{y}}_{t-1} - \eta_t \left( -\nabla_{\hat{\mathbf{y}}} h(\hat{\mathbf{x}}_{t-\frac{1}{2}}, \hat{\mathbf{y}}_{t-\frac{1}{2}}; \zeta_t) + \nabla g(\hat{\mathbf{y}}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right) \\ \Leftrightarrow \mathbf{y}_t &= \mathbf{y}_{t-1} - \eta_t \cdot \frac{\mu_F}{\mu_G} \left( -\nabla_{\mathbf{y}} h(\mathbf{x}_{t-\frac{1}{2}}, \mathbf{y}_{t-\frac{1}{2}}; \zeta_t) + \nabla g(\mathbf{y}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right) \end{aligned}$$

Therefore, it suffices to analyze Algorithm 3 for  $\widehat{\mathcal{F}}(\widehat{\mathbf{x}}, \widehat{\mathbf{y}})$  and due to this scaling reduction, we only need to prove all results for the case of  $\mu_F = \mu_G = \mu$ . It is also straightforward to justify corresponding scaling changes as:  $L = L_F \vee \frac{\mu_F}{\mu_G} L_G$ ,  $M = \sqrt{\frac{\mu_F}{\mu_G} \lambda_{\max}(\mathbf{B}^\top \mathbf{B})}$ , and  $\mu = \mu_F$ . The following corollary is rediscovered by reverting the scaling reduction from  $\widehat{\mathcal{F}}(\widehat{\mathbf{x}}, \widehat{\mathbf{y}})$  to  $\mathcal{F}(\mathbf{x}, \mathbf{y})$ .

**Corollary 3** (Iteration complexity of stochastic AG-EG on SC-SC saddle-point problem). *For solving Eq. (3), Algorithm 1 with an epoch length  $\mathcal{T}_s \asymp \sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}} + \frac{\sigma^2}{\mu_F^2 \Gamma_0^2 e^{1-s}}$  has a total iteration complexity of*

$$\mathcal{O} \left( \left( \sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}} \right) \log \left( \frac{1}{\varepsilon} \right) + \frac{\sigma^2}{\mu_F^2 \varepsilon^2} \right)$$

In the deterministic case, the iteration complexity in Theorem 1 matches the lower bound established by [ZHZ22], i.e.,  $\Omega \left( \left( \sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}} \right) \log \left( \frac{1}{\varepsilon} \right) \right)$ . Moreover, our algorithm achieves the optimal statistical rate of  $\frac{\sigma^2}{\mu_F^2 \varepsilon^2}$  up to a constant prefactor.

**Remark.** A well-known finding regarding the second scheme of Nesterov’s acceleration is its connection to the primal-dual method [LZ18, LLF20]. This finding has been incorporated into the design of the LPD algorithm [THO22], where a Chambolle-Pock-style primal-dual method is utilized as an approximation of proximal point methods, instead of the extragradient used in this paper. The LPD algorithm [THO22] also achieves the optimal complexity for the deterministic bilinearly-coupled saddle-point problem.

### 3.2 Bilinear Games

In this subsection, we consider the particular case of bilinear games. We assume  $n = m$  such that  $\mathbf{B}$  is a nonsingular square matrix,  $\nabla f(\mathbf{x}; \xi) = \mathbf{0}$  and  $\nabla g(\mathbf{y}; \xi) = \mathbf{0}$  a.s., so (3) reduces to

$$\min_{\mathbf{x}} \max_{\mathbf{y}} \mathcal{F}(\mathbf{x}, \mathbf{y}) = \mathbb{E}_\zeta [h(\mathbf{x}, \mathbf{y}; \zeta)] = H(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{B} \mathbf{y} - \mathbf{x}^\top \mathbf{u}_x + \mathbf{u}_y^\top \mathbf{y} \quad (20)$$

and Algorithm 3 reduces to the independent-sample extragradient descent-ascent algorithm for (20). The saddle point  $[\omega_x^*; \omega_y^*]$  in this case is the unique solution to the linear equation

$$\begin{bmatrix} \mathbf{0} & \mathbf{B} \\ -\mathbf{B}^\top & \mathbf{0} \end{bmatrix} \begin{bmatrix} \omega_x^* \\ \omega_y^* \end{bmatrix} = \begin{bmatrix} \mathbf{u}_x \\ \mathbf{u}_y \end{bmatrix}, \quad \text{which has a closed-form solution } \begin{bmatrix} \omega_x^* \\ \omega_y^* \end{bmatrix} = \begin{bmatrix} -(\mathbf{B}^\top)^{-1} \mathbf{u}_y \\ \mathbf{B}^{-1} \mathbf{u}_x \end{bmatrix}$$

Our results imply the following iteration complexity for solving stochastic bilinear games.

**Corollary 4** (Iteration complexity of stochastic AG-EG, bilinear games). *For solving Eq. (20), choose the stepsizes  $\alpha_t = \frac{2}{t+1}$  and  $\eta_t \equiv \frac{1}{\sqrt{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}}$ , Algorithm 1 with an epoch length  $\mathcal{T}_s \asymp \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B} \mathbf{B}^\top)}}$  has the total iteration complexity of*

$$\mathcal{O} \left( \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B} \mathbf{B}^\top)}} \log \left( \frac{\sqrt[4]{\lambda_{\min}(\mathbf{B} \mathbf{B}^\top) \lambda_{\max}(\mathbf{B}^\top \mathbf{B})}}{\sigma_{\text{Bil}}} \right) + \frac{\sigma_{\text{Bil}}^2}{\lambda_{\min}(\mathbf{B} \mathbf{B}^\top) \varepsilon^2} \right) \quad (21)$$

Note that our choice of the stepsize is maximal and is independent of the noise. In the deterministic setting, let  $\sigma_{\text{Bil}} \asymp \varepsilon \sqrt[4]{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)\lambda_{\max}(\mathbf{B}^\top\mathbf{B})}$ , the complexity bound in Corollary 4 reduces to  $\mathcal{O}\left(\sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top\mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}} \log\left(\frac{1}{\varepsilon}\right)\right)$ , which matches the lower bound in [IAGM20]. Notably, [ASM<sup>+</sup>20] proposed an algorithm achieving an upper bound that matches the lower bound in [IAGM20].<sup>4</sup> [LYL<sup>+</sup>22] also proposed a lower-bound matching SEG algorithm that uses a shared sample in both steps under unbounded noise assumption. In contrast, our algorithm is in the independent-sample setting with bounded noise variance.

**Remark.** Standard acceleration techniques do not attain the optimal nonasymptotic convergence rate for the bilinear minimax game [GHP<sup>+</sup>19]. This limitation applies to various algorithms, including the direct approach [§2.3], as well as several other acceleration techniques [THO22, KGR22, JST22], all of which fall short in achieving optimal acceleration for bilinear games. Therefore, matching both lower bounds in a single algorithm in the general stochastic setting has been an open problem. While [LYGJ22] presents an algorithm that achieves both lower bounds in a single algorithm, it relies on the use of optimistic gradients rather than extragradients on the bilinear coupling function. Furthermore, our algorithm and analysis is more general than those in [LYGJ22] as we can handle the general variational inequality with proximal operators.

## 4 Conclusion and Future Work

We have presented a stochastic extragradient-based acceleration algorithm, AG-EG, for solving stochastic monotone variational inequalities with separable structure. The iteration complexity of our algorithm matches the lower bound and is independent of the size of the feasible set. When specialized to solving the bilinearly coupled saddle-point problem (3), our AG-EG algorithm simultaneously matches lower bounds due to [ZHJ22] and [IAGM20] for strongly-convex-strongly-concave and bilinear games, respectively. To the best of our knowledge, this is the first time that all three lower bounds have been met by a single algorithm. There are some remaining issues to be addressed, however, including the case of one-sided non-strong convexity, the setting of unbounded noise variance, and the characterization of the full parameter regime dependency on  $\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)$ . These are left as important directions for future research.

## Acknowledgements

This work is supported in part by NSF Award’s IIS-2110170 and DMS-2134106 to SSD, by Canada CIFAR AI Chair to GG, by the Mathematical Data Science program of the Office of Naval Research under grant number N00014-18-1-2764 and also the Vannevar Bush Faculty Fellowship program under grant number N00014-21-1-2941 and NSF grant IIS-1901252 to MIJ.

## References

[AMLJG20] Waïss Azizian, Ioannis Mitliagkas, Simon Lacoste-Julien, and Gauthier Gidel. A tight and unified analysis of gradient-based methods for a whole spectrum of differentiable games. In

---

<sup>4</sup>[ASM<sup>+</sup>20], while achieving optimality in bilinear games, has a flavor of “aggressive extrapolation” [HIMM20] that results in a convergence behavior resembling accelerated Hamiltonian gradient descent. Thus, it cannot be tuned to match the complexity lower bound for the bilinearly coupled SC-SC case.

- International Conference on Artificial Intelligence and Statistics*, pages 2863–2873. PMLR, 2020.
- [ASM<sup>+</sup>20] Waïss Azizian, Damien Scieur, Ioannis Mitliagkas, Simon Lacoste-Julien, and Gauthier Gidel. Accelerating smooth games by manipulating spectral shapes. In *International Conference on Artificial Intelligence and Statistics*, pages 1705–1715. PMLR, 2020.
- [CLO14] Yunmei Chen, Guanghui Lan, and Yuyuan Ouyang. Optimal primal-dual methods for a class of saddle point problems. *SIAM Journal on Optimization*, 24(4):1779–1814, 2014.
- [CLO17] Yunmei Chen, Guanghui Lan, and Yuyuan Ouyang. Accelerated schemes for a class of variational inequalities. *Mathematical Programming*, 165(1):113–149, 2017.
- [CLZ21] Ziyi Chen, Qunwei Li, and Yi Zhou. Finding local minimax points via (stochastic) cubic-regularized gda: Global convergence and complexity. *arXiv preprint arXiv:2110.07098*, 2021.
- [CST21] Michael B Cohen, Aaron Sidford, and Kevin Tian. Relative Lipschitzness in extragradient methods and a direct recipe for acceleration. In *12th Innovations in Theoretical Computer Science Conference (ITCS 2021)*, volume 185, page 62, 2021.
- [DCL<sup>+</sup>17] Simon S Du, Jianshu Chen, Lihong Li, Lin Xiao, and Dengyong Zhou. Stochastic variance reduction methods for policy evaluation. In *International Conference on Machine Learning*, pages 1049–1058. PMLR, 2017.
- [DISZ18] Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training GANs with optimism. In *International Conference on Learning Representations*, 2018.
- [DNP<sup>+</sup>14] Christoph Dann, Gerhard Neumann, Jan Peters, et al. Policy evaluation with temporal differences: A survey and comparison. *Journal of Machine Learning Research*, 15:809–883, 2014.
- [FOP20] Alireza Fallah, Asuman Ozdaglar, and Sarath Pattathil. An optimal multistage stochastic gradient method for minimax problems. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pages 3573–3579. IEEE, 2020.
- [FP03] Francisco Facchinei and Jong-Shi Pang. *Finite-Dimensional Variational Inequalities and Complementarity Problems*. Springer, 2003.
- [GBV<sup>+</sup>19] Gauthier Gidel, Hugo Berard, Gaëtan Vignoud, Pascal Vincent, and Simon Lacoste-Julien. A variational inequality perspective on generative adversarial networks. In *International Conference on Learning Representations*, 2019.
- [GHP<sup>+</sup>19] Gauthier Gidel, Reyhane Askari Hemmat, Mohammad Pezeshki, Rémi Le Priol, Gabriel Huang, Simon Lacoste-Julien, and Ioannis Mitliagkas. Negative momentum for improved game dynamics. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1802–1811. PMLR, 2019.
- [GL12] Saeed Ghadimi and Guanghui Lan. Optimal stochastic approximation algorithms for strongly convex stochastic composite optimization i: A generic algorithmic framework. *SIAM Journal on Optimization*, 22(4):1469–1492, 2012.
- [GPAM<sup>+</sup>20] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [GYYY20] Zhishuai Guo, Zhuoning Yuan, Yan Yan, and Tianbao Yang. Fast objective & duality gap convergence for nonconvex-strongly-concave min-max problems. *arXiv preprint arXiv:2006.06889*, 2020.
- [HIMM20] Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. Explore aggressively, update conservatively: Stochastic extragradient methods with variable stepsize scaling. In *Advances in Neural Information Processing Systems*, volume 33, pages 16223–16234, 2020.

- [IAGM20] Adam Ibrahim, Waiss Azizian, Gauthier Gidel, and Ioannis Mitliagkas. Linear lower bounds and conditioning of differentiable games. In *International Conference on Machine Learning*, pages 4583–4593. PMLR, 2020.
- [IJOT17] Alfredo N Iusem, Alejandro Jofré, Roberto Imbuzeiro Oliveira, and Philip Thompson. Extragradient method with variance reduction for stochastic variational inequalities. *SIAM Journal on Optimization*, 27(2):686–724, 2017.
- [JDES<sup>+</sup>20] Samy Jelassi, Carles Domingo-Enrich, Damien Scieur, Arthur Mensch, and Joan Bruna. Extragradient with player sampling for faster convergence in n-player games. In *International Conference on Machine Learning*, pages 4736–4745. PMLR, 2020.
- [JLZ23] Michael I Jordan, Tianyi Lin, and Manolis Zampetakis. First-order algorithms for nonlinear generalized Nash equilibrium problems. *Journal of Machine Learning Research*, 24(38):1–46, 2023.
- [JNT11] Anatoli Juditsky, Arkadi Nemirovski, and Claire Tauvel. Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems*, 1(1):17–58, 2011.
- [JST22] Yujia Jin, Aaron Sidford, and Kevin Tian. Sharper rates for separable minimax and finite sum optimization via primal-dual extragradient methods. In *Conference on Learning Theory*, pages 4362–4415. PMLR, 2022.
- [KGR22] Dmitry Kovalev, Alexander Gasnikov, and Peter Richtárik. Accelerated primal-dual gradient method for smooth and convex-concave saddle-point problems with bilinear coupling. *Advances in Neural Information Processing Systems*, 35:21725–21737, 2022.
- [KLL20] Georgios Kotsalis, Guanghui Lan, and Tianjiao Li. Simple and optimal methods for stochastic variational inequalities, I: Operator extrapolation. *arXiv preprint arXiv:2011.02987*, 2020.
- [Kor76] Galina M Korpelevich. The extragradient method for finding saddle points and other problems. *Ekonomika i Matematicheskie Metody*, 12:747–756, 1976.
- [KS80] David Kinderlehrer and Guido Stampacchia. *An Introduction to Variational Inequalities and Their Applications*, volume 31. SIAM, 1980.
- [LBJM<sup>+</sup>20] Nicolas Loizou, Hugo Berard, Alexia Jolicoeur-Martineau, Pascal Vincent, Simon Lacoste-Julien, and Ioannis Mitliagkas. Stochastic hamiltonian gradient methods for smooth games. In *International Conference on Machine Learning*, pages 6370–6381. PMLR, 2020.
- [LJJ20a] Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pages 6083–6093. PMLR, 2020.
- [LJJ20b] Tianyi Lin, Chi Jin, and Michael I Jordan. Near-optimal algorithms for minimax optimization. In *Conference on Learning Theory*, pages 2738–2779. PMLR, 2020.
- [LLC22] Luo Luo, Yujun Li, and Cheng Chen. Finding second-order stationary points in nonconvex-strongly-concave minimax optimization. *Advances in Neural Information Processing Systems*, 35:36667–36679, 2022.
- [LLF20] Zhouchen Lin, Huan Li, and Cong Fang. Accelerated optimization for machine learning. *Nature Singapore: Springer*, 2020.
- [LO21] Guanghui Lan and Yuyuan Ouyang. Mirror-prox sliding methods for solving a class of monotone variational inequalities. *arXiv preprint arXiv:2111.00996*, 2021.
- [LS19] Tengyuan Liang and James Stokes. Interaction matters: A note on non-asymptotic local convergence of generative adversarial networks. In *International Conference on Artificial Intelligence and Statistics*, pages 907–915. PMLR, 2019.

- [LYGJ22] Chris Junchi Li, Angela Yuan, Gauthier Gidel, and Michael I Jordan. Nesterov meets optimism: Rate-optimal optimistic-gradient-based method for stochastic bilinearly-coupled minimax optimization. *arXiv preprint arXiv:2210.17550*, 2022.
- [LYL<sup>+</sup>22] Chris Junchi Li, Yaodong Yu, Nicolas Loizou, Gauthier Gidel, Yi Ma, Nicolas Le Roux, and Michael Jordan. On the convergence of stochastic extragradient for bilinear games using restarted iteration averaging. In *International Conference on Artificial Intelligence and Statistics*, pages 9793–9826. PMLR, 2022.
- [LZ18] Guanghui Lan and Yi Zhou. An optimal randomized incremental gradient method. *Mathematical Programming*, 171:167–215, 2018.
- [MKS<sup>+</sup>20] Konstantin Mishchenko, Dmitry Kovalev, Egor Shulgin, Peter Richtárik, and Yura Malitsky. Revisiting stochastic extragradient. In *International Conference on Artificial Intelligence and Statistics*, pages 4573–4582. PMLR, 2020.
- [MLZ<sup>+</sup>18] Panayotis Mertikopoulos, Bruno Lecouat, Houssam Zenati, Chuan-Sheng Foo, Vijay Chandrasekhar, and Georgios Piliouras. Optimistic mirror descent in saddle-point problems: Going the extra (gradient) mile. In *International Conference on Learning Representations*, 2018.
- [MNG17] Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. The numerics of GANs. *Advances in Neural Information Processing Systems*, 30, 2017.
- [MOP20] Aryan Mokhtari, Asuman Ozdaglar, and Sarath Pattathil. A unified analysis of extra-gradient and optimistic gradient methods for saddle point problems: Proximal point approach. In *International Conference on Artificial Intelligence and Statistics*, pages 1497–1507. PMLR, 2020.
- [MRGK22] Dmitriy Metev, Alexander Rogozin, Alexander Gasnikov, and Dmitry Kovalev. Decentralized saddle-point problems with different constants of strong convexity and strong concavity. *arXiv preprint arXiv:2206.00090*, 2022.
- [Nem04] Arkadi Nemirovski. Prox-method with rate of convergence  $\mathcal{O}(1/t)$  for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1):229–251, 2004.
- [Nes04] Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Springer Science & Business Media, 2004.
- [NJLS09] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19:1574–1609, 2009.
- [NK17] Vaishnavh Nagarajan and J Zico Kolter. Gradient descent GAN optimization is locally stable. *Advances in Neural Information Processing Systems*, 30, 2017.
- [NS11] Yurii Nesterov and Laura Scramali. Solving strongly monotone variational and quasi-variational inequalities. *Discrete & Continuous Dynamical Systems*, 31(4):1383, 2011.
- [OC15] Brendan O’donoghue and Emmanuel Candes. Adaptive restart for accelerated gradient schemes. *Foundations of computational mathematics*, 15(3):715–732, 2015.
- [OX21] Yuyuan Ouyang and Yangyang Xu. Lower complexity bounds of first-order methods for convex-concave bilinear saddle-point problems. *Mathematical Programming*, 185(1-2):1–35, 2021.
- [Rd17] Vincent Roulet and Alexandre d’Aspremont. Sharpness, restart and acceleration. *Advances in Neural Information Processing Systems*, 30, 2017.
- [RG22] James Renegar and Benjamin Grimmer. A simple nearly optimal restart scheme for speeding up first-order methods. *Foundations of Computational Mathematics*, 22(1):211–256, 2022.

- [RLLY21] Hassan Rafique, Mingrui Liu, Qihang Lin, and Tianbao Yang. Weakly-convex-concave min-max optimization: provable algorithms and applications in machine learning. *Optimization Methods and Software*, pages 1–35, 2021.
- [RYY19] Ernest K Ryu, Kun Yuan, and Wotao Yin. Ode analysis of stochastic gradient methods with optimism and anchoring for minimax problems. *arXiv preprint arXiv:1905.10899*, 2019.
- [SCP22] Othmane Sebbouh, Marco Cuturi, and Gabriel Peyré. Randomized stochastic gradient descent ascent. In *International Conference on Artificial Intelligence and Statistics*, pages 2941–2969. PMLR, 2022.
- [THO22] Kiran K Thekumparampil, Niao He, and Sewoong Oh. Lifted primal-dual method for bilinearly coupled smooth minimax optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 4281–4308. PMLR, 2022.
- [Tse95] Paul Tseng. On linear convergence of iterative methods for the variational inequality problem. *Journal of Computational and Applied Mathematics*, 60(1-2):237–252, 1995.
- [Tse08] Paul Tseng. On accelerated proximal gradient methods for convex-concave optimization. *submitted to SIAM Journal on Optimization*, 2(3), 2008.
- [WL20] Yuanhao Wang and Jian Li. Improved algorithms for convex-concave minimax optimization. *Advances in Neural Information Processing Systems*, 33:4800–4810, 2020.
- [WX17] Jialei Wang and Lin Xiao. Exploiting strong convexity from data with primal-dual first-order algorithms. In *International Conference on Machine Learning*, pages 3694–3702, 2017.
- [XHZ21] Guangzeng Xie, Yuze Han, and Zhihua Zhang. Dippa: An improved method for bilinear saddle point problems. *arXiv preprint arXiv:2103.08270*, 2021.
- [XYLC19] Lin Xiao, Adams Wei Yu, Qihang Lin, and Weizhu Chen. Dscovr: Randomized primal-dual block coordinate algorithms for asynchronous distributed optimization. *Journal of Machine Learning Research*, 20(1):1634–1691, 2019.
- [YOLH22] Junchi Yang, Antonio Orvieto, Aurelien Lucchi, and Niao He. Faster single-loop algorithms for minimax optimization without strong concavity. In *International Conference on Artificial Intelligence and Statistics*, pages 5485–5517. PMLR, 2022.
- [YXL<sup>+</sup>19] Yan Yan, Yi Xu, Qihang Lin, Lijun Zhang, and Tianbao Yang. Stochastic primal-dual algorithms with faster convergence than  $o(1/\sqrt{T})$  for problems without bilinear structure. *arXiv preprint arXiv:1904.10112*, 2019.
- [YXL<sup>+</sup>20] Yan Yan, Yi Xu, Qihang Lin, Wei Liu, and Tianbao Yang. Optimal epoch stochastic gradient descent ascent methods for min-max optimization. *Advances in Neural Information Processing Systems*, 33:5789–5800, 2020.
- [ZH22] Junyu Zhang, Mingyi Hong, and Shuzhong Zhang. On lower iteration complexity bounds for the convex concave saddle point problems. *Mathematical Programming*, 194(1-2):901–935, 2022.
- [ZX17] Yuchen Zhang and Lin Xiao. Stochastic primal-dual coordinate method for regularized empirical risk minimization. *Journal of Machine Learning Research*, 18:1–42, 2017.



The Appendix is organized as follows. Section A provides specific examples in our minimax optimization setting. Section B compares our work with prior related works. Section C discusses the stochastic AG-EG algorithms in detail. Section D proves the main results. Finally, Section E provides proofs of auxiliary lemmas that support the proofs of main results.

## A Examples

We conduct an overview of some applications in this section.

**Reinforcement learning.** Reinforcement learning problems can be formalized as Markov Decision Processes (MDPs) where, at each step  $t = 1, \dots, n$ , the learner receives a four-element tuple,  $\{s_t, a_t, r_t, s_{t+1}\}$ , where  $(s_t, a_t)$  is the current state-action pair,  $r_t$  is the reward received upon choosing  $a_t$ , and  $s_{t+1}$  is the next state drawn from a transition distribution. For example, policy evaluation with a linear function approximator can be formalized in terms of the minimization of the *mean squared projected Bellman-Error* (MSPBE) [DNP<sup>+</sup>14] based on a set of tuples:

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{A}\boldsymbol{\theta} - \mathbf{b}\|_{\mathbf{C}^{-1}}^2 + \frac{\rho}{2} \|\boldsymbol{\theta}\|^2 \quad (22)$$

where  $\mathbf{A} = \frac{1}{n} \sum_{t=1}^n \boldsymbol{\phi}(s_t)(\boldsymbol{\phi}(s_t) - \gamma\boldsymbol{\phi}(s_{t+1}))^\top$ ,  $\mathbf{b} = \frac{1}{n} \sum_{t=1}^n r_t \boldsymbol{\phi}(s_t)$ , and  $\mathbf{C} = \frac{1}{n} \sum_{t=1}^n \boldsymbol{\phi}(s_t)\boldsymbol{\phi}(s_t)^\top$  for a given feature mapping  $\boldsymbol{\phi}$ . To reduce the computational cost incurred by calculating the inverse of matrix  $\mathbf{C}$ , [DCL<sup>+</sup>17] propose an alternative minimax form of (22):

$$\min_{\boldsymbol{\theta}} \max_{\mathbf{w}} \frac{\rho}{2} \|\boldsymbol{\theta}\|^2 - \mathbf{w}^\top \mathbf{A}\boldsymbol{\theta} - \frac{1}{2} \|\mathbf{w}\|_{\mathbf{C}}^2 + \mathbf{w}^\top \mathbf{b}$$

which falls under the umbrella of problem (3) whenever  $\mathbf{C}$  is positive definite.

**Quadratic games.** Another class of examples arises in the setting of bilinear games, where the minimax objective is:

$$\mathcal{F}(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \mathbf{x}^\top \mathbf{M}_F \mathbf{x} + \mathbf{x}^\top \mathbf{B} \mathbf{y} - \frac{1}{2} \mathbf{y}^\top \mathbf{M}_G \mathbf{y} - \mathbf{x}^\top \mathbf{v}_x + \mathbf{v}_y^\top \mathbf{y} \quad (23)$$

where  $\mathbf{M}_F, \mathbf{M}_G$  are real-valued matrices of dimensions  $n \times n$  and  $m \times m$ . This has the form (3) with  $F(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{M}_F \mathbf{x} - \mathbf{x}^\top \mathbf{v}_x$ ,  $G(\mathbf{y}) = \frac{1}{2} \mathbf{y}^\top \mathbf{M}_G \mathbf{y} - \mathbf{v}_y^\top \mathbf{y}$  and  $H(\mathbf{x}, \mathbf{y}) \equiv \mathbf{x}^\top \mathbf{B} \mathbf{y}$ . A particular case we will be considering in §3.2 is the case of bilinear games; i.e., where there are no quadratic terms. We provide a detailed analysis of the nonasymptotic convergence in this setting in §3.2 and show that the upper bound on the convergence rate given by our algorithm matches the lower bound of [IAGM20, Theorem 3].

**Regularized empirical risk minimization.** The problem of the minimization of the regularized empirical risk for convex losses and linear predictors is a core problem in classical supervised learning:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{A}\mathbf{x}) + F(\mathbf{x}) \equiv \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{a}_i^\top \mathbf{x}) + F(\mathbf{x})$$

where  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]^\top \in \mathbb{R}^{n \times d}$  consists of feature vectors  $\{\mathbf{a}_i\}$ ,  $\mathcal{L}_i(\mathbf{y})$  is a univariate convex loss for the  $i$ th data point, and  $F(\mathbf{x})$  is a convex regularizer. A standard construction turns this empirical risk minimization problem into a saddle-point problem as follows:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathbb{R}^m} F(\mathbf{x}) + \mathbf{x}^\top \mathbf{A} \mathbf{y} - \underbrace{\mathcal{L}^*(\mathbf{y})}_{\text{Legendre dual function of } \mathcal{L}(\mathbf{y})} \equiv F(\mathbf{x}) + \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{y}^\top \mathbf{a}_i - \frac{1}{n} \sum_{i=1}^n \mathcal{L}^*(\mathbf{y}_i)$$

See [ZX17, WX17, XYLC19] for in-depth discussions of solving this problem under such a dual form of representation.

## B Related Work

Here we compare our results with related work on saddle-point (minimax) optimization in the machine learning and optimization literature.

**Bilinear game case, nonstochastic setting.** In the bilinear game case where  $L_F = \mu_F = L_G = \mu_G = 0$ , a lower bound has been established by [IAGM20]:  $\Omega\left(\sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B} \mathbf{B}^\top)}} \log\left(\frac{1}{\varepsilon}\right)\right)$ . The study of bilinear game has been initiated by [DISZ18] for understanding saddle-point optimization. They proposed the optimistic gradient descent-ascent (OGDA) algorithm and achieved sublinear convergence. Subsequently, the classical methods of ExtraGradient (EG) and Optimistic Gradient Descent Ascent (OGDA) algorithms were proven to have linear convergence rate for strongly monotone and Lipschitz operator with  $\mathcal{O}\left(\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B} \mathbf{B}^\top)} \log\left(\frac{1}{\varepsilon}\right)\right)$  iteration complexity [GHP<sup>+</sup>19, MOP20]. [AMLJG20] proved that by considering first-order methods using a fixed number of composed gradient evaluations and only the last iterate (this class of methods is called 1-SCLI and excludes momentum and restarting), the  $\mathcal{O}\left(\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B} \mathbf{B}^\top)} \log\left(\frac{1}{\varepsilon}\right)\right)$  iteration complexity for EG is optimal. In the absence of strong monotonicity assumption, [LBJM<sup>+</sup>20] provided the first set of global non-asymptotic last-iterate convergence guarantees for a stochastic game over a non-compact domain from a Hamiltonian viewpoint. In particular, the proposed stochastic Hamiltonian gradient method attains convergence in the finite-sum stochastic bilinear game as well. In a very recent work, when restricted to the bilinear minimax optimization, [KGR22] derived an iteration complexity that is essentially  $\mathcal{O}\left(\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B} \mathbf{B}^\top)} \log\left(\frac{1}{\varepsilon}\right)\right)$ . This is comparable to the rates in [DISZ18, LS19, GHP<sup>+</sup>19, MOP20, MKS<sup>+</sup>20]. To match the  $\Omega\left(\sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B} \mathbf{B}^\top)}} \log\left(\frac{1}{\varepsilon}\right)\right)$  lower bound provided by [IAGM20], [ASM<sup>+</sup>20] considered EG with momentum. They used a perturbed spectral analysis encompassing Polyak momentum. Nonetheless, [ASM<sup>+</sup>20] only provided accelerated rates in the regime where the condition number is large. [LYL<sup>+</sup>22] is the first to show that a variant of stochastic extragradient method converges at an accelerated convergence rate for bilinear games with unbounded domain and unbounded stochastic noise using restarted iterate averaging, and when focusing on the nonstochastic setting, matches the lower bound of [IAGM20].

**Smooth strongly convex-concave case, nonstochastic setting.** Lower bound has been recently studied by [OX21] for smooth convex-concave minimax optimization, and by [ZH22] for

strongly-convex-strongly-concave saddle point problems. The latter is of order  $\Omega\left(\left(\sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}\right) \log\left(\frac{1}{\varepsilon}\right)\right)$ . As for upper bounds, earlier extragradient-based methods [Tse95] and accelerated dual extrapolation algorithm [NS11] achieve, when restricted to the bilinearly coupled problem, an iteration complexity of  $\tilde{\mathcal{O}}\left(\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}\right)$ . The same iteration complexity has also been achieved by [GBV<sup>+</sup>19], [MOP20], [CST21] from a relative Lipschitz viewpoint.<sup>5</sup> Improving upon this result, [LJJ20b] achieved a complexity of  $\tilde{\mathcal{O}}\left(\sqrt{\frac{L_F L_G}{\mu_F \mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}\right)$  using proper acceleration methods, when restricted to the bilinearly coupled problem. [WL20] achieved<sup>6</sup>  $\tilde{\mathcal{O}}\left(\sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt[4]{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G} \cdot \frac{L_F L_G}{\mu_F \mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}\right)$  iteration complexity and a Hermitian-skew-based analysis nearly matches [ZHZ22] for the quadratic minimax game case. For the same problem, [XHZ21] achieved a complexity of  $\tilde{\mathcal{O}}\left(\sqrt[4]{\frac{L_F L_G}{\mu_F \mu_G} \left(\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}\right)} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}\right)$ . These works improve upon [LJJ20b] in a fine-grained fashion where separate Lipschitz constants on different parts of the objective are allowed. In early 2022, three concurrent works [KGR22, THO22, JST22] study the nonstochastic problem and independently match the lower bound by [ZHZ22]. The main novelty of this work is that both lower bounds [IAGM20] and [ZHZ22] are achieved in a single algorithm, plus an optimal statistical error term up to a constant prefactor in the stochastic setting. Recently, an independent work by [LYGJ22] also proposed a single algorithm that can achieve the optimal rates for both settings. However, their algorithm is based on optimistic gradient, and is less general than the variational inequality setting studied in this paper.

**Stochastic setting.** Stochastic minimax optimization has been studied intensively as a special case of the variational inequalities. It is widely assumed in classical literature on stochastic variational inequality [JNT11] that the set of parameters and the variance of the stochastic estimate of the vector field are bounded. [CLO17] extended the analysis of [JNT11] that accelerates the convergence rates for a class of variational inequalities. [IJOT17] proposed an analysis of stochastic extragradient using large batches to reduce the variance. [MLZ<sup>+</sup>18] showed almost sure convergence of SEG to a strictly coherent solution (a.k.a. star-strict monotone variational inequality problem). In a similar vein, [RYY19] showed that SGDA with anchoring almost surely converges to strictly convex-concave saddle points. [FOP20] developed a multistage variant of stochastic gradient descent ascent and stochastic optimistic gradient descent ascent with constant learning rate decay schedule. We improve upon their rates since their iteration complexity depends on a significantly larger condition number than our method and is infinite in the absence of strong convex-concavity. They achieved the optimal dependency on the noise variance but suboptimal dependency on the condition number. [HIMM20] developed a double stepsize extragradient method and proved the last-iterate convergence rates under an error bound condition

<sup>5</sup>[MOP20] report an  $\tilde{\mathcal{O}}\left(\frac{L_F \vee L_G + \sqrt{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}}{\mu_F \wedge \mu_G}\right)$  complexity, but the mentioned complexity can be obtained via a scaling-reduction argument: consider  $\mu_F = \mu_G$  case first, then consider the general case by rescaling the  $y$  variable by a factor of  $\sqrt{\frac{\mu_G}{\mu_F}}$ .

<sup>6</sup>Note the cross term here,  $\tilde{\mathcal{O}}\left(\sqrt[4]{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G} \cdot \frac{L_F L_G}{\mu_F \mu_G}}\right)$ , cannot be absorbed into the summation of the remaining terms.

---

**Algorithm 2** Stochastic Accelerated Gradient-Extra Gradient (AG-EG) Descent-Ascent Algorithm, Direct Approach

---

**Require:** Initialization  $\mathbf{z}_0$ , total number of iterates  $\mathcal{T}$ , stepsizes  $(\alpha_t, \eta_t : t = 1, 2, \dots)$

- 1: Set  $\mathbf{z}_{-\frac{1}{2}}^{\text{ag}} \leftarrow \mathbf{z}_0, \mathbf{z}_0^{\text{md}} \leftarrow \mathbf{z}_0$
  - 2: **for**  $t = 1, 2, \dots, \mathcal{T}$  **do**
  - 3:   Draw samples  $\xi_{t-\frac{1}{2}} \sim \mathcal{D}_\xi$  from oracle, and also  $\zeta_{t-\frac{1}{2}}, \zeta_t \sim \mathcal{D}_\zeta$  independently from oracle
  - 4:    $\mathbf{z}_{t-\frac{1}{2}} \leftarrow \mathbf{z}_{t-1} - \eta_t \left( \nabla \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}}) + \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) - \mu(\mathbf{z}_{t-1}^{\text{md}} - \mathbf{z}_{t-1}) \right)$
  - 5:    $\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \leftarrow (1 - \alpha_t) \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t \mathbf{z}_{t-\frac{1}{2}}$
  - 6:    $\mathbf{z}_t \leftarrow \mathbf{z}_{t-1} - \eta_t \left( \nabla \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) + \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) - \mu(\mathbf{z}_{t-1}^{\text{md}} - \mathbf{z}_{t-\frac{1}{2}}) \right)$
  - 7:    $\mathbf{z}_t^{\text{md}} \leftarrow (1 - \alpha_{t+1}) \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} + \alpha_{t+1} \mathbf{z}_t$
  - 8: **end for**
  - 9: **Output:**  $\mathbf{z}_{\mathcal{T}}$
- 

similar to star-strong monotonicity. [KLL20] proposed a simple and optimal scheme for a class of generalized strongly monotone (stochastic) variational inequalities. Due to the unconstrained nature of stochastic bilinear models, these two assumptions do not hold in this case because the noise increases with the value of the parameters. In a recent work, [MKS<sup>+</sup>20] has shown that stochastic extragradients can be computed under a different stepsize, which removes the bounded domain assumption, while still requiring the bounded noise assumption. The work also discussed the advantages of using the same mini-batch for the two stochastic gradients in stochastic extragradients. In another vein, [JDES<sup>+</sup>20] focused on stochastic extragradient in games with a large number of players. In that case they proposed an extragradient algorithm that randomly updates a small subset of the players at each iteration. [YXL<sup>+</sup>19, YXL<sup>+</sup>20, RLLY21] studied the nonsmooth setting and obtained fast rates. More recent works consider minimax optimization problems without convexity and/or concavity, where the goal is to find first-order and second-order stationary points [SCP22, LJJ20a, YOLH22, GYYY20, CLZ21, LLC22]. One interesting direction is to extend our algorithm to their settings and obtain a fine-grained complexity bound with optimal rates.

## C Algorithms

In this section we provide delayed algorithms for the AG-EG (direct approach) and the AG-EG with bounded domain and proximal operator.

### C.1 Stochastic AG-EG, Direct Approach

The full algorithm for AG-EG, direct approach is shown in Algorithm 2.

### C.2 Stochastic AG-EG, with Restarting and Projection

The full algorithm for AG-EG, with restarting and proximal operations is shown in algorithm 3.

---

**Algorithm 3** Stochastic Accelerated Gradient-Extra Gradient (AG-EG) Descent-Ascent Algorithm, with Scheduled Restarting

---

**Require:** Initialization  $\mathbf{z}_0^{[0]}$ , total number of epochs  $\mathcal{S} \geq 1$ , total number of per-epoch iterates  $(\mathcal{T}_s : s = 1, \dots, \mathcal{S})$ , stepsizes  $(\alpha_t, \eta_t : t = 1, 2, \dots)$ , ratio of strong-convexity params.  $\mathcal{R} = \frac{\mu_G}{\mu_F}$

**for**  $s = 1, 2, \dots, \mathcal{S}$  **do**

    Set  $\mathbf{z}_{-\frac{1}{2}}^{\text{ag}} \leftarrow \mathbf{z}_0^{[s-1]}, \mathbf{z}_0 \leftarrow \mathbf{z}_0^{[s-1]}, \mathbf{z}_0^{\text{md}} \leftarrow \mathbf{z}_0^{[s-1]}$

**for**  $t = 1, 2, \dots, \mathcal{T}_s$  **do**

        Draw samples  $\xi_{t-\frac{1}{2}} \sim \mathcal{D}_\xi$  from oracle, and also  $\zeta_{t-\frac{1}{2}}, \zeta_t \sim \mathcal{D}_\zeta$  independently from oracle

$\mathbf{z}_{t-\frac{1}{2}} \leftarrow \text{prox}_{\mathbf{z}_{t-1}}^{\eta_t J} \left( \eta_t \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}}) + \eta_t \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right)$

$\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \leftarrow (1 - \alpha_t) \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t \mathbf{z}_{t-\frac{1}{2}}$

$\mathbf{z}_t \leftarrow \text{prox}_{\mathbf{z}_{t-1}}^{\eta_t J} \left( \eta_t \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) + \eta_t \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right)$

$\mathbf{z}_t^{\text{md}} \leftarrow (1 - \alpha_{t+1}) \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} + \alpha_{t+1} \mathbf{z}_t$

**end for**

    Set  $\mathbf{z}_0^{[s]} \leftarrow \mathbf{z}_{t_s-\frac{1}{2}}^{\text{ag}}$  {//Warm-start using the output of the previous epoch}

**end for**

**Output:**  $\mathbf{z}_0^{[\mathcal{S}]}$

---

## D Proofs of Main Results

In this section we present the proofs of our main results. §D.1 illustrates the scaling reduction argument used in the instance of bilinearly-coupled saddle-point problem. §D.2 provides auxiliary lemmas. With a slight adjustment of their presentation order §D.3 proves Theorem 2, §D.4 proves Theorem 1, §D.5 proves Corollary 2 and finally §D.6 proves Corollary 4. Throughout the section, we assume that the Bregman divergence  $\mathcal{D}(\cdot, \cdot)$  is  $\mu_{\mathcal{D}}$ -strongly convex.

### D.1 Scaling Reduction Argument

Here we illustrate the scaling reduction argument that reduces our analysis of our AG-EG Algorithm 1 under bilinearly-coupled saddle-point problem to the one with equal strong-convexity parameters of  $F$  and  $G$  using a reparametrized objective function; the same argument applies to Algorithm 3 and we omit the details. The idea is in fact analogous to mirror descent-ascent with respect to a Bregman divergence, and our goal here is to detail this argument for our analysis.

In lieu to (3) we consider

$$\min_{\hat{\mathbf{x}}} \max_{\hat{\mathbf{y}}} \widehat{\mathcal{F}}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = F(\hat{\mathbf{x}}) + \widehat{H}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) - \widehat{G}(\hat{\mathbf{y}})$$

where we have  $\widehat{\mathcal{F}}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \mathcal{F}(\mathbf{x}, \mathbf{y})$  with the symbolic reparametrization  $\hat{\mathbf{x}} = \mathbf{x}$ ,  $\hat{\mathbf{y}} = \sqrt{\frac{\mu_G}{\mu_F}} \mathbf{y}$ ,  $\widehat{H}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = H(\mathbf{x}, \mathbf{y})$ ,  $\widehat{h}(\hat{\mathbf{x}}, \hat{\mathbf{y}}; \zeta) = h(\mathbf{x}, \mathbf{y}; \zeta)$ ,  $\widehat{G}(\hat{\mathbf{y}}) = G(\mathbf{y})$ ,  $\widehat{g}(\hat{\mathbf{y}}; \xi) = g(\mathbf{y}; \xi)$  and also their derivatives

$$\nabla_{\hat{\mathbf{y}}} \widehat{H}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \sqrt{\frac{\mu_F}{\mu_G}} \nabla_{\mathbf{y}} H(\mathbf{x}, \mathbf{y}), \quad \nabla_{\hat{\mathbf{y}}} \widehat{h}(\hat{\mathbf{x}}, \hat{\mathbf{y}}; \zeta) = \sqrt{\frac{\mu_F}{\mu_G}} \nabla_{\mathbf{y}} h(\mathbf{x}, \mathbf{y}; \zeta)$$

and

$$\nabla \widehat{G}(\hat{\mathbf{y}}) = \sqrt{\frac{\mu_F}{\mu_G}} \nabla G(\mathbf{y}), \quad \nabla \widehat{g}(\hat{\mathbf{y}}; \xi) = \sqrt{\frac{\mu_F}{\mu_G}} \nabla g(\mathbf{y}; \xi)$$

It is straightforward to verify  $\widehat{\mathcal{F}}(\widehat{\mathbf{x}}, \widehat{\mathbf{y}})$  is arguably  $\mu$ -strongly-convex- $\mu$ -strongly-concave. The essence of our update rules is captured by 8 lines corresponding to Lines 5–8 in Algorithm 1, which becomes:

$$\widehat{\mathbf{x}}_{t-\frac{1}{2}} = \widehat{\mathbf{x}}_{t-1} - \eta_t \left( \nabla f(\widehat{\mathbf{x}}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \nabla_{\widehat{\mathbf{x}}} h(\widehat{\mathbf{x}}_{t-1}, \widehat{\mathbf{y}}_{t-1}; \zeta_{t-\frac{1}{2}}) \right) \quad (24a)$$

$$\widehat{\mathbf{y}}_{t-\frac{1}{2}} = \widehat{\mathbf{y}}_{t-1} - \eta_t \left( -\nabla_{\widehat{\mathbf{y}}} h(\widehat{\mathbf{x}}_{t-1}, \widehat{\mathbf{y}}_{t-1}; \zeta_{t-\frac{1}{2}}) + \nabla g(\widehat{\mathbf{y}}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right) \quad (24b)$$

$$\widehat{\mathbf{x}}_{t-\frac{1}{2}}^{\text{ag}} = (1 - \alpha_t) \widehat{\mathbf{x}}_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t \widehat{\mathbf{x}}_{t-\frac{1}{2}} \quad (24c)$$

$$\widehat{\mathbf{y}}_{t-\frac{1}{2}}^{\text{ag}} = (1 - \alpha_t) \widehat{\mathbf{y}}_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t \widehat{\mathbf{y}}_{t-\frac{1}{2}} \quad (24d)$$

$$\widehat{\mathbf{x}}_t = \widehat{\mathbf{x}}_{t-1} - \eta_t \left( \nabla f(\widehat{\mathbf{x}}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \nabla_{\widehat{\mathbf{x}}} h(\widehat{\mathbf{x}}_{t-\frac{1}{2}}, \widehat{\mathbf{y}}_{t-\frac{1}{2}}; \zeta_t) \right) \quad (24e)$$

$$\widehat{\mathbf{y}}_t = \widehat{\mathbf{y}}_{t-1} - \eta_t \left( -\nabla_{\widehat{\mathbf{y}}} h(\widehat{\mathbf{x}}_{t-\frac{1}{2}}, \widehat{\mathbf{y}}_{t-\frac{1}{2}}; \zeta_t) + \nabla g(\widehat{\mathbf{y}}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right) \quad (24f)$$

$$\widehat{\mathbf{x}}_t^{\text{md}} = (1 - \alpha_{t+1}) \widehat{\mathbf{x}}_{t-\frac{1}{2}}^{\text{ag}} + \alpha_{t+1} \widehat{\mathbf{x}}_t \quad (24g)$$

$$\widehat{\mathbf{y}}_t^{\text{md}} = (1 - \alpha_{t+1}) \widehat{\mathbf{y}}_{t-\frac{1}{2}}^{\text{ag}} + \alpha_{t+1} \widehat{\mathbf{y}}_t \quad (24h)$$

The rest translations are also straightforward, represented by

$$\begin{aligned} \widehat{\mathbf{x}}_{t-\frac{1}{2}} &= \widehat{\mathbf{x}}_{t-1} - \eta_t \left( \nabla f(\widehat{\mathbf{x}}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \nabla_{\widehat{\mathbf{x}}} h(\widehat{\mathbf{x}}_{t-1}, \widehat{\mathbf{y}}_{t-1}; \zeta_{t-\frac{1}{2}}) \right) \\ \Leftrightarrow \mathbf{x}_{t-\frac{1}{2}} &= \mathbf{x}_{t-1} - \eta_t \left( \nabla f(\mathbf{x}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \nabla_{\mathbf{x}} h(\mathbf{x}_{t-1}, \mathbf{y}_{t-1}; \zeta_{t-\frac{1}{2}}) \right) \end{aligned}$$

as well as

$$\begin{aligned} \widehat{\mathbf{y}}_t &= \widehat{\mathbf{y}}_{t-1} - \eta_t \left( -\nabla_{\widehat{\mathbf{y}}} h(\widehat{\mathbf{x}}_{t-\frac{1}{2}}, \widehat{\mathbf{y}}_{t-\frac{1}{2}}; \zeta_t) + \nabla g(\widehat{\mathbf{y}}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right) \\ \Leftrightarrow \mathbf{y}_t &= \mathbf{y}_{t-1} - \eta_t \cdot \frac{\mu_F}{\mu_G} \left( -\nabla_{\mathbf{y}} h(\mathbf{x}_{t-\frac{1}{2}}, \mathbf{y}_{t-\frac{1}{2}}; \zeta_t) + \nabla g(\mathbf{y}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) \right) \end{aligned}$$

It is also straightforward to justify that Assumptions 1 and 2 are rediscovered by reverting the scaling reduction from  $\widehat{\mathcal{F}}(\widehat{\mathbf{x}}, \widehat{\mathbf{y}})$  to  $\mathcal{F}(\mathbf{x}, \mathbf{y})$ . Therefore, it suffices to analyze Algorithm 1 for  $\widehat{\mathcal{F}}(\widehat{\mathbf{x}}, \widehat{\mathbf{y}})$  and due to this scaling reduction, we only need to prove all results for the case of  $\frac{\mu_F}{\mu_G} = 1$ . To keep the notations simple, till the rest of this work we slightly abuse the notations and remove the hats in all symbols.

## D.2 Auxiliary Lemmas

We first state the following basic lemma to handle the inner-product induced terms for extragradient analysis:

**Lemma 1.** . Given  $\boldsymbol{\theta}, \boldsymbol{\varphi}_1, \boldsymbol{\varphi}_2 \in \mathcal{Z}$ , a simple and convex function  $J(\cdot)$ , and also  $\boldsymbol{\delta}_1, \boldsymbol{\delta}_2$  that satisfies

$$\boldsymbol{\varphi}_1 = \text{prox}_{\boldsymbol{\theta}}^J(\boldsymbol{\delta}_1), \quad \boldsymbol{\varphi}_2 = \text{prox}_{\boldsymbol{\theta}}^J(\boldsymbol{\delta}_2) \quad (25)$$

then for any  $\mathbf{z} \in \mathcal{Z}$  we have

$$\langle \boldsymbol{\delta}_2, \boldsymbol{\varphi}_1 - \mathbf{z} \rangle + J(\boldsymbol{\varphi}_1) - J(\mathbf{z}) \leq \frac{1}{2\mu_D} \|\boldsymbol{\delta}_2 - \boldsymbol{\delta}_1\|^2 + \mathcal{D}(\boldsymbol{\theta}, \mathbf{z}) - \mathcal{D}(\boldsymbol{\varphi}_2, \mathbf{z}) - \mathcal{D}(\boldsymbol{\theta}, \boldsymbol{\varphi}_1) \quad (26)$$

Furthermore, when taking  $J = 0$ ,  $\mathcal{Z} = \mathbf{R}^d$  and  $\mathcal{D}(\mathbf{z}, \mathbf{u}) = 1/2 \|\mathbf{z} - \mathbf{u}\|^2$ , Eq. (26) reduces to:

$$\langle \boldsymbol{\delta}_2, \boldsymbol{\varphi}_1 - \mathbf{z} \rangle \leq \frac{1}{2} \|\boldsymbol{\delta}_2 - \boldsymbol{\delta}_1\|^2 + \frac{1}{2} [\|\boldsymbol{\theta} - \mathbf{z}\|^2 - \|\boldsymbol{\varphi}_2 - \mathbf{z}\|^2 - \|\boldsymbol{\theta} - \boldsymbol{\varphi}_1\|^2] \quad (27)$$

Proof of Lemma 1 is provided in §E.1. Lemma 1 is standard and commonly adopted in extragradient-based analysis; see Lemma 2 of [CLO17] for one with similar flavor.

En route to our proofs of Theorems 2 and 1 we first introduce some notations. Let  $\tilde{\mathbf{z}} \in \mathcal{Z}$  and let the *pointwise primal-dual gap* function be

$$V(\mathbf{z} \mid \tilde{\mathbf{z}}) = \mathcal{F}(\mathbf{z}) - \mathcal{F}(\tilde{\mathbf{z}}) + \langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z} - \tilde{\mathbf{z}} \rangle \quad (28)$$

We prove that this quantity is lower bounded by a positive quadratic:

**Lemma 2.** *For  $L$ -smooth and  $\mu$ -strongly convex  $\mathcal{F}(\mathbf{z})$ , simple and convex  $J$ , and for any  $\mathbf{z} \in \mathcal{Z}$  we have*

$$V(\mathbf{z} \mid \mathbf{z}^*) = \mathcal{F}(\mathbf{z}) - \mathcal{F}(\mathbf{z}^*) + \langle \nabla \mathcal{H}(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle + J(\mathbf{z}) - J(\mathbf{z}^*) \geq \frac{\mu}{2} \|\mathbf{z} - \mathbf{z}^*\|^2 \quad (29)$$

Proof of Lemma 2 is provided in §E.2. Our final auxiliary lemma on the key properties on stepsizes spells as follows:

**Lemma 3.** *Our stepsize choice (13) satisfies (i)  $\eta_t \leq \frac{t}{B}$ ; (ii)  $\left(\frac{t}{\eta_t} : t \geq 1\right)$  is a nonnegative, nondecreasing arithmetic sequence with common difference  $\sqrt{\frac{1+\beta}{r}}M$ ; (iii)  $M\eta_t \leq 1$ , and (iv) the stepsize condition*

$$r - \frac{2L}{t+1}\eta_t - (1+\beta)M^2\eta_t^2 \geq 0 \quad (30)$$

Proof of Lemma 3 is provided in §E.3.

### D.3 Proof of Theorem 2

*Proof of Theorem 2.* We first introduce some notations. Denote the incurred stochastic noise terms as

$$\begin{aligned} \Delta_{\mathcal{F}}^{t-\frac{1}{2}} &\equiv \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) - \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), & \Delta_{\mathcal{H}}^{t-\frac{1}{2}} &\equiv \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}}) - \mathcal{H}(\mathbf{z}_{t-1}), \\ \Delta_{\mathcal{H}}^t &\equiv \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) - \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}) \end{aligned} \quad (31)$$

For our martingale analysis we adopt the filtrations  $\mathcal{F}_t^\xi \equiv \sigma(\xi_s : s = \frac{1}{2}, \frac{3}{2}, \dots, s \leq t)$  and  $\mathcal{F}_t^\zeta \equiv \sigma(\zeta_s : s = \frac{1}{2}, 1, \frac{3}{2}, \dots, s \leq t)$ , and also  $\mathcal{F}_t \equiv \sigma(\mathcal{F}_t^\xi \cup \mathcal{F}_t^\zeta)$  be the  $\sigma$ -algebra generated by the union of  $\mathcal{F}_t^\xi$  and  $\mathcal{F}_t^\zeta$ . We are ready for the proof which proceeds as the following steps:

**Step 1.** Estimating the primal-dual gap function difference sequence. We provide the following Lemma (4), whose proof is in §E.4:

**Lemma 4.** *For arbitrary  $\tilde{\mathbf{z}} \in \mathcal{Z}$  the iterates of Algorithm 1 satisfy for  $t = 1, \dots, \mathcal{T}$ , almost surely*

$$\begin{aligned} &V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}}) - (1 - \alpha_t)V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}}) \\ &\leq \alpha_t \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) + \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{\alpha_t^2 L}{2} \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \end{aligned} \quad (32)$$

Note the proof only relies on the interpolation updates in our algorithm as in Lines 6 and 8, and hence this result holds in a per-trajectory (almost-sure) fashion.

**Step 2.** We target to prove the following lemma, the complete proof is in §E.6

**Lemma 5.** *For our choice of  $\eta_t$  that satisfies, for a given  $r \in (0, 1)$ , (30) of Lemma 3(iv) that  $r - \frac{2L}{t+1}\eta_t - (1 + \beta)M^2\eta_t^2 \geq 0$ , we have for any  $\tilde{\mathbf{z}} \in \mathbb{R}^n$ ,  $\tilde{\mathbf{y}} \in \mathbb{R}^m$  and  $t = 1, \dots, \mathcal{T}$  that*

$$\begin{aligned} & t(t+1)\mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - (t-1)t\mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\ & \leq \frac{t}{\eta_t}\mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\| - \|\mathbf{z}_t - \tilde{\mathbf{z}}\|] + \left( \frac{1}{1-r}\sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta})\sigma_{\text{Bil}}^2 \right) t\eta_t \end{aligned} \quad (33)$$

Now for a given  $1 \leq \mathcal{T} \leq \mathcal{T}$ , we finish the proof by telescope the above recursion for  $t = 1, \dots, \mathcal{T}$ . We conclude from our choice of stepsize as in (13) that satisfies (30) so by denoting  $\sigma \equiv \frac{1}{\sqrt{3}}\sqrt{\frac{1}{1-r}\sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta})\sigma_{\text{Bil}}^2}$ , we have by Lemma 3(i)

$$\begin{aligned} & \left( \frac{1}{1-r}\sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta})\sigma_{\text{Bil}}^2 \right) \sum_{t=1}^{\mathcal{T}} t\eta_t = 3\sigma^2 \sum_{t=1}^{\mathcal{T}} t\eta_t \leq 3\sigma^2 \cdot \frac{1}{\mathcal{B}} \sum_{t=1}^{\mathcal{T}} t^2 \\ & = 3\sigma^2 \cdot \frac{C\mathbb{E}[\|\mathbf{z}_0 - \tilde{\mathbf{z}}\|]}{\sigma[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}} \cdot \frac{\mathcal{T}(\mathcal{T} + \frac{1}{2})(\mathcal{T}+1)}{3} = \frac{\mathcal{T}(\mathcal{T} + \frac{1}{2})(\mathcal{T}+1)}{[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}} \cdot \sigma C\mathbb{E}[\|\mathbf{z}_0 - \tilde{\mathbf{z}}\|] \end{aligned}$$

where we recall in Lemma 3 that  $\mathcal{B} \equiv \frac{\sigma[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}}{C\sqrt{\mathbb{E}[\|\mathbf{z}_0 - \mathbf{z}^*\|^2]}}$ . Finally by summing over  $t = 1, \dots, \mathcal{T}$ , we have

$$\begin{aligned} & \mathcal{T}(\mathcal{T}+1)\mathbb{E}[V(\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\ & \leq \sum_{t=1}^{\mathcal{T}} \frac{t}{\eta_t}\mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\| - \|\mathbf{z}_t - \tilde{\mathbf{z}}\|] + \left( \frac{1}{1-r}\sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta})\sigma_{\text{Bil}}^2 \right) \sum_{t=1}^{\mathcal{T}} t\eta_t \\ & = \frac{1}{\eta_1}\mathbb{E}\|\mathbf{z}_0 - \tilde{\mathbf{z}}\|^2 + \sum_{t=2}^{\mathcal{T}} \left( \frac{t}{\eta_t} - \frac{t-1}{\eta_{t-1}} \right) \mathbb{E}\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2 - \frac{\mathcal{T}}{\eta_{\mathcal{T}}}\mathbb{E}\|\mathbf{z}_{\mathcal{T}} - \tilde{\mathbf{z}}\|^2 \\ & \quad + \frac{\mathcal{T}(\mathcal{T} + \frac{1}{2})(\mathcal{T}+1)}{[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}} \cdot C\sigma\sqrt{\mathbb{E}\|\mathbf{z}_0 - \tilde{\mathbf{z}}\|^2} \end{aligned}$$

Following the above derivations and apply Lemma 3(ii) we obtain  $\frac{t}{\eta_t} - \frac{t-1}{\eta_{t-1}} = \sqrt{\frac{1+\beta}{r}}M$ . Rearranging the terms along with Jensen's inequality, and noting that

$$\frac{\mathcal{T}(\mathcal{T} + \frac{1}{2})(\mathcal{T}+1)}{[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}} \leq \frac{\mathcal{T}(\mathcal{T} + \frac{1}{2})(\mathcal{T}+1)}{[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}} \leq [\mathcal{T}(\mathcal{T}+1)^2]^{1/2}$$

proves the following inequality (34).

$$\begin{aligned} & \mathcal{T}(\mathcal{T}+1)\mathbb{E}[V(\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] + \frac{\mathcal{T}}{\eta_{\mathcal{T}}}\mathbb{E}\|\mathbf{z}_{\mathcal{T}} - \tilde{\mathbf{z}}\|^2 \\ & \leq \frac{1}{\eta_1}\mathbb{E}\|\mathbf{z}_0 - \tilde{\mathbf{z}}\|^2 + \sqrt{\frac{1+\beta}{r}}M \sum_{t=2}^{\mathcal{T}} \mathbb{E}\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2 + [\mathcal{T}(\mathcal{T}+1)^2]^{1/2} \cdot C\sigma\sqrt{\mathbb{E}\|\mathbf{z}_0 - \tilde{\mathbf{z}}\|^2} \end{aligned} \quad (34)$$



**Step 3. Bounded Iterates** We conduct the following “bootstrapping” argument to arrive at our final theorem. Starting from the recursion (34) we have by setting  $\tilde{z} = z^*$ , Lemma 2 implies that its first summand on the left hand  $\mathcal{T}(\mathcal{T}+1)\mathbb{E}[V(z_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} | z^*)]$  is nonnegative, and hence we can drop it and have for any  $\mathcal{T} = 1, \dots, \mathcal{T}$

$$\begin{aligned}
& \frac{\mathcal{T}}{\eta_{\mathcal{T}}} \mathbb{E} \|z_{\mathcal{T}} - \tilde{z}\|^2 \\
& \leq \frac{1}{\eta_1} \mathbb{E} \|z_0 - \tilde{z}\|^2 + \sqrt{\frac{1+\beta}{r}} M \sum_{t=2}^{\mathcal{T}} \mathbb{E} \|z_{t-1} - \tilde{z}\|^2 + [\mathcal{T}(\mathcal{T}+1)^2]^{1/2} \cdot C\sigma \sqrt{\mathbb{E} \|z_0 - \tilde{z}\|^2} \\
& = \left(\frac{2}{r} L \vee \mathcal{B}\right) \mathbb{E} \|z_0 - z^*\|^2 + \underbrace{\sqrt{\frac{1+\beta}{r}} M \sum_{t=1}^{\mathcal{T}} \mathbb{E} \|z_{t-1} - z^*\|^2}_{\equiv \mathcal{Q}_{\mathcal{T}-1}} + \underbrace{[\mathcal{T}(\mathcal{T}+1)^2]^{1/2} \cdot C\sigma \sqrt{\mathbb{E} \|z_0 - \tilde{z}\|^2}}_{\mathcal{R}_0}
\end{aligned} \tag{35}$$

Converting (35) to a version of partial sum  $\mathcal{Q}_{\mathcal{T}-1} \equiv \sum_{t=1}^{\mathcal{T}} \mathbb{E} \|z_{t-1} - z^*\|^2$  that for all  $\mathcal{T} = 1, \dots, \mathcal{T}$

$$\frac{\mathcal{T}}{\eta_{\mathcal{T}}} \mathbb{E} \|z_{\mathcal{T}} - z^*\|^2 = \frac{\mathcal{T}}{\eta_{\mathcal{T}}} (\mathcal{Q}_{\mathcal{T}} - \mathcal{Q}_{\mathcal{T}-1}) \leq \sqrt{\frac{1+\beta}{r}} M \mathcal{Q}_{\mathcal{T}-1} + \underbrace{\mathcal{R}_0 + \left(\frac{2}{r} L \vee \mathcal{B}\right) \mathcal{Q}_0}_{\mathcal{D}_0} \tag{36}$$

(36) is equivalently written as

$$\frac{\mathcal{T}}{\eta_{\mathcal{T}}} \mathcal{Q}_{\mathcal{T}} \leq \frac{\mathcal{T}+1}{\eta_{\mathcal{T}+1}} \mathcal{Q}_{\mathcal{T}-1} + \mathcal{D}_0$$

From here and onwards, we denote  $\kappa_t \equiv \frac{t}{\eta_t} = \frac{2}{r} L \vee \mathcal{B} + \sqrt{\frac{1+\beta}{r}} M t$  for each  $t = 1, \dots, \mathcal{T}$ . Dividing both sides of the above display by  $\kappa_{\mathcal{T}} \kappa_{\mathcal{T}+1} = \frac{\mathcal{T}}{\eta_{\mathcal{T}}} \cdot \frac{\mathcal{T}+1}{\eta_{\mathcal{T}+1}}$  gives

$$\frac{\mathcal{Q}_{\mathcal{T}}}{\kappa_{\mathcal{T}+1}} \leq \frac{\mathcal{Q}_{\mathcal{T}-1}}{\kappa_{\mathcal{T}}} + \frac{\mathcal{D}_0}{\kappa_{\mathcal{T}} \cdot \kappa_{\mathcal{T}+1}}$$

Telescoping up from  $1, \dots, \mathcal{T}-1$  for  $1 \leq \mathcal{T} \leq \mathcal{T}$  yields

$$\frac{\mathcal{Q}_{\mathcal{T}-1}}{\kappa_{\mathcal{T}}} \leq \frac{\mathcal{Q}_0}{\kappa_1} + \sum_{t=1}^{\mathcal{T}-1} \frac{\mathcal{D}_0}{\kappa_t \cdot \kappa_{t+1}} \leq \frac{\mathcal{Q}_0}{\kappa_1} + \mathcal{D}_0 \sum_{t=1}^{\mathcal{T}-1} \frac{1}{\kappa_t \cdot \kappa_{t+1}}$$

where we applied Lemma 3(ii) that for all  $t = 1, \dots, \mathcal{T}-1$  we have  $\kappa_{t+1} - \kappa_t = \sqrt{\frac{1+\beta}{r}} M$ . This yields

$$\sqrt{\frac{1+\beta}{r}} M \sum_{t=1}^{\mathcal{T}-1} \frac{1}{\kappa_t \cdot \kappa_{t+1}} = \sum_{t=1}^{\mathcal{T}-1} \left[ \frac{1}{\kappa_t} - \frac{1}{\kappa_{t+1}} \right] = \frac{1}{\kappa_1} - \frac{1}{\kappa_{\mathcal{T}}}$$

and hence

$$\sqrt{\frac{1+\beta}{r}} M \frac{\mathcal{Q}_{\mathcal{T}-1}}{\kappa_{\mathcal{T}}} \leq \sqrt{\frac{1+\beta}{r}} M \frac{\mathcal{Q}_0}{\kappa_1} + \mathcal{D}_0 \sqrt{\frac{1+\beta}{r}} M \sum_{t=1}^{\mathcal{T}-1} \frac{1}{\kappa_t \cdot \kappa_{t+1}} = \sqrt{\frac{1+\beta}{r}} M \frac{\mathcal{Q}_0}{\kappa_1} + \mathcal{D}_0 \left( \frac{1}{\kappa_1} - \frac{1}{\kappa_{\mathcal{T}}} \right)$$

Next, we rearrange the above quantity and derive

$$\frac{\sqrt{\frac{1+\beta}{r}}M\mathcal{Q}_0 + \mathcal{D}_0}{\kappa_1} - \frac{\mathcal{D}_0}{\kappa_{\mathcal{T}}} = \frac{\sqrt{\frac{1+\beta}{r}}M\mathcal{Q}_0 + (\mathcal{R}_0 + (\frac{2}{r}L \vee \mathcal{B})\mathcal{Q}_0)}{\kappa_1} - \frac{\mathcal{D}_0}{\kappa_{\mathcal{T}}} = \mathcal{Q}_0 + \frac{\mathcal{R}_0}{\kappa_1} - \frac{\mathcal{D}_0}{\kappa_{\mathcal{T}}}$$

Plugging this into (36) we have for all iterates  $1 \leq \mathcal{T} \leq \mathcal{T}$

$$\begin{aligned} \mathbb{E} \|\mathbf{z}_{\mathcal{T}} - \mathbf{z}^*\|^2 &\leq \sqrt{\frac{1+\beta}{r}}M \frac{\mathcal{Q}_{\mathcal{T}-1}}{\kappa_{\mathcal{T}}} + \frac{\mathcal{D}_0}{\kappa_{\mathcal{T}}} \leq \mathcal{Q}_0 + \frac{\mathcal{R}_0}{\kappa_1} \leq \left(1 + \frac{C\sigma[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}}{\kappa_1\sqrt{\mathcal{Q}_0}}\right) \mathcal{Q}_0 \\ &= \underbrace{(1 + C^2\mathcal{B}\eta_1)}_{\mathcal{A}} \mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2 \end{aligned} \quad (37)$$

where the prefactor  $\mathcal{A}$  lies in  $[1, 1 + C^2]$  and reduces to 1 when the argument is set as 0.

Now we drop the second summand on the left hand of (34) with  $\tilde{\mathbf{z}} = \mathbf{z}^*$ ,  $\mathcal{T} = \mathcal{T}$ . Combining with (37) gives

$$\begin{aligned} &\mathcal{T}(\mathcal{T}+1)\mathbb{E}[V(\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\ &\leq \frac{1}{\eta_1}\mathbb{E} \|\mathbf{z}_0 - \tilde{\mathbf{z}}\|^2 + \sqrt{\frac{1+\beta}{r}}M \sum_{t=2}^{\mathcal{T}} \mathbb{E} \|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2 + [\mathcal{T}(\mathcal{T}+1)]^{1/2} \cdot C\sigma\sqrt{\mathbb{E} \|\mathbf{z}_0 - \tilde{\mathbf{z}}\|^2} \\ &\leq \left(\frac{2}{r}L \vee \mathcal{B} + \sqrt{\frac{1+\beta}{r}}M\right) \mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2 \\ &\quad + \sqrt{\frac{1+\beta}{r}}M(\mathcal{T}-1) \cdot \mathcal{A} \cdot \mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2 + C\sigma[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}\sqrt{\mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2} \\ &\leq \left(\frac{2}{r}L + \mathcal{A}\sqrt{\frac{1+\beta}{r}}M\mathcal{T}\right) \mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2 + \left(\frac{1}{C} + C\right)\sigma[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}\sqrt{\mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2} \end{aligned}$$

Using (29) in Lemma 2 again lower bounds the left hand in the last display as

$$\mathcal{T}(\mathcal{T}+1)\mathbb{E}[V(\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} \mid \mathbf{z}^*)] \geq \frac{\mu}{2}\mathcal{T}(\mathcal{T}+1)\mathbb{E} \left\| \mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} - \mathbf{z}^* \right\|^2 \geq 0$$

Dividing both sides by  $\frac{\mu}{2}\mathcal{T}(\mathcal{T}+1)$  concludes

$$\mathbb{E} \left\| \mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} - \mathbf{z}^* \right\|^2 \leq \frac{2\left(\frac{2}{r}L + \mathcal{A}\sqrt{\frac{1+\beta}{r}}M\mathcal{T}\right)}{\mu\mathcal{T}(\mathcal{T}+1)} \mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2 + \frac{2(\frac{1}{C} + C)\sigma}{\mu\mathcal{T}^{1/2}} \sqrt{\mathbb{E} \|\mathbf{z}_0 - \mathbf{z}^*\|^2}$$

and hence concludes (14) and the whole proof of Theorem 2.  $\square$

#### D.4 Proof of Theorem 1

We overload function notations  $\mathcal{F}, \mathcal{H}$  to the new group accordingly where  $\mathcal{F}(\mathbf{z}) \leftarrow \mathcal{F}(\mathbf{z}) - \frac{\mu_*}{2}\|\mathbf{z} - \mathbf{z}_0\|^2$  is nonstrongly convex and  $\mathcal{H}(\mathbf{z}) \leftarrow \mathcal{H}(\mathbf{z}) + \frac{\mu_*}{2}\|\mathbf{z} - \mathbf{z}_0\|^2$  is an isotropic quadratic, and analogously for their stochastic estimates. For convenience we repeat the iterates of Algorithm 2 as

$$\mathbf{z}_{t-\frac{1}{2}} = \mathbf{z}_{t-1} - \eta_t \left( \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}}) \right)$$

$$\begin{aligned}
\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} &= (1 - \alpha_t) \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t \mathbf{z}_{t-\frac{1}{2}} \\
\mathbf{z}_t &= \mathbf{z}_{t-1} - \eta_t \left( \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) \right) \\
\mathbf{z}_t^{\text{md}} &= (1 - \alpha_{t+1}) \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} + \alpha_{t+1} \mathbf{z}_t
\end{aligned}$$

with the initialization  $\mathbf{z}_0 = \mathbf{z}_0^{\text{md}} = \mathbf{z}_{-\frac{1}{2}}^{\text{ag}} \in \mathbb{R}^{n+m}$ . We continue to assume the noise-related setting as in (31). Our proof proceeds in the following steps:

**Step 1.** We prove the following generalization of Lemma 4:

**Lemma 6.** *For arbitrary  $\tilde{\mathbf{z}} \in \mathbb{R}^{n+m}$  and  $\alpha_t \in (0, 1]$  the iterates of Algorithm 2 satisfy almost surely*

$$\begin{aligned}
&V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}}) - (1 - \alpha_t) V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}}) \\
&\leq \alpha_t \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) + \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{\alpha_t^2 L}{2} \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 - \alpha_t \mu_\star \left\| \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \right\|^2
\end{aligned} \tag{38}$$

The proof goes in an analogous fashion as the proof of Lemma 4, except that the display above (55) is replaced by

$$\langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \leq \langle \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle - \mu_\star \left\| \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \right\|^2$$

due to our  $H$  being a  $\mu_\star$ -strongly-convex- $\mu_\star$ -strongly-concave isotropic quadratic function after scaling reduction. Hence (55) becomes

$$\langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \tilde{\mathbf{z}} \rangle - (1 - \alpha_t) \langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} - \tilde{\mathbf{z}} \rangle \leq \alpha_t \left[ \langle \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle - \mu_\star \left\| \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \right\|^2 \right] \tag{39}$$

This concludes (38) and the whole lemma.

**Step 2.** Analogous to Step 2 in the proof of Theorem 2 in §D.3 we conclude for all  $\mathbf{z} \in \mathbb{R}^n$ ,  $\mathbf{z} \in \mathbb{R}^m$ ,

$$\begin{aligned}
&\eta_t \mathbb{E} \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \\
&\leq \frac{1}{2} \left( \mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2] - \mathbb{E}[\|\mathbf{z}_t - \tilde{\mathbf{z}}\|^2] \right) - \frac{1 - (1 + \beta) M^2 \eta_t^2}{2} \mathbb{E}[\left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2] + \frac{\eta_t^2}{2} \left( 2 + \frac{1}{\beta} \right) \sigma_{\text{Bil}}^2
\end{aligned}$$

To show this, note that

$$\begin{aligned}
&\eta_t \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \\
&\leq \frac{1}{2} \left( \|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2 - \|\mathbf{z}_t - \tilde{\mathbf{z}}\|^2 - \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2 \right) + \frac{\eta_t^2}{2} \left\| \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) - \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}}) \right\|^2
\end{aligned}$$

To handle the stochastic terms, Young's inequality combined with the martingale structure, along with the definition of  $M$ , indicates

$$\begin{aligned}
&\mathbb{E} \left\| \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) - \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}}) \right\|^2 = \mathbb{E} \left\| \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}) - \mathcal{H}(\mathbf{z}_{t-1}) - \boldsymbol{\Delta}_{\mathcal{H}}^{t-\frac{1}{2}} \right\|^2 + \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^t\|^2 \\
&\leq (1 + \beta) M^2 \mathbb{E} \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2 + \left( 1 + \frac{1}{\beta} \right) \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^t\|^2
\end{aligned}$$

Combining the last three displays gives

$$\begin{aligned}
& \eta_t \mathbb{E} \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \\
& \leq \frac{1}{2} \left( \mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2] - \mathbb{E}[\|\mathbf{z}_t - \tilde{\mathbf{z}}\|^2] \right) - \frac{1 - (1 + \beta)M^2\eta_t^2}{2} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2] \\
& \quad + \frac{\eta_t^2}{2} \left( \left(1 + \frac{1}{\beta}\right) \mathbb{E}\|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E}\|\Delta_{\mathcal{H}}^t\|^2 \right)
\end{aligned} \tag{40}$$

Combining this with Lemma 6, we have

$$\begin{aligned}
& \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - (1 - \alpha_t) \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\
& \leq \alpha_t \mathbb{E} \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) + \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{\alpha_t^2 L}{2} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2] - \alpha_t \mu_{\star} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}}\|^2] \\
& = \alpha_t \mathbb{E} \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \\
& \quad - \alpha_t \mathbb{E} \langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}} + \Delta_{\mathcal{H}}^t, \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{\alpha_t^2 L}{2} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2] - \alpha_t \mu_{\star} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}}\|^2] \\
& \leq \frac{\alpha_t}{\eta_t} \left( \frac{1}{2} \left( \mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2] - \mathbb{E}[\|\mathbf{z}_t - \tilde{\mathbf{z}}\|^2] \right) - \frac{1 - (1 + \beta)M^2\eta_t^2}{2} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2] + \frac{\eta_t^2}{2} \left(2 + \frac{1}{\beta}\right) \sigma_{\text{Bil}}^2 \right) \\
& \quad - \alpha_t \mathbb{E} \langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}} + \Delta_{\mathcal{H}}^t, \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{\alpha_t^2 L}{2} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2] - \alpha_t \mu_{\star} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}}\|^2]
\end{aligned}$$

Continuing this estimation gives (note Young's inequality applies, and  $\mathbb{E} \langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}} + \Delta_{\mathcal{H}}^t, \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle = \mathbb{E} \langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}}, \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \rangle$ )

$$\begin{aligned}
& \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - (1 - \alpha_t) \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\
& \leq \frac{\alpha_t}{2\eta_t} \left( \mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2] - \mathbb{E}[\|\mathbf{z}_t - \tilde{\mathbf{z}}\|^2] \right) - \frac{\alpha_t}{2\eta_t} (r - \alpha_t L \eta_t - (1 + \beta)M^2\eta_t^2) \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2] \\
& \quad + \frac{\alpha_t \eta_t}{2} \left(2 + \frac{1}{\beta}\right) \sigma_{\text{Bil}}^2 - \frac{\alpha_t(1-r)}{2\eta_t} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2] \\
& \quad - \alpha_t \mathbb{E} \langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}}, \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \rangle - \alpha_t \mu_{\star} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}}\|^2] \\
& \leq \frac{\alpha_t}{2\eta_t} \left( \mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2] - \mathbb{E}[\|\mathbf{z}_t - \tilde{\mathbf{z}}\|^2] \right) - \frac{\alpha_t}{2\eta_t} (r - \alpha_t L \eta_t - (1 + \beta)M^2\eta_t^2) \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2] \\
& \quad + \frac{\alpha_t \eta_t}{2} \left(2 + \frac{1}{\beta}\right) \sigma_{\text{Bil}}^2 + \frac{\alpha_t \eta_t}{2(1-r)} \mathbb{E}\|\Delta_{\mathcal{F}}^{t-\frac{1}{2}}\|^2 - \alpha_t \mu_{\star} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}}\|^2] \\
& \leq \frac{\alpha_t}{2\eta_t} \left( \mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2] - \mathbb{E}[\|\mathbf{z}_t - \tilde{\mathbf{z}}\|^2] \right) - \frac{\alpha_t}{2\eta_t} (r - \alpha_t L \eta_t - (1 + \beta)M^2\eta_t^2) \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2] \\
& \quad - \alpha_t \mu_{\star} \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}}\|^2] + \frac{\alpha_t \eta_t}{2} \left( \frac{1}{1-r} \sigma_{\text{Str}}^2 + \left(2 + \frac{1}{\beta}\right) \sigma_{\text{Bil}}^2 \right)
\end{aligned}$$

This yields, applying Young's inequality,

$$\begin{aligned}
& \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - (1 - \alpha_t) \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\
& \leq \frac{\alpha_t}{2\eta_t} \left( \mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2] - \mathbb{E}[\|\mathbf{z}_t - \tilde{\mathbf{z}}\|^2] \right) - \frac{\alpha_t}{2\eta_t} (r - \alpha_t L \eta_t - (1 + \beta)M^2\eta_t^2) \mathbb{E}[\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2]
\end{aligned}$$

$$\begin{aligned}
& -\alpha_t \mu_\star \mathbb{E}[\|z_{t-\frac{1}{2}} - \tilde{z}\|^2] + \frac{\alpha_t \eta_t}{2} \left( \frac{1}{1-r} \sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta}) \sigma_{\text{Bil}}^2 \right) \\
& \leq \frac{\alpha_t}{2\eta_t} \left( (1 - \alpha_t) \mathbb{E}[\|z_{t-1} - \tilde{z}\|^2] - \mathbb{E}[\|z_t - \tilde{z}\|^2] \right) - \frac{\alpha_t}{2\eta_t} (r - \alpha_t L \eta_t - (1 + \beta) M^2 \eta_t^2) \mathbb{E}[\|z_{t-\frac{1}{2}} - z_{t-1}\|^2] \\
& \quad + \frac{\alpha_t^2}{2\eta_t} \mathbb{E}[\|z_{t-1} - \tilde{z}\|^2] - \alpha_t \mu_\star \mathbb{E}[\|z_{t-\frac{1}{2}} - \tilde{z}\|^2] + \frac{\alpha_t \eta_t}{2} \left( \frac{1}{1-r} \sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta}) \sigma_{\text{Bil}}^2 \right) \\
& \leq \frac{\alpha_t}{2\eta_t} \left( (1 - \alpha_t) \mathbb{E}[\|z_{t-1} - \tilde{z}\|^2] - \mathbb{E}[\|z_t - \tilde{z}\|^2] \right) - \frac{\alpha_t}{2\eta_t} (r - \alpha_t L \eta_t - (1 + \beta) M^2 \eta_t^2) \mathbb{E}[\|z_{t-\frac{1}{2}} - z_{t-1}\|^2] \\
& \quad + \eta_t \mu_\star^2 \mathbb{E}[\|z_{t-\frac{1}{2}} - z_{t-1}\|^2] + \frac{\alpha_t \eta_t}{2} \left( \frac{1}{1-r} \sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta}) \sigma_{\text{Bil}}^2 \right)
\end{aligned}$$

Setting  $\eta_t = \frac{\alpha_t}{\mu_\star}$  we have

$$\begin{aligned}
& \mathbb{E}[V(z_{t-\frac{1}{2}}^{\text{ag}} | \tilde{z})] + \frac{\mu_\star}{2} \mathbb{E}[\|z_t - \tilde{z}\|^2] - (1 - \alpha_t) \left( \mathbb{E}[V(z_{t-\frac{3}{2}}^{\text{ag}} | \tilde{z})] + \frac{\mu_\star}{2} \mathbb{E}[\|z_{t-1} - \tilde{z}\|^2] \right) \\
& \leq -\frac{\mu_\star}{2} \left( r - 2\alpha_t - \left( \frac{L}{\mu_\star} + \frac{(1+\beta)M^2}{\mu_\star^2} \right) \alpha_t^2 \right) \mathbb{E}[\|z_{t-\frac{1}{2}} - z_{t-1}\|^2] \\
& \quad + \frac{\alpha_t^2}{2\mu_\star} \left( \frac{1}{1-r} \sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta}) \sigma_{\text{Bil}}^2 \right)
\end{aligned}$$

**Step 3.** By the definition  $\alpha_t$  we have  $r - 2\alpha_t - \left( \frac{L}{\mu_\star} + \frac{(1+\beta)M^2}{\mu_\star^2} \right) \alpha_t^2 \geq 0$ , so we obtain regularity condition  $\alpha_t \leq \bar{\alpha} = \frac{r}{1 + \sqrt{1 + r \left( \frac{L}{\mu_\star} + \frac{(1+\beta)M^2}{\mu_\star^2} \right)}}$  of Theorem 1. Since we assumed both  $F$  and  $G$  are nonstrongly convex and  $H$  is a  $\mu_\star$ -strongly-convex- $\mu_\star$ -strongly-concave isotropic quadratic, this implies

$$\mathbb{E}[V(z_{t-\frac{1}{2}}^{\text{ag}} | \tilde{z})] + \frac{\mu_\star}{2} \mathbb{E}[\|z_t - \tilde{z}\|^2] \leq (1 - \alpha_t) \left( \mathbb{E}[V(z_{t-\frac{3}{2}}^{\text{ag}} | \tilde{z})] + \frac{\mu_\star}{2} \mathbb{E}[\|z_{t-1} - \tilde{z}\|^2] \right) + \frac{3\alpha_t^2}{2\mu_\star} \sigma^2$$

Plugging in  $\tilde{z} = z^*$  gives

$$\mathbb{E}[V(\tilde{z} | z^*)] = \mathcal{F}(\tilde{z}) - \mathcal{F}(z^*) + \langle \mathcal{H}(z^*), \tilde{z} - z^* \rangle \geq \langle \nabla \mathcal{F}(z^*) + \mathcal{H}(z^*), \tilde{z} - z^* \rangle = 0$$

and also

$$\mathbb{E}[V(\tilde{z} | z^*)] \leq \langle \nabla \mathcal{F}(z^*) + \mathcal{H}(z^*), \tilde{z} - z^* \rangle + \frac{L}{2} \|\tilde{z} - z^*\|^2 = \frac{L}{2} \|\tilde{z} - z^*\|^2$$

so (by the fact that  $z_{-\frac{1}{2}}^{\text{ag}} = z_0$  and  $z_{-\frac{1}{2}}^{\text{ag}} = z_0$ )

$$\begin{aligned}
& \frac{\mu_\star}{2} \mathbb{E}[\|z_t - z^*\|^2] \leq \mathbb{E}[V(z_{t-\frac{1}{2}}^{\text{ag}} | z^*)] + \frac{\mu_\star}{2} \mathbb{E}[\|z_t - z^*\|^2] \\
& \leq \left( V(z_{-\frac{1}{2}}^{\text{ag}} | z^*) + \frac{\mu_\star}{2} \|z_0 - z^*\|^2 \right) \prod_{\tau=1}^t (1 - \alpha_\tau) + \sum_{\tau=1}^t \frac{3\alpha_\tau^2}{2\mu_\star} \left[ \prod_{\tau'=\tau+1}^t (1 - \alpha_{\tau'}) \right] \sigma^2 \\
& \leq \|z_0 - z^*\|^2 \frac{L + \mu_\star}{2} \prod_{\tau=1}^t (1 - \alpha_\tau) + \frac{3\sigma^2}{2\mu_\star} \sum_{\tau=1}^t \alpha_\tau^2 \prod_{\tau'=\tau+1}^t (1 - \alpha_{\tau'})
\end{aligned}$$

Dividing both sides by  $\frac{\mu_\star}{2}$  gives (11) and our theorem.

## D.5 Proof of Corollary 2

The proof of Corollary 2 mostly follows the proof of Theorem 2 and Corollary 1, except that we modify some steps to adapt to the proximal operator. The proof is as follows:

**Step 1.** Estimating the primal-dual gap function difference sequence. We have the following Lemma (7), whose proof is in §E.5:

**Lemma 7.** *For arbitrary  $\tilde{\mathbf{z}} \in \mathcal{Z}$  the iterates of Algorithm 1 satisfy for  $t = 1, \dots, \mathcal{T}$ , almost surely*

$$\begin{aligned} & V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}}) - (1 - \alpha_t)V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}}) \\ & \leq \alpha_t \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) + \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{\alpha_t^2 L}{2} \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 + \alpha_t \left( J(\mathbf{z}_{t-\frac{1}{2}}) - J(\tilde{\mathbf{z}}) \right) \end{aligned} \quad (41)$$

Note the proof only relies on the interpolation updates in our algorithm as in Lines 6 and 8, and hence this result holds in a per-trajectory (almost-sure) fashion.

**Step 2.** We target to prove the following lemma, the complete proof is in §E.7:

**Lemma 8.** *For our choice of  $\eta_t$  that satisfies, for a given  $r \in (0, 1)$ , that  $r\mu_{\mathcal{D}} - \frac{2L}{t+1}\eta_t - \frac{(1+\beta)M^2\eta_t^2}{\mu_{\mathcal{D}}} \geq 0$ , we have for any  $\tilde{\mathbf{z}} \in \mathbb{R}^n$ ,  $\tilde{\mathbf{y}} \in \mathbb{R}^m$  and  $t = 1, \dots, \mathcal{T}$  that*

$$\begin{aligned} & t(t+1)\mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - (t-1)t\mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\ & \leq \frac{2t}{\eta_t} (\mathbb{E}[\mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}})] - \mathbb{E}[\mathcal{D}(\mathbf{z}_t, \tilde{\mathbf{z}})]) + \left( \frac{1}{1-r}\sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta})\sigma_{\text{Bil}}^2 \right) \frac{t\eta_t}{\mu_{\mathcal{D}}} \end{aligned} \quad (42)$$

We note that Eq. (42) in Lemma 8 only differs with Eq. (33) in Lemma 5 by the use of Bregman distance  $\mathcal{D}$  and a factor of  $1/\mu_{\mathcal{D}}$  on the variance term. Following similar derivations as in the Proof of Theorem 2, we obtain

$$\begin{aligned} & \mathcal{T}(\mathcal{T}+1)\mathbb{E}[V(\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] + \frac{\mathcal{T}}{\eta_{\mathcal{T}}}\mathbb{E}\|\mathbf{z}_{\mathcal{T}} - \tilde{\mathbf{z}}\|^2 \\ & \leq \frac{2}{\eta_1}\mathcal{D}(\mathbf{z}_0, \tilde{\mathbf{z}}) + 2\sqrt{\frac{1+\beta}{r}}M \sum_{t=2}^{\mathcal{T}} \mathbb{E}\mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}}) + \frac{[\mathcal{T}(\mathcal{T}+1)^2]^{1/2}}{\mu_{\mathcal{D}}} \cdot C\sigma\sqrt{\mathbb{E}\|\mathbf{z}_0 - \tilde{\mathbf{z}}\|^2} \end{aligned} \quad (43)$$

The rest of the proof follows the same bounded iterates argument and the restarting argument exactly as in the previous proof of Theorem 2 with only a difference in a factor of  $\mu_{\mathcal{D}}$  that does not change the terms within  $\mathcal{O}$ . Similar derivatives gives us a total iteration complexity of

$$\mathcal{O} \left( \left( \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} \right) \log \left( \frac{1}{\varepsilon} \right) + \frac{\sigma^2}{\mu^2 \varepsilon^2} \right)$$

with epoch length  $\mathcal{T}_s \asymp \sqrt{\frac{L}{\mu}} + \frac{M}{\mu} + \frac{\sigma^2}{\mu^2 \Gamma_0^2 e^{1-s}}$ .

## D.6 Proof of Corollary 4

Before the proof we first adopt the scaling reduction argument as in §D.1, to argue that we only need to prove the result for the case of bilinear games centered at zero, i.e.  $F(\mathbf{x}) = 0 = G(\mathbf{y})$  we have  $L = \mu = \mu_F = 0$ . We set the iteration symbol  $\mathbf{z} \equiv \begin{bmatrix} \hat{\mathbf{x}} \\ \hat{\mathbf{y}} \end{bmatrix} = \begin{bmatrix} \mathbf{x} - \mathbf{z}^* \\ \mathbf{y} - \boldsymbol{\omega}_y^* \end{bmatrix}$  and also  $\widehat{\mathcal{F}}(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \hat{\mathbf{x}}^\top \mathbf{B} \hat{\mathbf{y}}$ , with  $\widehat{\mathcal{F}}(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  being equal to  $\mathcal{F}(\mathbf{x}, \mathbf{y})$  defined as in (20) up to an additive constant. Our scaling-reduction argument hence applies.

*Proof of Corollary 4.* From the update rule we have

$$\mathbf{z}_{t-\frac{1}{2}} = \mathbf{z}_{t-1} - \eta \mathbf{H}_J \mathbf{z}_{t-1} + \eta \boldsymbol{\varepsilon}_{t-\frac{1}{2}} \quad (44a)$$

$$\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} = \frac{t-1}{t+1} \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} + \frac{2}{t+1} \mathbf{z}_{t-\frac{1}{2}} \quad (44b)$$

$$\mathbf{z}_t = \mathbf{z}_{t-1} - \eta \mathbf{H}_J \mathbf{z}_{t-\frac{1}{2}} + \eta \boldsymbol{\varepsilon}_t \quad (44c)$$

Note the  $[\mathbf{x}_t^{\text{md}}; \mathbf{y}_t^{\text{md}}]$  sequence becomes irrelevant in this update;  $\mathbf{H}_J \equiv \begin{bmatrix} \mathbf{0} & \mathbf{B} \\ -\mathbf{B}^\top & \mathbf{0} \end{bmatrix}$  is skew-symmetric with  $\mathbf{H}_J^\top = -\mathbf{H}_J$ , so  $\mathbf{H}_J^2 = -\mathbf{H}_J^\top \mathbf{H}_J$  is symmetric and negative semidefinite. We proceed with the proof in steps:

**Step 1.** We target to show the last-iterate bound

$$\mathbb{E} \|\mathbf{z}_t\|^2 \leq \mathbb{E} \|\mathbf{z}_0\|^2 + 2t\eta^2 \sigma_{\text{Bil}}^2 \quad (45)$$

Note (44a) and (44c) together gives

$$\mathbf{z}_t = (\mathbf{I} - \eta \mathbf{H}_J + \eta^2 \mathbf{H}_J^2) \mathbf{z}_{t-1} - \eta^2 \mathbf{H}_J \boldsymbol{\varepsilon}_{t-\frac{1}{2}} + \eta \boldsymbol{\varepsilon}_t \quad (46)$$

Taking squared norm on both sides of (46), we have when  $\eta \leq \frac{1}{\sqrt{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}}$ ,  $\mathbf{z}_t$  does not expand in Euclidean norm in the nonstochastic case, so

$$\begin{aligned} \mathbb{E} \|\mathbf{z}_t\|^2 &= \mathbb{E} \left[ (\mathbf{z}_{t-1})^\top (\mathbf{I} + \eta^2 \mathbf{H}_J^2 + \eta^4 \mathbf{H}_J^4) \mathbf{z}_{t-1} \right] + \mathbb{E} \left\| -\eta^2 \mathbf{H}_J \boldsymbol{\varepsilon}_{t-\frac{1}{2}} + \eta \boldsymbol{\varepsilon}_t \right\|^2 \\ &\leq \mathbb{E} \|\mathbf{z}_{t-1}\|^2 + \mathbb{E} \left\| \eta^2 \mathbf{H}_J \boldsymbol{\varepsilon}_{t-\frac{1}{2}} \right\|^2 + \mathbb{E} \|\eta \boldsymbol{\varepsilon}_t\|^2 \\ &\leq \mathbb{E} \|\mathbf{z}_{t-1}\|^2 + \eta^2 \left( 1 + \eta^2 \lambda_{\max}(\mathbf{B}^\top \mathbf{B}) \right) \sigma_{\text{Bil}}^2 \leq \mathbb{E} \|\mathbf{z}_{t-1}\|^2 + 2\eta^2 \sigma_{\text{Bil}}^2 \end{aligned} \quad (47)$$

Recursively applying the above concludes (45).

**Step 2.** We start from the update rule (44b) which implies  $(t+1)t\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} = t(t-1)\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} + 2t\mathbf{z}_{t-\frac{1}{2}}$  holds for  $t = 1, \dots, \mathcal{T}$ , so

$$(\mathcal{T} + 1)\mathcal{T}\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} = 2 \sum_{t=1}^{\mathcal{T}} t\mathbf{z}_{t-\frac{1}{2}} \quad \Rightarrow \quad \mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} = \frac{2}{(\mathcal{T} + 1)\mathcal{T}} \sum_{t=1}^{\mathcal{T}} t\mathbf{z}_{t-\frac{1}{2}}$$

Using this to analyze our algorithm:

$$t\mathbf{z}_t - (t-1)\mathbf{z}_{t-1} - \mathbf{z}_{t-1} = t(\mathbf{z}_t - \mathbf{z}_{t-1}) = -\eta \mathbf{H}_J \left[ t\mathbf{z}_{t-\frac{1}{2}} \right] + \eta t \boldsymbol{\varepsilon}_t$$

so telescoping gives

$$\mathcal{T}\mathbf{z}_{\mathcal{T}} - \sum_{t=1}^{\mathcal{T}} \mathbf{z}_{t-1} = -\eta \mathbf{H}_J \sum_{t=1}^{\mathcal{T}} t \mathbf{z}_{t-\frac{1}{2}} + \eta \sum_{t=1}^{\mathcal{T}} t \boldsymbol{\varepsilon}_t$$

which yields

$$\mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} = \frac{2}{(\mathcal{T}+1)\mathcal{T}} \sum_{t=1}^{\mathcal{T}} t \mathbf{z}_{t-\frac{1}{2}} = \frac{2}{-\eta(\mathcal{T}+1)\mathcal{T}} \mathbf{H}_J^{-1} \left( \mathcal{T}\mathbf{z}_{\mathcal{T}} - \sum_{t=1}^{\mathcal{T}} \mathbf{z}_{t-1} - \eta \sum_{t=1}^{\mathcal{T}} t \boldsymbol{\varepsilon}_t \right) \quad (48)$$

Obviously the least singular value of the matrix  $\mathbf{H}_J$  can be lower bounded as  $\sigma_{\min}(\mathbf{H}_J) \geq \sqrt{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}$ . We conclude from (48) along with Young's inequality that

$$\begin{aligned} \lambda_{\min}(\mathbf{B}\mathbf{B}^\top) \mathbb{E} \left\| \mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} \right\|^2 &\leq \mathbb{E} \left\| \mathbf{H}_J \mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} \right\|^2 \\ &= (1+\gamma) \frac{4}{\eta^2(\mathcal{T}+1)^2 \mathcal{T}^2} \mathbb{E} \left\| \sum_{t=1}^{\mathcal{T}} (\mathbf{z}_{\mathcal{T}} - \mathbf{z}_{t-1}) \right\|^2 + (1+\frac{1}{\gamma}) \frac{4}{\eta^2(\mathcal{T}+1)^2 \mathcal{T}^2} \mathbb{E} \left\| \eta \sum_{t=1}^{\mathcal{T}} t \boldsymbol{\varepsilon}_t \right\|^2 \\ &\equiv (1+\gamma)\text{I} + (1+\frac{1}{\gamma})\text{II} \end{aligned}$$

where applying the last-iterate bound (45) together with some elementary estimates leads to

$$\begin{aligned} \text{I} &\leq \frac{4}{\eta^2(\mathcal{T}+1)^2 \mathcal{T}^2} \cdot \mathcal{T} \sum_{t=1}^{\mathcal{T}} \left[ 2\mathbb{E} \|\mathbf{z}_{\mathcal{T}}\|^2 + 2\mathbb{E} \|\mathbf{z}_{t-1}\|^2 \right] \\ &\leq \frac{4}{\eta^2(\mathcal{T}+1)^2 \mathcal{T}^2} \cdot \mathcal{T} \sum_{t=1}^{\mathcal{T}} \left[ 4\mathbb{E} \|\mathbf{z}_0\|^2 + 4(\mathcal{T}+t-1)\eta^2 \sigma_{\text{Bil}}^2 \right] \\ &\leq \frac{16\mathbb{E} \|\mathbf{z}_0\|^2 + 24\eta^2 \sigma_{\text{Bil}}^2 \mathcal{T}}{\eta^2(\mathcal{T}+1)^2} \leq \frac{16\lambda_{\max}(\mathbf{B}^\top \mathbf{B}) \mathbb{E} \|\mathbf{z}_0\|^2}{(\mathcal{T}+1)^2} + \frac{24\sigma_{\text{Bil}}^2}{\mathcal{T}+1} \end{aligned}$$

and, using the property of square-integrable martingales,

$$\begin{aligned} \text{II} &\leq \frac{4}{\eta^2(\mathcal{T}+1)^2 \mathcal{T}^2} \mathbb{E} \left\| \eta \sum_{t=1}^{\mathcal{T}} t \boldsymbol{\varepsilon}_t \right\|^2 = \frac{4}{\eta^2(\mathcal{T}+1)^2 \mathcal{T}^2} \cdot \eta^2 \sum_{t=1}^{\mathcal{T}} t^2 \mathbb{E} \|\boldsymbol{\varepsilon}_t\|^2 \\ &\leq \frac{4\sigma_{\text{Bil}}^2}{\eta^2(\mathcal{T}+1)^2 \mathcal{T}^2} \cdot \eta^2 \frac{\mathcal{T}(\mathcal{T}+\frac{1}{2})(\mathcal{T}+1)}{3} \leq \frac{4\sigma_{\text{Bil}}^2}{3\mathcal{T}} \end{aligned}$$

To summarize we have for arbitrary  $\gamma \in (0, \infty)$

$$\lambda_{\min}(\mathbf{B}\mathbf{B}^\top) \mathbb{E} \left\| \mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} \right\|^2 \leq (1+\gamma) \left( \frac{16\lambda_{\max}(\mathbf{B}^\top \mathbf{B}) \mathbb{E} \|\mathbf{z}_0\|^2}{(\mathcal{T}+1)^2} + \frac{24\sigma_{\text{Bil}}^2}{\mathcal{T}+1} \right) + (1+\frac{1}{\gamma}) \frac{4\sigma_{\text{Bil}}^2}{3\mathcal{T}}$$

Optimizing  $\gamma$  gives along with  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for nonnegatives  $a$  and  $b$ :

$$\begin{aligned} \sqrt{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)} \sqrt{\mathbb{E} \left\| \mathbf{z}_{\mathcal{T}-\frac{1}{2}}^{\text{ag}} \right\|^2} &\leq \sqrt{\frac{16\lambda_{\max}(\mathbf{B}^\top \mathbf{B}) \mathbb{E} \|\mathbf{z}_0\|^2}{(\mathcal{T}+1)^2} + \frac{24\sigma_{\text{Bil}}^2}{\mathcal{T}+1}} + \sqrt{\frac{4\sigma_{\text{Bil}}^2}{3\mathcal{T}}} \\ &\leq \sqrt{\frac{16\lambda_{\max}(\mathbf{B}^\top \mathbf{B}) \mathbb{E} \|\mathbf{z}_0\|^2}{(\mathcal{T}+1)^2}} + \sqrt{\frac{24\sigma_{\text{Bil}}^2}{\mathcal{T}+1}} + \sqrt{\frac{4\sigma_{\text{Bil}}^2}{3\mathcal{T}}} \leq \frac{4\sqrt{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}}{\mathcal{T}+1} \sqrt{\mathbb{E} \|\mathbf{z}_0\|^2} + \frac{7\sigma_{\text{Bil}}}{\sqrt{\mathcal{T}}} \end{aligned}$$

Dividing both sides by  $\sqrt{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}$  and taking squares conclude the result.  $\square$



## E Proof of Auxiliary Lemmas

### E.1 Proof of Lemma 1

The analysis in this subsection is partially motivated by Lemma 2 of [CLO17].

*Proof of Lemma 1.* We first introduce the following lemma on the operator prox:

**Lemma 9** (Lemma 2 in [GL12] and Lemma 1 in [CLO17]). *If  $\phi = \text{prox}_{\theta}(\delta)$  for arbitrarily chosen  $\theta, \delta \in \mathbb{R}^d$ , then for  $\forall \mathbf{z} \in \mathcal{Z}$ , we have the following inequality*

$$\langle \delta, \phi - \mathbf{z} \rangle + J(\phi) - J(\mathbf{z}) \leq \mathcal{D}(\theta, \mathbf{z}) - \mathcal{D}(\theta, \phi) - V(\phi, \mathbf{z})$$

By applying Lemma 9 to Eq. (25), we have for any  $\mathbf{z} \in \mathcal{Z}$

$$\langle \delta_1, \varphi_1 - \mathbf{z} \rangle + J(\varphi_1) - J(\mathbf{z}) \leq \mathcal{D}(\theta, \mathbf{z}) - \mathcal{D}(\theta, \varphi_1) - \mathcal{D}(\varphi_1, \mathbf{z}), \quad (49)$$

$$\langle \delta_2, \varphi_2 - \mathbf{z} \rangle + J(\varphi_2) - J(\mathbf{z}) \leq \mathcal{D}(\theta, \mathbf{z}) - \mathcal{D}(\theta, \varphi_2) - \mathcal{D}(\varphi_2, \mathbf{z}) \quad (50)$$

Specifically, letting  $\mathbf{z} = \varphi_2$  in (49) we have

$$\langle \delta_1, \varphi_1 - \varphi_2 \rangle + J(\varphi_1) - J(\varphi_2) = \mathcal{D}(\theta, \varphi_2) - \mathcal{D}(\theta, \varphi_1) - \mathcal{D}(\varphi_1, \varphi_2) \quad (51)$$

Now, combining inequalities (50) and (51) we have

$$\langle \delta_2, \varphi_2 - \mathbf{z} \rangle + \langle \delta_1, \varphi_1 - \varphi_2 \rangle + J(\varphi_1) - J(\mathbf{z}) \leq \mathcal{D}(\theta, \mathbf{z}) - \mathcal{D}(\varphi_2, \mathbf{z}) - \mathcal{D}(\theta, \varphi_1) - \mathcal{D}(\varphi_1, \varphi_2)$$

which in turn gives

$$\begin{aligned} & \langle \delta_2, \varphi_1 - \mathbf{z} \rangle + J(\varphi_1) - J(\mathbf{z}) \\ & \leq \langle \delta_2 - \delta_1, \varphi_1 - \varphi_2 \rangle + \mathcal{D}(\theta, \mathbf{z}) - \mathcal{D}(\varphi_2, \mathbf{z}) - \mathcal{D}(\theta, \varphi_1) - \mathcal{D}(\varphi_1, \varphi_2) \end{aligned}$$

An application of the Young and Cauchy-Schwartz inequalities gives

$$\begin{aligned} & \langle \delta_2, \varphi_1 - \mathbf{z} \rangle + J(\varphi_1) - J(\mathbf{z}) \\ & \leq \|\delta_2 - \delta_1\| \|\varphi_1 - \varphi_2\| + \mathcal{D}(\theta, \mathbf{z}) - \mathcal{D}(\varphi_2, \mathbf{z}) - \mathcal{D}(\theta, \varphi_1) - \mathcal{D}(\varphi_1, \varphi_2) \\ & \leq \frac{1}{2\mu_{\mathcal{D}}} \|\delta_2 - \delta_1\|^2 + \frac{\mu_{\mathcal{D}}}{2} \|\varphi_1 - \varphi_2\|^2 + \mathcal{D}(\theta, \mathbf{z}) - \mathcal{D}(\varphi_2, \mathbf{z}) - \mathcal{D}(\theta, \varphi_1) - \mathcal{D}(\varphi_1, \varphi_2) \quad (52) \\ & \leq \frac{1}{2\mu_{\mathcal{D}}} \|\delta_2 - \delta_1\|^2 + \mathcal{D}(\theta, \mathbf{z}) - \mathcal{D}(\varphi_2, \mathbf{z}) - \mathcal{D}(\theta, \varphi_1) \end{aligned}$$

In the last inequality, we use the fact that

$$\frac{\mu_{\mathcal{D}}}{2} \|\varphi_1 - \varphi_2\|^2 \leq \mathcal{D}(\varphi_1, \varphi_2)$$

This establishes (27) and hence Lemma 1. □

## E.2 Proof of Lemma 2

*Proof of Lemma 2.* Since  $\mathcal{F}(\mathbf{z})$  is  $L$ -smooth and  $\mu$ -strongly convex. For the rest of this proof, we observe that the saddle definition of  $\mathbf{z}^*$  satisfies the first-order stationary condition for problem (3):

$$\nabla \mathcal{F}(\mathbf{z}^*) + \mathcal{H}(\mathbf{z}^*) + J'(\mathbf{z}^*) = 0 \quad (53)$$

Furthermore, we have

$$\begin{aligned} & \mathcal{F}(\mathbf{z}) - \mathcal{F}(\mathbf{z}^*) + \langle \mathcal{H}(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle + J(\mathbf{z}) - J(\mathbf{z}^*) \\ & \geq \langle \nabla \mathcal{F}(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle + \frac{\mu}{2} \|\mathbf{z} - \mathbf{z}^*\|^2 + \langle \mathcal{H}(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle + \langle J'(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle \\ & = \langle \nabla \mathcal{F}(\mathbf{z}^*) + \mathcal{H}(\mathbf{z}^*) + J'(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle + \frac{\mu}{2} \|\mathbf{z} - \mathbf{z}^*\|^2 = \frac{\mu}{2} \|\mathbf{z} - \mathbf{z}^*\|^2 \end{aligned}$$

where in both of the two displays, the inequality holds due to the  $\mu$ -strong convexity of  $\mathcal{F}$ , and the equality holds due to the first-order stationary condition (53). This completes the proof.  $\square$

## E.3 Proof of Lemma 3

*Proof of Lemma 3.* Items (i)–(iii) are straightforward. For the proof of (30) in item (iv), we note that  $\eta_t = \bar{\eta}_t(\sigma; \mathcal{T}, C, r, \beta) \leq \frac{t}{\frac{2}{r}L + \sqrt{\frac{1+\beta}{r}}Mt} \leq \frac{1}{\sqrt{\frac{1+\beta}{r}}M}$  which gives

$$r - \frac{2L}{t+1}\eta_t - (1+\beta)M^2\eta_t^2 \geq \frac{r}{t} \left( t - \left( \frac{2}{r}L + \sqrt{\frac{1+\beta}{r}}Mt \right) \eta_t \right) \geq 0$$

and hence completes the proof.  $\square$

## E.4 Proof of Lemma 4

*Proof of Lemma 4.* From the convexity and  $L$ -smoothness of  $F$  as in Assumption 1, we know that for arbitrary  $\tilde{\mathbf{z}}$ :

$$\begin{aligned} & \mathcal{F}(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}}) - \mathcal{F}(\tilde{\mathbf{z}}) = \mathcal{F}(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}}) - \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) - \left( \mathcal{F}(\tilde{\mathbf{z}}) - \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) \right) \\ & \leq \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \rangle + \frac{L}{2} \left\| \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \right\|^2 - \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \tilde{\mathbf{z}} - \mathbf{z}_{t-1}^{\text{md}} \rangle \end{aligned}$$

Taking  $\tilde{\mathbf{z}} = \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}}$  in the above inequality, we have

$$\begin{aligned} & \mathcal{F}(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}}) - \mathcal{F}(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}}) = \mathcal{F}(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}}) - \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) - \left( \mathcal{F}(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}}) - \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) \right) \\ & \leq \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \rangle + \frac{L}{2} \left\| \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \right\|^2 - \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \rangle \end{aligned}$$

Multiplying the first display by  $\alpha_t$  and the second display by  $(1 - \alpha_t)$  and adding them up, we have

$$\begin{aligned}
& \mathcal{F}(z_{t-\frac{1}{2}}^{\text{ag}}) - (1 - \alpha_t)\mathcal{F}(z_{t-\frac{3}{2}}^{\text{ag}}) - \alpha_t\mathcal{F}(\tilde{z}) \\
& \leq \langle \nabla \mathcal{F}(z_{t-1}^{\text{md}}), z_{t-\frac{1}{2}}^{\text{ag}} - z_{t-1}^{\text{md}} \rangle + \frac{L}{2} \left\| z_{t-\frac{1}{2}}^{\text{ag}} - z_{t-1}^{\text{md}} \right\|^2 - \langle \nabla \mathcal{F}(z_{t-1}^{\text{md}}), (1 - \alpha_t)z_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t\tilde{z} - z_{t-1}^{\text{md}} \rangle \\
& \leq \langle \nabla \mathcal{F}(z_{t-1}^{\text{md}}), \alpha_t(z_{t-\frac{1}{2}} - z_{t-1}) \rangle + \frac{L}{2} \left\| \alpha_t(z_{t-\frac{1}{2}} - z_{t-1}) \right\|^2 - \langle \nabla \mathcal{F}(z_{t-1}^{\text{md}}), \alpha_t(\tilde{z} - z_{t-1}) \rangle \\
& = \alpha_t \langle \nabla \mathcal{F}(z_{t-1}^{\text{md}}), z_{t-\frac{1}{2}} - \tilde{z} \rangle + \frac{\alpha_t^2 L}{2} \left\| z_{t-\frac{1}{2}} - z_{t-1} \right\|^2
\end{aligned} \tag{54}$$

where we applied the fact from our update rules that  $z_{t-\frac{1}{2}}^{\text{ag}} - z_{t-1}^{\text{md}} = \alpha_t(z_{t-\frac{1}{2}} - z_{t-1})$ .

On the other hand, due to Line (6) in Algorithm 1 we have

$$\begin{aligned}
\langle \mathcal{H}(\tilde{z}), z_{t-\frac{1}{2}}^{\text{ag}} - \tilde{z} \rangle - (1 - \alpha_t) \langle \mathcal{H}(\tilde{z}), z_{t-\frac{3}{2}}^{\text{ag}} - \tilde{z} \rangle &= \langle \mathcal{H}(\tilde{z}), z_{t-\frac{1}{2}}^{\text{ag}} - \tilde{z} - (1 - \alpha_t)(z_{t-\frac{3}{2}}^{\text{ag}} - \tilde{z}) \rangle \\
&= \alpha_t \langle \mathcal{H}(\tilde{z}), z_{t-\frac{1}{2}} - \tilde{z} \rangle
\end{aligned}$$

Further, due to our monotonicity assumption on  $\mathcal{H}$  we have

$$\langle \mathcal{H}(\tilde{z}), z_{t-\frac{1}{2}} - \tilde{z} \rangle \leq \langle \mathcal{H}(z_{t-\frac{1}{2}}), z_{t-\frac{1}{2}} - \tilde{z} \rangle$$

Combining the above two displays together yields

$$\langle \mathcal{H}(\tilde{z}), z_{t-\frac{1}{2}}^{\text{ag}} - \tilde{z} \rangle - (1 - \alpha_t) \langle \mathcal{H}(\tilde{z}), z_{t-\frac{3}{2}}^{\text{ag}} - \tilde{z} \rangle \leq \alpha_t \langle \mathcal{H}(z_{t-\frac{1}{2}}), z_{t-\frac{1}{2}} - \tilde{z} \rangle \tag{55}$$

Now, summing up Eqs. (54), (55) and recalling the definition of  $V$  in (28), we conclude that

$$\begin{aligned}
& V(z_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{z}) - (1 - \alpha_t)V(z_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{z}) \\
& = \mathcal{F}(z_{t-\frac{1}{2}}^{\text{ag}}) - (1 - \alpha_t)\mathcal{F}(z_{t-\frac{3}{2}}^{\text{ag}}) - \alpha_t\mathcal{F}(\tilde{z}) + \langle \mathcal{H}(\tilde{z}), z_{t-\frac{1}{2}}^{\text{ag}} - \tilde{z} \rangle - (1 - \alpha_t) \langle \mathcal{H}(\tilde{z}), z_{t-\frac{3}{2}}^{\text{ag}} - \tilde{z} \rangle \\
& \leq \alpha_t \langle \nabla \mathcal{F}(z_{t-1}^{\text{md}}) + \mathcal{H}(z_{t-\frac{1}{2}}), z_{t-\frac{1}{2}} - \tilde{z} \rangle + \frac{\alpha_t^2 L}{2} \left\| z_{t-\frac{1}{2}} - z_{t-1} \right\|^2
\end{aligned}$$

and hence conclude (32) and Lemma 4.  $\square$

## E.5 Proof of Lemma 7

*Proof of Lemma 7.* From the convexity and  $L$ -smoothness of  $F$  as in Assumption 1, we know that for arbitrary  $\tilde{z}$ :

$$\begin{aligned}
\mathcal{F}(z_{t-\frac{1}{2}}^{\text{ag}}) - \mathcal{F}(\tilde{z}) &= \mathcal{F}(z_{t-\frac{1}{2}}^{\text{ag}}) - \mathcal{F}(z_{t-1}^{\text{md}}) - \left( \mathcal{F}(\tilde{z}) - \mathcal{F}(z_{t-1}^{\text{md}}) \right) \\
&\leq \langle \nabla \mathcal{F}(z_{t-1}^{\text{md}}), z_{t-\frac{1}{2}}^{\text{ag}} - z_{t-1}^{\text{md}} \rangle + \frac{L}{2} \left\| z_{t-\frac{1}{2}}^{\text{ag}} - z_{t-1}^{\text{md}} \right\|^2 - \langle \nabla \mathcal{F}(z_{t-1}^{\text{md}}), \tilde{z} - z_{t-1}^{\text{md}} \rangle
\end{aligned}$$

Taking  $\tilde{z} = z_{t-\frac{3}{2}}^{\text{ag}}$  in the above inequality, we have

$$\mathcal{F}(z_{t-\frac{1}{2}}^{\text{ag}}) - \mathcal{F}(z_{t-\frac{3}{2}}^{\text{ag}}) = \mathcal{F}(z_{t-\frac{1}{2}}^{\text{ag}}) - \mathcal{F}(z_{t-1}^{\text{md}}) - \left( \mathcal{F}(z_{t-\frac{3}{2}}^{\text{ag}}) - \mathcal{F}(z_{t-1}^{\text{md}}) \right)$$

$$\leq \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \rangle + \frac{L}{2} \left\| \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \right\|^2 - \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \rangle$$

Multiplying the first display by  $\alpha_t$  and the second display by  $(1 - \alpha_t)$  and adding them up, we have

$$\begin{aligned} & \mathcal{F}(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}}) - (1 - \alpha_t)\mathcal{F}(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}}) - \alpha_t\mathcal{F}(\tilde{\mathbf{z}}) \\ & \leq \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \rangle + \frac{L}{2} \left\| \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} \right\|^2 - \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), (1 - \alpha_t)\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} + \alpha_t\tilde{\mathbf{z}} - \mathbf{z}_{t-1}^{\text{md}} \rangle \\ & \leq \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \alpha_t(\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}) \rangle + \frac{L}{2} \left\| \alpha_t(\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}) \right\|^2 - \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \alpha_t(\tilde{\mathbf{z}} - \mathbf{z}_{t-1}) \rangle \\ & = \alpha_t \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{\alpha_t^2 L}{2} \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \end{aligned} \tag{56}$$

where we applied the fact from our update rules that  $\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \mathbf{z}_{t-1}^{\text{md}} = \alpha_t(\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1})$ .

On the other hand, due to Line (6) in Algorithm 3 we have

$$\begin{aligned} \langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \tilde{\mathbf{z}} \rangle - (1 - \alpha_t) \langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} - \tilde{\mathbf{z}} \rangle &= \langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \tilde{\mathbf{z}} - (1 - \alpha_t)(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} - \tilde{\mathbf{z}}) \rangle \\ &= \alpha_t \langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \end{aligned}$$

Further, due to our monotonicity assumption on  $\mathcal{H}$  we have

$$\langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \leq \langle \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle$$

Combining the above two displays together yields

$$\langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \tilde{\mathbf{z}} \rangle - (1 - \alpha_t) \langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} - \tilde{\mathbf{z}} \rangle \leq \alpha_t \langle \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \tag{57}$$

Moreover, we have

$$\begin{aligned} J(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}}) - J(\tilde{\mathbf{z}}) - (1 - \alpha_t) \left( J(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}}) - J(\tilde{\mathbf{z}}) \right) &= J(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}}) - (1 - \alpha_t)J(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}}) - \alpha_t J(\tilde{\mathbf{z}}) \\ &\leq \alpha_t J(\mathbf{z}_{t-\frac{1}{2}}) - \alpha_t J(\tilde{\mathbf{z}}) \end{aligned} \tag{58}$$

Now, summing up Eqs. (56), (57) and recalling the definition of  $V$ , we conclude that

$$\begin{aligned} & V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}}) - (1 - \alpha_t)V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}}) \\ &= \mathcal{F}(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}}) - (1 - \alpha_t)\mathcal{F}(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}}) - \alpha_t\mathcal{F}(\tilde{\mathbf{z}}) + \langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} - \tilde{\mathbf{z}} \rangle - (1 - \alpha_t) \langle \mathcal{H}(\tilde{\mathbf{z}}), \mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} - \tilde{\mathbf{z}} \rangle \\ &\leq \alpha_t \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) + \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{\alpha_t^2 L}{2} \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \end{aligned}$$

and hence conclude (41) and Lemma 7.  $\square$

## E.6 Proof of Lemma 5

*Proof of Lemma 5.* To bound the inner-product terms in (32), by setting  $\varphi_1 = \mathbf{z}_{t-\frac{1}{2}}$ ,  $\boldsymbol{\theta} = \mathbf{z}_{t-1}$ ,  $\varphi_2 = \mathbf{z}_t$ ,  $\boldsymbol{\delta}_1 = \eta_t \left( \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}}) \right)$ ,  $\boldsymbol{\delta}_2 = \eta_t \left( \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) \right)$  as in Lemma 1 (with  $\mathbf{z} = \tilde{\mathbf{z}}$ ), we have

$$\begin{aligned} & \eta_t \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \\ & \leq \frac{1}{2} \left[ \|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2 - \|\mathbf{x}_t - \tilde{\mathbf{z}}\|^2 - \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \right] + \frac{\eta_t^2}{2} \|\tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) - \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}})\|^2 \end{aligned}$$

where Young's inequality combined with the martingale structure yields (also noting (31))

$$\begin{aligned} & \mathbb{E} \|\tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) - \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}})\|^2 \\ & = \mathbb{E} \|\mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}) - \mathcal{H}(\mathbf{z}_{t-1}) - \boldsymbol{\Delta}_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^t\|^2 \\ & \leq (1 + \beta) M^2 \mathbb{E} \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2 + (1 + \frac{1}{\beta}) \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^t\|^2 \end{aligned}$$

Combining the above two displays with expectation taken gives

$$\begin{aligned} & \mathbb{E} \left[ \eta_t \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \right] \\ & \leq \frac{1}{2} \mathbb{E} \left[ \|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\|^2 - \|\mathbf{x}_t - \tilde{\mathbf{z}}\|^2 - \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \right] \\ & \quad + \frac{\eta_t^2}{2} \left( (1 + \beta) M^2 \mathbb{E} \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2 + (1 + \frac{1}{\beta}) \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^t\|^2 \right) \end{aligned} \tag{59}$$

Further, by definition of the primal-dual gap function and the definition of the noisy terms (31), by taking  $\alpha_t = \frac{2}{t+1}$  in (32) of Lemma 4 and taking expectations on both sides, we have

$$\begin{aligned} & \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} | \tilde{\mathbf{z}})] - \frac{t-1}{t+1} \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} | \tilde{\mathbf{z}})] \\ & \leq \frac{2}{t+1} \mathbb{E} \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) + \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{2L}{(t+1)^2} \mathbb{E} \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \\ & = \frac{2}{t+1} \mathbb{E} \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{2L}{(t+1)^2} \mathbb{E} \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \\ & \quad - \frac{2}{t+1} \mathbb{E} \langle \boldsymbol{\Delta}_{\mathcal{F}}^{t-\frac{1}{2}} + \boldsymbol{\Delta}_{\mathcal{H}}^t, \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \end{aligned}$$

Bringing in (59) into the above derivation, we obtain

$$\begin{aligned} & \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} | \tilde{\mathbf{z}})] - \frac{t-1}{t+1} \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} | \tilde{\mathbf{z}})] \\ & \leq \frac{1}{(t+1)\eta_t} \mathbb{E} [\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\| - \|\mathbf{z}_t - \tilde{\mathbf{z}}\|] \\ & \quad - \frac{1}{(t+1)\eta_t} \left( 1 - \frac{2L}{t+1} \eta_t - (1 + \beta) M^2 \eta_t^2 \right) \mathbb{E} \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \\ & \quad + \frac{\eta_t}{t+1} \left( (1 + \frac{1}{\beta}) \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\boldsymbol{\Delta}_{\mathcal{H}}^t\|^2 \right) - \frac{2}{t+1} \mathbb{E} \langle \boldsymbol{\Delta}_{\mathcal{F}}^{t-\frac{1}{2}} + \boldsymbol{\Delta}_{\mathcal{H}}^t, \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle \end{aligned}$$

Recalling that we use the choice of  $\eta_t$  that satisfies for a given  $r \in (0, 1)$  that  $r - \frac{2L}{t+1}\eta_t - (1 + \beta)M^2\eta_t^2 \geq 0$ . With some manipulations we obtain

$$\begin{aligned}
& \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - \frac{t-1}{t+1}\mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\
& \leq \frac{1}{(t+1)\eta_t}\mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\| - \|\mathbf{z}_t - \tilde{\mathbf{z}}\|] \\
& \quad + \frac{\eta_t}{(t+1)}\left((1 + \frac{1}{\beta})\mathbb{E}\|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E}\|\Delta_{\mathcal{H}}^t\|^2\right) - \frac{(1-r)}{(t+1)\eta_t}\mathbb{E}\left[\left\|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\right\|^2\right] \\
& \quad - \frac{2}{t+1}\mathbb{E}\langle\Delta_{\mathcal{F}}^{t-\frac{1}{2}}, \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\rangle - \underbrace{\frac{2}{t+1}\mathbb{E}\langle\Delta_{\mathcal{F}}^{t-\frac{1}{2}}, \mathbf{z}_{t-1} - \tilde{\mathbf{z}}\rangle - \frac{2}{t+1}\mathbb{E}\langle\Delta_{\mathcal{H}}^t, \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}}\rangle}_{\text{I}} \quad (60)
\end{aligned}$$

Due to the law of iterated expectation applied to martingale difference conditions  $\mathbb{E}[\Delta_{\mathcal{F}}^{t-\frac{1}{2}} \mid \mathcal{F}_{t-1}] = \mathbf{0}$  and  $\mathbb{E}[\Delta_{\mathcal{H}}^t \mid \mathcal{F}_{t-\frac{1}{2}}] = \mathbf{0}$ ,  $i = 1, 2$ , we have

$$\text{I} = \mathbf{0}$$

Moreover, for the rest of the terms in Eq. (60), we note that there is a basic quadratic inequality that  $-\frac{1-r}{\eta_t}\|\mathbf{z}_{t-1} - \mathbf{x}_{t-\frac{1}{2}}\|^2 - 2\langle\Delta_{\mathcal{F}}^{t-\frac{1}{2}}, \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\rangle \leq \frac{\eta_t}{1-r}\|\Delta_{\mathcal{F}}^{t-\frac{1}{2}}\|^2$ . Eq. (60) reduces to

$$\begin{aligned}
& \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - \frac{t-1}{t+1}\mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \leq \frac{1}{(t+1)\eta_t}\mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\| - \|\mathbf{z}_t - \tilde{\mathbf{z}}\|] \\
& \quad + \frac{\eta_t}{t+1}\left((1 + \frac{1}{\beta})\mathbb{E}\|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E}\|\Delta_{\mathcal{H}}^t\|^2\right) + \frac{\eta_t}{(1-r)(t+1)}\mathbb{E}\left\|\Delta_{\mathcal{F}}^{t-\frac{1}{2}}\right\|^2 \quad (61)
\end{aligned}$$

Multiplying both sides of (61) by  $t(t+1)$ , we obtain for all  $t = 1, \dots, \mathcal{T}$

$$\begin{aligned}
& t(t+1)\mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - (t-1)t\mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\
& \leq \frac{t}{\eta_t}\mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\| - \|\mathbf{z}_t - \tilde{\mathbf{z}}\|] + t\eta_t\left(\frac{1}{1-r}\mathbb{E}\|\Delta_{\mathcal{F}}^{t-\frac{1}{2}}\|^2 + (1 + \frac{1}{\beta})\mathbb{E}\|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E}\|\Delta_{\mathcal{H}}^t\|^2\right) \\
& \leq \frac{t}{\eta_t}\mathbb{E}[\|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}\| - \|\mathbf{z}_t - \tilde{\mathbf{z}}\|] + \left(\frac{1}{1-r}\sigma_{\text{Str}}^2 + (2 + \frac{1}{\beta})\sigma_{\text{Bil}}^2\right)t\eta_t
\end{aligned}$$

where in the last line above we applied Assumption 2, so by law of iterated expectations

$$\begin{aligned}
& \mathbb{E}\|\Delta_{\mathcal{F}}^{t-\frac{1}{2}}\|^2 = \mathbb{E}\left[\|\nabla\tilde{F}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) - \nabla F(\mathbf{z}_{t-1}^{\text{md}})\|^2\right] \leq \sigma_{\text{Str}}^2, \\
& \mathbb{E}\|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 = \mathbb{E}\left[\|\tilde{H}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}}) - H(\mathbf{z}_{t-1})\|^2\right] \leq \sigma_{\text{Bil}}^2, \\
& \mathbb{E}\|\Delta_{\mathcal{H}}^t\|^2 = \mathbb{E}\left[\|\tilde{H}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) - H(\mathbf{z}_{t-\frac{1}{2}})\|^2\right] \leq \sigma_{\text{Bil}}^2
\end{aligned} \quad (62)$$

□

## E.7 Proof of Lemma 8

*Proof of Lemma 8.* To bound the inner-product terms in (41), by setting  $\boldsymbol{\varphi}_1 = \mathbf{z}_{t-\frac{1}{2}}$ ,  $\boldsymbol{\theta} = \mathbf{z}_{t-1}$ ,  $\boldsymbol{\varphi}_2 = \mathbf{z}_t$ ,  $\boldsymbol{\delta}_1 = \eta_t\left(\nabla\tilde{F}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}})\right)$ ,  $\boldsymbol{\delta}_2 = \eta_t\left(\nabla\tilde{F}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t)\right)$  and

$J = \eta_t J$  as in Lemma 1 (with  $\mathbf{z} = \tilde{\mathbf{z}}$ ), as in Lemma 1 (with  $\mathbf{z} = \tilde{\mathbf{z}}$ ), we have

$$\begin{aligned} & \eta_t \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \eta_t \left( J(\mathbf{z}_{t-\frac{1}{2}}) - J(\tilde{\mathbf{z}}) \right) \\ & \leq \mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}}) - \mathcal{D}(\mathbf{x}_t, \tilde{\mathbf{z}}) - \mathcal{D}(\mathbf{z}_{t-\frac{1}{2}}, \mathbf{z}_{t-1}) + \frac{\eta_t^2}{2\mu_{\mathcal{D}}} \|\tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) - \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}})\|^2 \end{aligned}$$

where Young's inequality combined with the martingale structure yields (also noting (31))

$$\begin{aligned} & \mathbb{E} \|\tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t) - \tilde{\mathcal{H}}(\mathbf{z}_{t-1}; \zeta_{t-\frac{1}{2}})\|^2 \\ & = \mathbb{E} \|\mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}) - \mathcal{H}(\mathbf{z}_{t-1}) - \Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\Delta_{\mathcal{H}}^t\|^2 \\ & \leq (1 + \beta) M^2 \mathbb{E} \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2 + (1 + \frac{1}{\beta}) \mathbb{E} \|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\Delta_{\mathcal{H}}^t\|^2 \end{aligned}$$

Combining the above two displays with expectation taken gives

$$\begin{aligned} & \mathbb{E} \left[ \eta_t \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \eta_t \left( J(\mathbf{z}_{t-\frac{1}{2}}) - J(\tilde{\mathbf{z}}) \right) \right] \\ & \leq \mathbb{E} \left[ \mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}}) - \mathcal{D}(\mathbf{x}_t, \tilde{\mathbf{z}}) - \mathcal{D}(\mathbf{z}_{t-\frac{1}{2}}, \mathbf{z}_{t-1}) \right] \\ & \quad + \frac{\eta_t^2}{2\mu_{\mathcal{D}}} \left( (1 + \beta) M^2 \mathbb{E} \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2 + (1 + \frac{1}{\beta}) \mathbb{E} \|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\Delta_{\mathcal{H}}^t\|^2 \right) \end{aligned} \tag{63}$$

Applying the inequality  $\mathcal{D}(\mathbf{z}_{t-\frac{1}{2}}, \mathbf{z}_{t-1}) \geq \frac{\mu_{\mathcal{D}}}{2} \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2$  again gives

$$\begin{aligned} & \mathbb{E} \left[ \eta_t \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \eta_t \left( J(\mathbf{z}_{t-\frac{1}{2}}) - J(\tilde{\mathbf{z}}) \right) \right] \\ & \leq \mathbb{E} [\mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}}) - \mathcal{D}(\mathbf{x}_t, \tilde{\mathbf{z}})] \\ & \quad - \left( \frac{\mu_{\mathcal{D}}}{2} - \frac{\eta_t^2}{2\mu_{\mathcal{D}}} (1 + \beta) M^2 \right) \mathbb{E} \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2 + \frac{\eta_t^2}{2\mu_{\mathcal{D}}} \left( (1 + \frac{1}{\beta}) \mathbb{E} \|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\Delta_{\mathcal{H}}^t\|^2 \right) \end{aligned}$$

Further, by definition of the primal-dual gap function and the definition of the noisy terms (31), by taking  $\alpha_t = \frac{2}{t+1}$  in (41) of Lemma 7 and taking expectations on both sides, we have

$$\begin{aligned} & \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - \frac{t-1}{t+1} \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\ & \leq \frac{2}{t+1} \mathbb{E} \langle \nabla \mathcal{F}(\mathbf{z}_{t-1}^{\text{md}}) + \mathcal{H}(\mathbf{z}_{t-\frac{1}{2}}), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{2L}{(t+1)^2} \mathbb{E} \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2 \\ & \quad + \frac{2}{t+1} \mathbb{E} \left( J(\mathbf{z}_{t-\frac{1}{2}}) - J(\tilde{\mathbf{z}}) \right) \\ & = \frac{2}{t+1} \mathbb{E} \langle \nabla \tilde{\mathcal{F}}(\mathbf{z}_{t-1}^{\text{md}}; \xi_{t-\frac{1}{2}}) + \tilde{\mathcal{H}}(\mathbf{z}_{t-\frac{1}{2}}; \zeta_t), \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{2L}{(t+1)^2} \mathbb{E} \|\mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1}\|^2 \\ & \quad - \frac{2}{t+1} \mathbb{E} \langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}} + \Delta_{\mathcal{H}}^t, \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle + \frac{2}{t+1} \mathbb{E} \left( J(\mathbf{z}_{t-\frac{1}{2}}) - J(\tilde{\mathbf{z}}) \right) \end{aligned}$$

Bringing in (63) into the above derivation, we obtain

$$\begin{aligned}
& \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - \frac{t-1}{t+1} \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\
& \leq \frac{2}{(t+1)\eta_t} (\mathbb{E}[\mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}})] - \mathbb{E}[\mathcal{D}(\mathbf{z}_t, \tilde{\mathbf{z}})]) \\
& \quad - \frac{1}{(t+1)\eta_t} \left( \mu_{\mathcal{D}} - \frac{2L}{t+1} \eta_t - \frac{(1+\beta)M^2\eta_t^2}{\mu_{\mathcal{D}}} \right) \mathbb{E} \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \\
& \quad + \frac{\eta_t}{(t+1)\mu_{\mathcal{D}}} \left( \left(1 + \frac{1}{\beta}\right) \mathbb{E} \|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\Delta_{\mathcal{H}}^t\|^2 \right) - \frac{2}{t+1} \mathbb{E} \langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}} + \Delta_{\mathcal{H}}^t, \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle
\end{aligned}$$

Recalling that we use the choice of  $\eta_t$  that satisfies for a given  $r \in (0, 1)$  that  $r\mu_{\mathcal{D}} - \frac{2L}{t+1}\eta_t - \frac{(1+\beta)M^2\eta_t^2}{\mu_{\mathcal{D}}} \geq 0$ . With some manipulations we obtain

$$\begin{aligned}
& \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - \frac{t-1}{t+1} \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\
& \leq \frac{2}{(t+1)\eta_t} (\mathbb{E}[\mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}})] - \mathbb{E}[\mathcal{D}(\mathbf{z}_t, \tilde{\mathbf{z}})]) \\
& \quad + \frac{\eta_t}{(t+1)\mu_{\mathcal{D}}} \left( \left(1 + \frac{1}{\beta}\right) \mathbb{E} \|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\Delta_{\mathcal{H}}^t\|^2 \right) - \frac{(1-r)\mu_{\mathcal{D}}}{(t+1)\eta_t} \mathbb{E} \left[ \left\| \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \right\|^2 \right] \\
& \quad - \frac{2}{t+1} \mathbb{E} \langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}}, \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \rangle - \underbrace{\frac{2}{t+1} \mathbb{E} \langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}}, \mathbf{z}_{t-1} - \tilde{\mathbf{z}} \rangle + \frac{2}{t+1} \mathbb{E} \langle \Delta_{\mathcal{H}}^t, \mathbf{z}_{t-\frac{1}{2}} - \tilde{\mathbf{z}} \rangle}_{\text{I}} \quad (64)
\end{aligned}$$

Due to the law of iterated expectation applied to martingale difference conditions  $\mathbb{E}[\Delta_{\mathcal{F}}^{t-\frac{1}{2}} \mid \mathcal{F}_{t-1}] = \mathbf{0}$  and  $\mathbb{E}[\Delta_{\mathcal{H}}^t \mid \mathcal{F}_{t-\frac{1}{2}}] = \mathbf{0}$ ,  $i = 1, 2$ , we have

$$\text{I} = \mathbf{0}$$

Moreover, for the rest of the terms in Eq. (64), we note that there is a basic quadratic inequality that  $-\frac{(1-r)\mu_{\mathcal{D}}}{\eta_t} \|\mathbf{z}_{t-1} - \mathbf{x}_{t-\frac{1}{2}}\|^2 - 2\langle \Delta_{\mathcal{F}}^{t-\frac{1}{2}}, \mathbf{z}_{t-\frac{1}{2}} - \mathbf{z}_{t-1} \rangle \leq \frac{\eta_t}{(1-r)\mu_{\mathcal{D}}} \|\Delta_{\mathcal{F}}^{t-\frac{1}{2}}\|^2$ . Eq. (64) reduces to

$$\begin{aligned}
& \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - \frac{t-1}{t+1} \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \leq \frac{2}{(t+1)\eta_t} (\mathbb{E}[\mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}})] - \mathbb{E}[\mathcal{D}(\mathbf{z}_t, \tilde{\mathbf{z}})]) \\
& \quad + \frac{\eta_t}{(t+1)\mu_{\mathcal{D}}} \left( \left(1 + \frac{1}{\beta}\right) \mathbb{E} \|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\Delta_{\mathcal{H}}^t\|^2 \right) + \frac{\eta_t}{(1-r)(t+1)\mu_{\mathcal{D}}} \mathbb{E} \left\| \Delta_{\mathcal{F}}^{t-\frac{1}{2}} \right\|^2 \quad (65)
\end{aligned}$$

Multiplying both sides of (65) by  $t(t+1)$ , we obtain for all  $t = 1, \dots, \mathcal{T}$

$$\begin{aligned}
& t(t+1) \mathbb{E}[V(\mathbf{z}_{t-\frac{1}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] - (t-1)t \mathbb{E}[V(\mathbf{z}_{t-\frac{3}{2}}^{\text{ag}} \mid \tilde{\mathbf{z}})] \\
& \leq \frac{2t}{\eta_t} (\mathbb{E}[\mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}})] - \mathbb{E}[\mathcal{D}(\mathbf{z}_t, \tilde{\mathbf{z}})]) + \frac{t\eta_t}{\mu_{\mathcal{D}}} \left( \frac{1}{1-r} \mathbb{E} \|\Delta_{\mathcal{F}}^{t-\frac{1}{2}}\|^2 + \left(1 + \frac{1}{\beta}\right) \mathbb{E} \|\Delta_{\mathcal{H}}^{t-\frac{1}{2}}\|^2 + \mathbb{E} \|\Delta_{\mathcal{H}}^t\|^2 \right) \\
& \leq \frac{2t}{\eta_t} (\mathbb{E}[\mathcal{D}(\mathbf{z}_{t-1}, \tilde{\mathbf{z}})] - \mathbb{E}[\mathcal{D}(\mathbf{z}_t, \tilde{\mathbf{z}})]) + \left( \frac{1}{1-r} \sigma_{\text{Str}}^2 + \left(2 + \frac{1}{\beta}\right) \sigma_{\text{Bil}}^2 \right) \frac{t\eta_t}{\mu_{\mathcal{D}}}
\end{aligned}$$

where in the last line above we applied Assumption 2 and law of iterated expectations. This completes the proof of Eq. (42) and hence Lemma 8.  $\square$



Method \ Setting	Bilinearly-coupled SC-SC	Bilinear Game	Stochastic VI
EG / OGDA [MOP20, CST21]	$\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}$	$\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}$	✓
Minimax-APPA [LJJ20b]	$\sqrt{\frac{L_F L_G}{\mu_F \mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}$	—	✗
Proximal Best Response [WL20]	$\sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt[4]{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}} \cdot \frac{L_F L_G}{\mu_F \mu_G} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}$	—	✗
DIPPA [XHZ21]	$\sqrt[4]{\frac{L_F L_G}{\mu_F \mu_G}} \left( \frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G} \right) + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}$	—	✗
LPD [THO22]	$\sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}$	$\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}$	✗
APDG [KGR22]	$\sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}$	$\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}$	✗
PD-EG [JST22]	$\sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}$	—	✗
EG+Momentum [ASM <sup>+</sup> 20]	—	$\sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}}$	✗
SEG with Restarting [LYL <sup>+</sup> 22]	—	$\sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}}$	✓
AG-EG-Direct (this work)	$\sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}$	$\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}$	✓
AG-EG with Restarting (this work)	$\sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}}$	$\sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}}$	✓
Lower Bound [ZHZ22, IAGM20]	$\Omega \left( \left( \sqrt{\frac{L_F}{\mu_F} \vee \frac{L_G}{\mu_G}} + \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\mu_F \mu_G}} \right) \log \left( \frac{1}{\varepsilon} \right) \right)$	$\tilde{\Omega} \left( \sqrt{\frac{\lambda_{\max}(\mathbf{B}^\top \mathbf{B})}{\lambda_{\min}(\mathbf{B}\mathbf{B}^\top)}} \log \left( \frac{1}{\varepsilon} \right) \right)$	—

**Table 1.** Comparison table for solving saddle-point problem (3): first-order oracle complexities for finding an  $\varepsilon$ -optimal point of our proposed AG-EG algorithms with prevailing algorithms under a variety of settings: general bilinearly-coupled SC-SC, bilinear games, as well as a column indicating whether the stochastic variational inequality (VI) case is discussed. The row in the LightCyan background is the convergence result presented in this paper. The “—” indicates that the complexity does not apply to the given case. A polylogarithmic factor in each upper bound in the table is ignored.

Reference	Stochastic variational inequality	No bounded domain assumption	No bounded noise assumption
[Kor76, JNT11, HIMM20]	$\frac{L \vee M}{\varepsilon} \mathcal{D}^2 + \frac{\sigma^2}{\varepsilon^2} \mathcal{D}^2$	✗	✗
[LYL <sup>+</sup> 22] for bilinear games	$\frac{L \vee M}{\varepsilon} \Gamma_0^2 + \frac{\sigma^2}{\varepsilon^2} \Gamma_0^2$	✓	✓
[CLO17, LO21]	$\sqrt{\frac{L}{\varepsilon}} \mathcal{D} + \frac{M}{\varepsilon} \mathcal{D}^2 + \frac{\sigma^2}{\varepsilon^2} \mathcal{D}^2$	✗	✗
This work	$\sqrt{\frac{L}{\varepsilon}} \Gamma_0 + \frac{M}{\varepsilon} \Gamma_0^2 + \frac{\sigma^2}{\varepsilon^2} \Gamma_0^2$	✓	✗
Lower Bound [OX21, Nes04]	$\Omega \left( \sqrt{\frac{L}{\varepsilon}} \mathcal{D} + \frac{M}{\varepsilon} \mathcal{D}^2 + \frac{\sigma^2}{\varepsilon^2} \mathcal{D}^2 \right)$	✗	✗

**Table 2.** Oracle complexities for stochastic VI problem (1)+(2) for finding a point of  $O(\varepsilon)$  primal-dual gap, as well as columns of domain/noise assumptions (note that  $\Gamma_0 \leq \mathcal{D}$ ).