# A LAZY HESSIAN EVALUATION FRAMEWORK FOR AC CELERATING STOCHASTIC BILEVEL OPTIMIZATION

Anonymous authors

004

005

010 011

012

013

014

015

016

017

018

019

020

021

022

023

024

025

026

027

028 029

031 032

033

034

Paper under double-blind review

#### Abstract

Bilevel optimization has recently gained popularity because of its applicability in many machine learning applications. Hypergradient-based algorithms have been widely used for solving bilevel optimization problems because of their strong theoretical and empirical performance in many applications. However, computing these hypergradients requires the evaluation of Hessians (or Hessian-vector products) of the lower-level objective, which presents a major computational bottleneck. To address this challenge, in this paper, we propose LazyBLO (Lazy Hessian Evaluation in Bilevel Optimization), an algorithmic framework that allows infrequent Hessian computation during the execution of the algorithm for solving stochastic bilevel problems. This allows the algorithm to execute faster compared to the stateof-the-art (SOTA) algorithms that evaluate either a single or multiple Hessians in each iteration. We theoretically establish the performance of vanilla SGD-based LazyBLO and show that, despite the additional errors incurred by the infrequent Hessian evaluations, LazyBLO surprisingly matches the computation complexity of the existing SGD-based bilevel algorithms. Extensive experiments further demonstrate that LazyBLO enjoys significant gains in numerical performance compared to the SOTA approaches. To our knowledge, this is the first work to theoretically establish that multiple Hessian computations are not necessary within each iteration to guarantee the convergence of stochastic bilevel algorithms.

#### 1 INTRODUCTION

Bilevel optimization refers to the class of problems with two levels of hierarchy, wherein the solution of the upper-level problem depends on the minimizer of the lower-level problem. Formally, a bilevel problem is stated as:

035 036 037

038

039

$$\min_{\mathbf{x}\in\mathbb{R}^{u}} \left\{ \ell(\mathbf{x}) \triangleq f(\mathbf{x}, \mathbf{y}^{*}(\mathbf{x})) \triangleq \mathbb{E}_{\xi\sim\pi_{f}} \left[ f(\mathbf{x}, \mathbf{y}^{*}(\mathbf{x}); \xi) \right] \right\}$$
s.t.  $\mathbf{y}^{*}(\mathbf{x}) = \arg\min_{\mathbf{y}\in\mathbb{R}^{l}} \left\{ g(\mathbf{x}, \mathbf{y}) \triangleq \mathbb{E}_{\zeta\sim\pi_{g}} \left[ g(\mathbf{x}, \mathbf{y}; \zeta) \right] \right\},$ 
(1)

where  $f(\mathbf{x}, \mathbf{y}) : \mathbb{R}^u \times \mathbb{R}^l \to \mathbb{R}$  and  $g(\mathbf{x}, \mathbf{y}) : \mathbb{R}^u \times \mathbb{R}^l \to \mathbb{R}$  are upper (UL) and lower-level (LL) objectives, respectively. Both the UL and LL objectives are assumed to be smooth while the LL objective is strongly convex with respect to  $\mathbf{y}$ . Moreover,  $\xi \sim \pi_f$  (resp.  $\zeta \sim \pi_g$ ) represents a sample of the UL (resp. LL) objective from distribution  $\pi_f$  (resp.  $\pi_g$ ).

Stochastic bilevel problems in (1) have recently gained prominence as many popular machine 045 learning problems can be modeled in this form. A few typical examples include hyperparameter 046 optimization (Franceschi et al., 2018; Shaban et al., 2019; Bao et al., 2021), meta-learning (Franceschi 047 et al., 2018; Rajeswaran et al., 2019; Ji et al., 2020), adversarial training (Li et al., 2019; Tian et al., 048 2021; Zhang et al., 2022), reinforcement learning (Konda & Tsitsiklis, 1999; Hong et al., 2020), neural architecture search (Liu et al., 2018; Hu et al., 2020; Lian et al., 2019), data hyper-cleaning (Franceschi et al., 2018; Shaban et al., 2019), dictionary learning (Mairal et al., 2011; Lecouat et al., 2020a;b), 051 and more recently, the pretraining-finetuning pipeline (Li et al., 2024; Wu et al., 2024) and data reweighting (Pan et al., 2024) in large language models (LLMs). Consequently, a major research 052 effort has been focused on developing efficient algorithms for solving stochastic bilevel optimization problems.

090

091

Among all existing methods for stochastic bilevel optimization (see Section 2 for detailed discussions), a state-of-the-art (SOTA) approach is the approximate implicit differentiation (AID) method, which relies on directly computing the approximate implicit gradient of the objective  $\ell(\cdot)$  using the implicit function theorem (Ghadimi & Wang, 2018). Because of its ease of implementation, AID is usually the algorithm of choice for many machine learning applications. A typical AID algorithm updates the LL variable using standard stochastic gradient descent (SGD) while the UL variable is updated in each iteration using:  $\mathbf{x}^+ = \mathbf{x} - \alpha h^f$ , where the descent direction  $h^f$  (also often referred to as hypergradient) is an approximation of the implicit gradient, i.e.,

$$h^{f} \approx \nabla \ell(\mathbf{x}) = \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}^{*}(\mathbf{x})) - \nabla_{\mathbf{xy}}^{2} g(\mathbf{x}, \mathbf{y}^{*}(\mathbf{x})) \left[ \nabla_{\mathbf{yy}}^{2} g(\mathbf{x}, \mathbf{y}^{*}(\mathbf{x})) \right]^{-1} \nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}^{*}(\mathbf{x})).$$
(2)

Although AID has been widely adopted for stochastic bilevel optimization in the literature, the computation of the hypergradient  $h^f$  in AID faces two major challenges:

- ① The hypergradient in Eq. (2) requires **multiple** Hessian-vector product (HVP) evaluations for 066 approximating the Hessian inverse in each iteration. This creates a major computational bottleneck 067 for solving the problem in Eq. (1) since the explicit Hessian evaluations are computationally 068 expensive. For example, the Hessian contains one million elements even for a moderately sized 069 problem of dimension d = 1000. What is worse is that inverting such a Hessian typically 070 has a computation complexity of  $\mathcal{O}(d^3)$ , which is time-consuming even for a moderately sized 071 problem. Some modern automatic differentiation tools (e.g., Pearlmutter trick (Pearlmutter, 1994) and Jax (Bradbury et al., 2018)) have been proposed to accelerate the Hessian computation, 073 and HVP computation may not be a major computational bottleneck in some situations where 074 extremely computationally powerful GPUs are available. However, for many resource-constrained 075 and computation-constrained settings (e.g., using small or edge-based devices without GPUs), HVP computation is still a computational bottleneck. For example, each HVP computation 076 could be at least two to six times more expensive than gradient computation using Jax when 077 performed on CPUs, which is still non-trivial, and the cost due to HVP remains not negligible in such systems. Moreover, we note that one Hessian inverse estimation needs multiple HVP 079 computations (Ghadimi & Wang, 2018; Hong et al., 2020). As a result, the total cost of the HVP computation depends on the Hessian-inverse estimation accuracy. This would make the 081 computational cost even higher. 082
- (2) The hypergradient in Eq. (2) depends on the optimal solution of the LL problem  $y^*(x)$ . However, solving the LL problem often requires an iterative method. Thus, solving the LL problem to optimality to obtain an exact value of  $y^*(x)$  may be expensive or even infeasible in practice.

We note that, although Challenge @ has been intensively studied in the literature and addressed to some extent (e.g., the hypergradient is approximated with  $\mathbf{y}^*(\mathbf{x})$  being replaced by  $\mathbf{y}^+ \approx \mathbf{y}^*(\mathbf{x})$ ), Challenge ① remains under-explored. So far, a foundational open problem in the theory of stochastic bilevel optimization naturally arises:

(Q): Can we design algorithms that require fewer Hessian evaluations compared to SOTA, and is it feasible to guarantee any theoretical performance for such algorithms?

092 In this paper, we answer the above question by developing a new algorithmic framework called LazyBLO (Lazy Hessian Evaluation in Bilevel Optimization), which allows infrequent 094 Hessian (Hessian-vector product) evaluations in solving stochastic bilevel problems. Thus, 095 LazyBLO alleviates the computational bottleneck in stochastic bilevel optimization. Specifically, in 096 our LazyBLO approach, a stale version of Hessian is used for multiple iterations while new gradients 097 are computed at each step, thus leading to computational savings. The intuition behind LazyBLO is 098 that, for iterations that are not separated too far from each other, the parameter values usually do not vary significantly. This implies that the Hessians evaluated at these points are highly correlated. Thus, 099 a stale Hessian can still be used to approximate a new one. 100

However, due to the additional errors accumulated because of the use of these stale Hessians, approximate Hessian (HVP) evaluations, and the coupling hierarchical structure of the bilevel problems, it is unclear whether LazyBLO will converge or not. Somewhat surprisingly, we prove that, despite the previously mentioned accumulated errors, LazyBLO not only converges but also achieves the *same* convergence rate as those of the SOTA non-lazy bilevel algorithms. To our knowledge, this is the first work that uses infrequent Hessian computations for computational savings but still can achieve convergence guarantee in solving stochastic bilevel problems.

Our major contributions in this work are summarized as follows:

• We develop a new algorithmic framework LazyBLO that allows the stochastic bilevel algorithms to compute HVPs infrequently. Specifically, the proposed framework updates the HVPs only over a subset of training iterations, while using stale Hessian information in the rest of the iterations.

• We theoretically establish the performance of LazyBLO when the UL and LL updates are performed using vanilla SGD-type updates. We show that the proposed lazy approach, which is supposed to perform worse due to stale Hessian information, can actually *match* the convergence performance of the SOTA bilevel algorithms. Specifically, we show that to achieve an  $\epsilon$ -stationary point, LazyBLO requires  $\mathcal{O}(\epsilon^{-2})$  partial gradient and HVP evaluations. Moreover, thanks to the less frequent Hessian evaluations, the *wall-clock time* of LazyBLO is significantly reduced compared to the SOTA approaches.

• We corroborate our theoretical findings via numerical experiments on data hyper-cleaning and deep hyper-representation tasks with real-world datasets. Our numerical results verify that the infrequent evaluations of HVP lead to considerable computational savings.

120 121 122

108

109

110

111

112

113

114

115 116

117 118

119

#### 2 RELATED WORK

123 124

125

126

127

In this section, we provide a brief overview of several areas of the most related work: ①AID-based bilevel optimization, ② Hessian-free bilevel optimization, and ③ other uses of infrequent Hessian evaluations, thus putting our work into comparative perspective to highlight our novelty. Due to space limitation, we give a summary of other related bilevel optimization methods in Appendix A.

128 129 ble 1, we compare existing AID-based stochas-130 tic bilevel algorithms. BSA (Ghadimi & 131 Wang, 2018) provided the first finite-time con-132 The stochastic bilevel algorithms (e.g., BSA 133 in (Ghadimi & Wang, 2018)), stocBiO in (Ji 134 et al., 2021), AmIGO in (Arbel & Mairal, 135 2022)) that use vanilla-SGD updates require 136  $\mathcal{O}(\epsilon^{-2})$  for both partial gradient evaluations 137 and HVP evaluations to reach an  $\epsilon$ -stationary 138 point. Meanwhile, several works (e.g., SUS-139 TAIN in (Khanduri et al., 2021b), SVRB 140 in (Guo et al., 2021), MRBO and VRBO 141 in (Yang et al., 2021)) utilize momentum-142 based approaches and/or variance reduction approaches to accelerate the convergence 143 of vanilla SGD-based algorithms, achieving 144  $\mathcal{O}(\epsilon^{-1.5})$  for both partial gradient evaluations 145 and HVP evaluations. Although these works 146 guarantee finite-time convergence, the prac-147 tical numerical performance of these bilevel 148 algorithms is slow since they require multi-149 ple Hessian (or HVP) evaluations of the LL 150 objective in each iteration to approximate the 151 Hessian inverse. In this work, we show that 152 the Hessian computations can be skipped and 153 stale Hessian information computed from the previous rounds can be used without hurting 154 the convergence performance while allowing 155 the algorithms to execute much faster. 156

**(D) AID-Based Bilevel Optimization:** In Table 1: Comparison of stochastic bilevel algorithms. BSA (Ghadimi & Wang, 2018) provided the first finite-time convergence guarantees for bilevel optimization. The stochastic bilevel algorithms (e.g., BSA in (Ghadimi & Wang, 2018)), stocBiO in (Ji et al., 2021), AmIGO in (Arbel & Mairal, 2022)) that use vanilla-SGD updates require  $\mathcal{O}\left(\epsilon^{-2}\right)$  for both partial gradient evaluations

|         | # of PG   | # of HVP  | Update   |
|---------|---|---|----------|
| TTSA    | $\mathcal{O}\left(\epsilon^{-2.5}\right)$         | $\mathcal{O}\left(\epsilon^{-2.5} ight)$          | SGD      |
| BSA     | $\mathcal{O}\left(\epsilon^{-2}\right)$           | $\tilde{\mathcal{O}}\left(\epsilon^{-2}\right)$   | SGD      |
| stocBiO | $\mathcal{O}(\epsilon^{-2})$                      | $\tilde{\mathcal{O}}(\epsilon^{-2})$              | SGD      |
| SOBA    | $\mathcal{O}\left(\epsilon^{-2}\right)$           | $\mathcal{O}\left(\epsilon^{-2}\right)$           | SGD      |
| ALSET   | $\mathcal{O}\left(\epsilon^{-2}\right)$           | $\mathcal{O}\left(\epsilon^{-2}\right)$           | SGD      |
| AmIGO   | $\mathcal{O}\left(\epsilon^{-2}\right)$           | $\mathcal{O}\left(\epsilon^{-2}\right)$           | SGD      |
| LazyBLO | $\mathcal{O}\left(\epsilon^{-2}\right)$           | $\mathcal{O}\left(\epsilon^{-2}\right)$           | SGD      |
| MSTSA   | $\mathcal{O}\left(\epsilon^{-2} ight)$            | $\tilde{\mathcal{O}}\left(\epsilon^{-2} ight)$    | Momentum |
| SUSTAIN | $\tilde{\mathcal{O}}\left(\epsilon^{-1.5}\right)$ | $\tilde{\mathcal{O}}\left(\epsilon^{-1.5}\right)$ | Momentum |
| MRBO    | $\mathcal{O}\left(\epsilon^{-1.5}\right)$         | $\tilde{\mathcal{O}}\left(\epsilon^{-1.5}\right)$ | Momentum |
| SEMA    | $\tilde{\mathcal{O}}(\epsilon^{-2})$              | $\tilde{\mathcal{O}}(\epsilon^{-2})$              | Momentum |
| SVRB    | $\mathcal{O}\left(\epsilon^{-1.5}\right)$         | $\mathcal{O}\left(\epsilon^{-1.5}\right)$         | Momentum |
| MA-SOBA | $\mathcal{O}(\epsilon^{-2})$                      | $\mathcal{O}(\epsilon^{-2})$                      | Momentum |
| VRBO    | $\tilde{\mathcal{O}}\left(\epsilon^{-1.5}\right)$ | $\tilde{\mathcal{O}}\left(\epsilon^{-1.5}\right)$ | VR       |
| FSLA    | $\mathcal{O}\left(\epsilon^{-2}\right)'$          | $\mathcal{O}\left(\epsilon^{-2}\right)$           | VR       |

PG: Partial gradient evaluation VR: Variance Reduction

(2) Hessian-Free Bilevel Optimization: To avoid the expensive Hessian evaluations, several Hessian-free bilevel algorithms have been proposed. For example, FO-MAML (Finn et al., 2017; Nichol et al., 2018) ignores the Hessian computation but does not offer any performance guarantee (Antoniou et al., 2018; Fallah et al., 2020). Several approaches have also been proposed to replace the LL problem with optimality-based constraints (Chen et al., 2023b; Liu et al., 2022a; Shen & Chen, 2023). However, these methods mostly focus on deterministic settings rather than stochastic ones. Several zeroth-order

methods have been proposed to approximate the hypergraident (e.g., ES-MAML (Song et al., 2019), 163 HOZOG (Gu et al., 2021), and PZOBO (Sow et al., 2022)). However, ES-MAML and HOZOG do 164 not provide any theoretical convergence guarantee, while PZOBO achieves  $\mathcal{O}\left(d^2\epsilon^{-2}\right)$  to reach an 165  $\epsilon$ -stationary point, where d is the problem dimension. Recently, F<sup>2</sup>SA and F<sup>3</sup>SA (momentum-based 166 version of  $F^2SA$  (Kwon et al., 2023) have been proposed, which are two first-order methods based 167 on the value-function-based lower-level problem reformulation. To reach an  $\epsilon$ -stationary point, F<sup>2</sup>SA and  $F^3SA$  require  $\mathcal{O}(\epsilon^{-3.5})$  and  $\mathcal{O}(\epsilon^{-2.5})$  iterations, respectively. The work in (Chen et al., 2023a) 168 improves the convergence rate for F<sup>2</sup>SA, resulting in a rate of  $\mathcal{O}\left(\epsilon^{-2}\log(1/\epsilon)\right)$ . However, this 169 170 improved rate is still slower than that of our proposed LazyBLO approach by a logarithmic factor. Compared to (Kwon et al., 2023), our proposed LazyBLO algorithm strikes a good balance in terms 171 of the use of Hessian information: On one hand, we leverage Hessian information to maintain good 172 convergence performance; on the other hand, we infrequently use Hessian information to significantly 173 reduce the wall-clock time. 174

3 Other Uses of Infrequent Hessian Evaluations: We note that infrequent Hessian evaluations have also been used for speeding up second-order methods for single-level optimization problems (Shaman-skii, 1967; Adler et al., 2020; Doikov et al., 2022; Lampariello & Sciandrone, 2001; Wang et al., 2006; Fan, 2013). However, in bilevel optimization, the Hessian information *necessarily* emerges due to the hypergradient computation, rather than as a "second-order" option in single-level optimization. Also, due to the complex problem structure, analyzing the use of infrequent Hessian in bilevel optimization is far more challenging than in a single-level setting.

182 183

188

189

190

191

213 214

#### **3** PRELIMINARIES

In this section, we provide some preliminaries for solving Problem (1). We first state a set of assumptions that are needed to establish the convergence of LazyBLO:

187 Assumption 3.1 (UL Objective).  $f(\mathbf{x}, \mathbf{y})$  satisfies:

- 1) For any  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^u \times \mathbb{R}^l$ ,  $\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$  is Lipschitz continuous (w.r.t.  $\mathbf{y}$ ) with constant  $L_{f_x} \ge 0$ , and  $\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$  is Lipschitz continuous (w.r.t. both  $\mathbf{x}$  and  $\mathbf{y}$ ) with constant  $L_{f_y} \ge 0$ .
- 2) For any  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{u} \times \mathbb{R}^{l}$ , we have  $\|\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})\| \leq B_{f_{u}}$  for some constant  $B_{f_{u}} \geq 0$ .
- 192 Assumption 3.2 (LL Objective).  $g(\mathbf{x}, \mathbf{y})$  satisfies:
- 193 *I*) For any  $\mathbf{x} \in \mathbb{R}^{u}$ ,  $g(\mathbf{x}, \cdot)$  is  $\mu_{g}$ -strongly convex with respect to  $\mathbf{y}$  for some  $\mu_{g} > 0$ .
- 195 2) For any  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^u \times \mathbb{R}^l$ ,  $\nabla_{\mathbf{y}} g(\mathbf{x}, \mathbf{y})$  is Lipschitz continuous (w.r.t. y) with constant  $L_g \ge 0$ , 196 and  $\nabla^2_{\mathbf{xy}} g(\mathbf{x}, \mathbf{y})$  and  $\nabla^2_{\mathbf{yy}} g(\mathbf{x}, \mathbf{y})$  are Lipschitz continuous (w.r.t. both x and y) with constants 197  $L_{g_{xy}} \ge 0$  and  $L_{g_{yy}} \ge 0$ , respectively.
- 3) For any  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{u} \times \mathbb{R}^{l}$ , we have  $\left\| \nabla_{\mathbf{x}\mathbf{y}}^{2} g\left(\mathbf{x}, \mathbf{y}\right) \right\| \leq B_{g_{xy}}$  for some constant  $B_{g_{xy}} > 0$ .

Note that all the above assumptions are standard in the analysis of bilevel optimization problems (e.g., Ghadimi & Wang (2018); Hong et al. (2020); Khanduri et al. (2021b); Liu et al. (2022b); Qiu et al. (2022)). With the above assumptions and using implicit function theorem (Rudin et al., 1976), the hypergradient of  $\ell(\cdot)$  can be computed as  $\nabla \ell(\mathbf{x}) = \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) - \nabla_{\mathbf{xy}}^2 g(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) \left[\nabla_{\mathbf{yy}}^2 g(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))\right]^{-1} \nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}^*(\mathbf{x})).$ 

Instead of computing the Hessian inverse explicitly, there are different ways to approximate the Hessian inverse or HVPs in bilevel optimization, such as conjugate gradient (CG) (Pedregosa, 2016) and Neumann series (Ghadimi & Wang, 2018) methods. For example, stocBiO (Ji et al., 2021) uses Neumann series, while AID-BiO (Ji et al., 2021), AID-CG (Grazzi et al., 2020) and AmIGO (Arbel & Mairal, 2022) implement CG. In this paper, we use CG to efficiently estimate the HVPs  $([\nabla^2_{yy}g(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))]^{-1} \nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}^*(\mathbf{x})))$ , which finds the minimizer of a quadratic function by solving a linear system derived from the hypergradient. The quadratic optimization problem is formulated as follows:

$$\min_{\mathbf{z}\in\mathbb{R}^l} q(\mathbf{x},\mathbf{y},\mathbf{z}) \triangleq \frac{1}{2} \mathbf{z}^\top \nabla_{\mathbf{y}\mathbf{y}}^2 g(\mathbf{x},\mathbf{y}) \mathbf{z} + \mathbf{z}^\top \nabla_{\mathbf{y}} f(\mathbf{x},\mathbf{y}).$$
(3)

For the function  $q(\cdot, \cdot, \cdot)$  defined in Eq. (3), the following lemma together with Assumption 3.2 implies that  $q(\mathbf{x}, \mathbf{y}, \mathbf{z})$  is  $\mu_q$ -strongly convex and  $L_q$ -Lipschitz smooth.

Lemma 3.3 (Quadratic Problem). For any  $(\mathbf{x}, \mathbf{y}, \mathbf{z})$ , the quadratic problem  $q(\mathbf{x}, \mathbf{y}, \mathbf{z})$  with respect to  $\mathbf{z}$  is Lipschitz-smooth with constant  $L_q \ge 0$ .

The admitted unique minimizer  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$  of Eq. (3) can then be utilized to compute the hypergradient estimate as  $\nabla \ell(\mathbf{x}) = \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) + \nabla^2_{\mathbf{xy}} g(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) \mathbf{z}^*(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))$ . Since it is challenging to obtain  $\mathbf{y}^*(\mathbf{x})$  and  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$  in closed form, it is natural to consider their approximations. Specifically, let  $\bar{\mathbf{y}}$  and  $\bar{\mathbf{z}}$  be some approximations of  $\mathbf{y}^*(\mathbf{x})$  and  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$ , respectively. Then, we have the approximation for  $\nabla \ell(\mathbf{x})$  defined as follows:

224 225

233

$$\nabla f(\mathbf{x}, \bar{\mathbf{y}}, \bar{\mathbf{z}}) = \nabla_{\mathbf{x}} f(\mathbf{x}, \bar{\mathbf{y}}) + \nabla_{\mathbf{x}\mathbf{y}}^2 g(\mathbf{x}, \bar{\mathbf{y}}) \bar{\mathbf{z}}.$$
(4)

Since Problem (1) can potentially be a large-scale stochastic optimization problem, computing a full gradient approximation in Eq. (4) can be computationally expensive. To address this challenge, a common approach for evaluating Eq. (4) is to build a stochastic gradient estimator. Define stochastic approximations as  $f(\mathbf{x}, \mathbf{y}; \mathcal{D}^f) \triangleq \frac{1}{|\mathcal{D}^f|} \sum_{\xi \in \mathcal{D}^f} f(\mathbf{x}, \mathbf{y}; \xi)$  and  $g(\mathbf{x}, \mathbf{y}; \mathcal{D}^g) \triangleq \frac{1}{|\mathcal{D}^g|} \sum_{\zeta \in \mathcal{D}^g} g(\mathbf{x}, \mathbf{y}; \zeta)$ , where  $\mathcal{D}^f$  and  $\mathcal{D}^g$  are the batches of independent and identically distributed samples with sizes  $|\mathcal{D}^f| \ge 1$  and  $|\mathcal{D}^g| \ge 1$ , respectively. Then, a stochastic estimator of Eq. (4) can be computed as:

$$\nabla f(\mathbf{x}, \mathbf{y}, \mathbf{z}; \mathcal{D}^{f_x}, \mathcal{D}^{g_{xy}}) = \nabla_{\mathbf{x}} f\left(\mathbf{x}, \mathbf{y}; \mathcal{D}^{f_x}\right) + \nabla_{\mathbf{xy}}^2 g\left(\mathbf{x}, \mathbf{y}; \mathcal{D}^{g_{xy}}\right) \mathbf{z}.$$

Here, for simplicity, we slightly abuse the notations  $\bar{\mathbf{y}}$  and  $\bar{\mathbf{z}}$  as  $\mathbf{y}$  and  $\mathbf{z}$  in the above equation and the rest of the paper as long as there is no confusion from the context. For  $\nabla f(\mathbf{x}, \mathbf{y}, \mathbf{z}; \mathcal{D}^{f_x}, \mathcal{D}^{g_{xy}})$  and  $\nabla_{\mathbf{y}}g(\mathbf{x}, \mathbf{y}; \mathcal{D}^{g_y})$ , we make the following typical assumption in stochastic optimization analysis.

Assumption 3.4 (Stochastic Gradients). For any  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{u} \times \mathbb{R}^{l}$  and data batch  $\mathcal{D}^{f_{x}}, \mathcal{D}^{f_{y}}, \mathcal{D}^{g_{y}}, \mathcal{D}^{g_{yy}}$  and  $\mathcal{D}^{g_{yy}}, define <math>\sigma_{f_{x}}^{2} \triangleq \tilde{\sigma}_{f_{x}}^{2} |\mathcal{D}^{f_{x}}|^{-1}, \sigma_{f_{y}}^{2} \triangleq \tilde{\sigma}_{f_{y}}^{2} |\mathcal{D}^{f_{y}}|^{-1}, \sigma_{g_{y}}^{2} \triangleq \tilde{\sigma}_{g_{y}}^{2} |\mathcal{D}^{g_{y}}|^{-1}, \sigma_{g_{y}}^{2} \equiv \tilde{\sigma}_{g_{y}}^{2} |\mathcal{D}^{g_{y}}|^{-1}, \sigma_{g_{y}}^{2} = \tilde{\sigma}_{g_{y}}^{2} |\mathcal{D}^{g_{y}}|^{-1}, \sigma_{g_{y}}^{2} |\mathcal{D}^{g_{y}}|^{-1}$ 

$$\begin{split} \mathbb{E}[\|\nabla_{\mathbf{x}}f(\mathbf{x},\mathbf{y};\mathcal{D}^{f_x}) - \nabla_{\mathbf{y}}f(\mathbf{x},\mathbf{y})\|^2] &\leq \sigma_{f_x}^2, \qquad \mathbb{E}[\|\nabla_{\mathbf{y}}f(\mathbf{x},\mathbf{y};\mathcal{D}^{f_y}) - \nabla_{\mathbf{y}}f(\mathbf{x},\mathbf{y})\|^2] \leq \sigma_{f_y}^2, \\ \mathbb{E}[\|\nabla_{\mathbf{y}}g(\mathbf{x},\mathbf{y};\mathcal{D}^{g_y}) - \nabla_{\mathbf{y}}g(\mathbf{x},\mathbf{y})\|^2] &\leq \sigma_{g_y}^2, \qquad \mathbb{E}[\|\nabla_{\mathbf{x}\mathbf{y}}^2g(\mathbf{x},\mathbf{y};\mathcal{D}^{g_{xy}}) - \nabla_{\mathbf{x}\mathbf{y}}^2g(\mathbf{x},\mathbf{y})\|^2] \leq \sigma_{g_{xy}}^2, \\ \mathbb{E}[\|\nabla_{\mathbf{y}\mathbf{y}}^2g(\mathbf{x},\mathbf{y};\mathcal{D}^{g_{yy}}) - \nabla_{\mathbf{y}\mathbf{y}}^2g(\mathbf{x},\mathbf{y})\|^2] \leq \sigma_{g_{xy}}^2. \end{split}$$

Lastly, we define the following performance metrics for solving the Problem (1).

**Definition 3.5** ( $\epsilon$ -Stationary Point). Point **x** is an  $\epsilon$ -stationary point if  $\mathbb{E}[\|\nabla \ell(\mathbf{x})\|^2] \leq \epsilon$ , where **x** is the output of a stochastic algorithm, and the expectation is taken over all randomness of the algorithm.

**Definition 3.6** ( $\epsilon$ -Optimal Point). Point **x** is an  $\epsilon$ -optimal point if  $\mathbb{E} \left[ \ell(\mathbf{x}) - \ell^* \right] \leq \epsilon$ , where  $\ell^* \triangleq \min_{\mathbf{x} \in \mathbb{R}^u} \ell(\mathbf{x})$ , and **x** is the output of a stochastic algorithm. The expectation is taken over all randomness of the algorithm.

256 257 258

259

254

255

### 4 THE LazyBLO ALGORITHM

In this section, we propose a new algorithmic framework called LazyBLO to solve the bilevel
 optimization problem in Eq. (1). Our goal is to reduce the computation of the HVPs, and our key
 idea is to update the HVP periodically on a subset of the entire training iterations while using stale
 Hessian information in the remaining iterations.

The most basic algorithm in the LazyBLO framework incorporates SGD-style updates, which is illustrated in Algorithm 1. We note that more sophisticated algorithms in the LazyBLO framework can include advanced algorithmic techniques, such as momentum and/or variance reduction to accelerate the convergence and enhance other performances. As shown in Algorithm 1, the LazyBLO framework uses a double-loop structure and constructs iterates  $\mathbf{x}_t^n$ ,  $\mathbf{y}_t^n$  and  $\mathbf{z}_t$ , where the inner iteration counter n goes from 0 to N - 1 and the outer iteration counter t runs from 0 to T - 1, so that  $\mathbf{x}_t^n$  approaches a stationary point of  $\ell(\cdot)$ , and  $\mathbf{y}_t^n$  and  $\mathbf{z}_t$  keep track of the quantities  $\mathbf{y}^*(\mathbf{x}_t^n)$  and  $\mathbf{z}^*(\mathbf{x}_t^N, \mathbf{y}_t^N)$ . In

| Algorithm 1 The LazyBLO Algorithmic Framwork with Basic SGD-type Updates.  |
|--|
| <b>Input:</b> Initial parameters $\mathbf{x}_0^0, \mathbf{v}_0^0, \mathbf{z}_0$ , and stepsize $\{\alpha_t\}_{t=0}^{T-1}, \{\beta_t\}_{t=0}^{T-1}, \{\gamma_t\}_{t=0}^{T-1}$ |
| for $t = 0$ to $T - 1$ do  |
| for $n = 0$ to $N - 1$ do  |
| Initialize $\mathbf{x}_t^0 = \mathbf{x}_{t-1}^N$ and $\mathbf{y}_t^0 = \mathbf{y}_{t-1}^N$   |
| Sample data batches $\mathcal{D}_{t,n}^{g}$ , $\mathcal{D}_{t,n}^{f_x}$ , and $\mathcal{D}_{t,n}^{g_{xy}}$   |
| Compute the gradient estimate $h_{t,n}^g$ using (6) and update $\mathbf{y}_t^{n+1} = \mathbf{y}_t^n - \beta_t h_{t,n}^g$   |
| Compute the gradient estimate $h_{t,n}^f$ using (5) and update $\mathbf{x}_{t,n}^{n+1} = \mathbf{x}_{t,n}^n - \alpha_t h_{t,n}^f$  |
| end for  |
| Sample data batches $\mathcal{D}_t^{g_{yy}}$ and $\mathcal{D}_t^{f_y}$   |
| Compute the gradient estimate $h_{t,n}^{q}$ using (7) and update $\mathbf{z}_{t+1} = \mathbf{z}_t - \gamma_t h_t^q$  |
| end for  |

the inner loop, the algorithm updates  $\mathbf{x}_t^n$  and  $\mathbf{y}_t^n$  using the stochastic gradient estimators  $h_{t,n}^f$  and  $h_{t,n}^g$  defined as:

$$h_{t,n}^{f} = \nabla_{\mathbf{x}} f\left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}; \mathcal{D}_{t,n}^{f_{x}}\right) + \nabla_{\mathbf{x}\mathbf{y}}^{2} g\left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}; \mathcal{D}_{t,n}^{g_{xy}}\right) \mathbf{z}_{t},$$
(5)

$$h_{t,n}^g = \nabla_{\mathbf{y}} g\left(\mathbf{x}_t^n, \mathbf{y}_t^n; \mathcal{D}_{t,n}^g\right).$$
(6)

The variable  $z_t$  in Eq. (5) is updated in the outer loop using a stochastic gradient estimator  $h_t^q$  as:

$$h_t^q = \nabla_{\mathbf{yy}}^2 g(\mathbf{x}_t^N, \mathbf{y}_t^N; \mathcal{D}_t^{g_{yy}}) \mathbf{z}_t + \nabla_{\mathbf{y}} f(\mathbf{x}_t^N, \mathbf{y}_t^N; \mathcal{D}_t^{f_y}).$$
(7)

Note that, compared to  $h_{t,n}^f$  and  $h_{t,n}^g$ , only  $h_t^q$  contains the HVP, and is computed *infrequently* after every N inner loop iterations. In addition, N needs to be chosen with a tolerable approximation error of the HVP. If N gets too large, the error of the HVP approximation would also increase, thus inevitably degrading the performance of LazyBLO. With less frequent Hessian computations, LazyBLO executes faster per iteration in terms of *wall-clock time* compared to standard bilevel algorithms that require multiple Hessian/vector evaluations in each round of updates (Ghadimi & Wang, 2018; Arbel & Mairal, 2022; Ji et al., 2021; Chen et al., 2021), resulting in a significant reduction in computational cost and savings in implementation time.

299 Another insightful remark on the Jacobian-vector product in (5) is also in order. To date, most 300 of the existing bilevel algorithms compute only one single Jacobian-vector product (JVP) in each 301 iteration, whereas HVPs are computed multiple times in each iteration even in some single-loop 302 bilevel algorithms (e.g., SUSTAIN (Khanduri et al., 2021b), TTSA (Hong et al., 2020), BSA (Ghadimi & Wang, 2018), and ALSET (Chen et al., 2021)). Due to this difference between JVP and HVP in 303 304 bilevel optimization algorithms, reducing the number of HVP computations is far more important than reducing the computations of JVPs. Therefore, we only focus on reducing the HVP in this paper. 305 We further note that reducing the computation of JVPs can be done in a similar manner as the HVPs 306 established in our work. 307

308

283

284

287 288

289 290

### 5 THEORETICAL PERFORMANCE ANALYSIS

309 310

In this section, we conduct the theoretical convergence analysis for the LazyBLO framework for 311 solving the bilevel optimization problem in Eq. (1). Note that, although LazyBLO executes faster per 312 iteration, we have a noisier hypergradient due to the use of stale Hessian information. As a result, 313 it remains unclear whether LazyBLO can converge and, if yes, what theoretical convergence rate 314 (i.e., iteration complexity) it can achieve. Intuitively, due to the lazy Hessian information updates, 315 one can expect that the theoretical convergence rate of LazyBLO cannot outperform its non-lazy 316 counterpart. Surprisingly, in this paper, we show that LazyBLO achieves the same convergence rate as 317 their non-lazy counterpart. This, together with the much lower per-iteration wall-clock time, implies 318 that LazyBLO will enjoy a much faster speed in terms of wall-clock time. This will also be verified 319 by our experiments in Section 6. 320

- 321 The convergence analysis for LazyBLO is highly non-trivial due to the following technical challenges:
- *i)* The use of lazy Hessian evaluation increases the error of stochastic gradient estimator  $h_{t,n}^{\dagger}$  for the upper-level function; *ii)* Due to the hierarchical and coupled structure of bilevel optimization
- problems, the error resulting from the stochastic gradient estimator  $h_{t,n}^f$  with stale Hessian information

further propagates to and increases the approximation error of  $\mathbf{y}^*(\mathbf{x})$  and the approximation error of z<sup>\*</sup> (x, y). What is even worse is that  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$  is also associated with  $\mathbf{y}^*(\mathbf{x})$ . All the above complex couplings of *laziness-induced errors* and the complications associated with these approximation errors are *unseen* in bilevel optimization algorithm analysis, which significantly increases the difficulty of analyzing the convergence of LazyBLO and necessitate *new* proof techniques.

#### 330 5.1 SUPPORTING LEMMAS

Toward this end, we first state two basic lemmas needed for the convergence analysis of LazyBLO.

Lemma 5.1 (Lemma 2.2 in (Ghadimi & Wang, 2018), Proposition 6 in (Arbel & Mairal, 2022)). *Under Assumptions 3.1 and 3.2, we have* 

329

331

332 333

337 338

344

345

346

347

348

349

350

351

352 353 354

355

356

357 358

359 360

361

362

 $\|\mathbf{y}^{*}(\mathbf{x}_{1}) - \mathbf{y}^{*}(\mathbf{x}_{2})\| \le L_{y} \|\mathbf{x}_{1} - \mathbf{x}_{2}\|, \qquad \|\nabla \ell(\mathbf{x}_{1}) - \nabla \ell(\mathbf{x}_{2})\| \le L_{l} \|\mathbf{x}_{1} - \mathbf{x}_{2}\|,$ 

for all  $\mathbf{x}, \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^u$ , and  $\bar{\mathbf{y}}, \bar{\mathbf{z}} \in \mathbb{R}^l$ , where the Lipschitz constants above are defined as:

$$L_{f} = \max\left\{L_{f_{x}} + \left(L_{g_{xy}}B_{f_{y}}/\mu_{g}\right) + B_{g_{xy}}L_{z}, B_{g_{xy}}\right\}, \quad L_{l} = L_{f}^{'} + \left(L_{f}^{'}B_{g_{xy}}^{2}/\mu_{g}\right), \quad L_{y} = B_{g_{xy}}^{2}/\mu_{g},$$

 $\left\|\nabla f(\mathbf{x}, \bar{\mathbf{y}}, \bar{\mathbf{z}}) - \nabla \ell(\mathbf{x})\right\| \le L_f\left(\left\|\bar{\mathbf{y}} - \mathbf{y}^*(\mathbf{x})\right\| + \left\|\bar{\mathbf{z}} - \mathbf{z}^*(\mathbf{x}, \bar{\mathbf{y}})\right\|\right),$ 

and where 
$$L_{f}^{'} = L_{f_x} + (L_{f_y}B_{g_{xy}}^2/\mu_g) + B_{f_y} [(L_{g_{xy}}/\mu_g) + (L_{g_{yy}}B_{g_{xy}}^2/\mu_g^2)]$$
.

We note that Lemma 5.1 plays a key role in the analysis of AID-based bilevel algorithms. First of all, it characterizes the bias of the implicit gradient as a function of approximation error in  $\bar{\mathbf{y}}$  and  $\bar{\mathbf{z}}$  (see Eq. (4)). It also ensures the Lipschitzness of the mapping  $\mathbf{y}^*(\mathbf{x})$  in characterizing the behavior of the LL problem's iterates. Most importantly, Lemma 5.1 establishes the Lipschitz-smoothness of the implicit function  $\ell(\cdot)$ , which allows the development of SGD-type algorithms for solving stochastic bilevel problems. To complement Lemma 5.1, next result states the properties of the optimal solution  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$  of the quadratic problem in Eq. (3).

**Lemma 5.2** (Proposition 6 in(Arbel & Mairal, 2022)). Under Assumptions 3.1 and 3.2,  $\forall \mathbf{x}, \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^u$  and  $\mathbf{y}, \mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^l$ , we have

$$\|\mathbf{z}^{*}(\mathbf{x}_{1},\mathbf{y}_{1}) - \mathbf{z}^{*}(\mathbf{x}_{2},\mathbf{y}_{2})\| \le L_{z}\left(\|\mathbf{x}_{1} - \mathbf{x}_{2}\| + \|\mathbf{y}_{1} - \mathbf{y}_{2}\|\right), \qquad \|\mathbf{z}^{*}(\mathbf{x},\mathbf{y})\| \le B_{f_{y}}/\mu_{g},$$

where  $L_{z} = (L_{g_{yy}}B_{f_{y}}/\mu_{q}^{2}) + L_{f_{y}}/\mu_{g}$ .

Lemma 5.2 also plays a key role in the analysis of LazyBLO as it is utilized to bound the drift in the Hessain vector product estimates (see Eq. (3)). Next, we present the main results of the paper.

5.2 MAIN RESULTS

① The Non-convex  $\ell(\mathbf{x})$  Setting: By leveraging Lemmas 5.1 and 5.2, we establish the main convergence result of the proposed LazyBLO for non-convex  $\ell(\mathbf{x})$  in Theorem 5.3.

**Theorem 5.3** (Non-Convex  $\ell(\mathbf{x})$ ). Under Assumptions 3.1–3.4, with step-sizes  $\alpha_t = \alpha = \mathcal{O}\left(\frac{1}{N^2}\right)$ ,  $\beta_t \triangleq c_{\beta}\alpha = \mathcal{O}\left(\frac{1}{N^2}\right)$ , and  $\gamma_t \triangleq c_{\gamma}\alpha = \mathcal{O}(1)$  for all  $t \in \{0, 1, \dots, T-1\}$ , where  $c_{\beta}$  and  $c_{\gamma}$  are defined in Eq. 22 in Appendix C. Then, the iterates generated by LazyBLO satisfy:

366 367 368

369

370

$$\frac{1}{TN}\sum_{t=0}^{T-1}\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\nabla\ell\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] = \mathcal{O}\left(\frac{N\Delta_{0}}{T}\right) + \mathcal{O}\left(\sigma_{g_{y}}^{2} + \sigma_{g_{xy}}^{2} + \sigma_{f_{x}}^{2} + \sigma_{g_{yy}}^{2} + \sigma_{f_{y}}^{2}\right),$$

where  $\Delta_0 = (\ell(\mathbf{x}_0^0) - \ell^*) + \|\mathbf{y}_0^0 - \mathbf{y}^*(\mathbf{x}_0^0)\|^2 + \|\mathbf{z}_0 - \mathbf{z}^*(\mathbf{x}_0^0, \mathbf{y}_0^0)\|^2.$ 

The proof of Theorem 5.3 can be found in Appendix C. Theorem 5.3 establishes the convergence of LazyBLO under the most general setting, where the implicit function  $\ell(\cdot)$  can be non-convex. The result characterizes the effect of different parameters on the convergence of LazyBLO. Specifically, as N increases, the performance of LazyBLO degrades. This is unsurprising since more stale Hessian information is expected to slow the convergence. Hence, N should be chosen below a certain threshold to maintain the accuracy of the hypergradient estimations. On the other hand, to enjoy the benefits of the LazyBLO approach, N is supposed to be strictly larger than 1. We can potentially choose N = 1, and our algorithm, which becomes fully single-loop, recovers standard results for bilevel algorithms under the same assumptions as ours (e.g., the guarantees achieved in (Arbel & Mairal, 2022)). Interestingly, under an appropriate *N*-value, the *N*-dependent slowdown effect in LazyBLO can be offset by Hessian computations skippings, allowing LazyBLO to run even faster than non-lazy approaches in terms of wall-clock time.

<sup>382</sup> Our next result characterizes the computation complexity of LazyBLO.

**Corollary 5.4** (Computation Complexity). Under the setting of Theorem 5.3, choose  $|\mathcal{D}^{f_x}|, |\mathcal{D}^{f_y}|, |\mathcal{D}^{g_y}|, |\mathcal{D}^{g_{xy}}|, |\mathcal{D}^{g_{yy}}| = \Theta(\epsilon^{-1})$ . Then, LazyBLO requires  $\mathcal{O}(\epsilon^{-2})$  partial gradient and HVP evaluations to reach an  $\epsilon$ -stationary point.

380

387 We note that, when  $\epsilon$  is small, the batch sizes in Corollary 5.4 could be large. However, it is worth 388 noting that the use of large batch sizes is not a consequence of the proposed LazyBLO algorithm 389 design; rather, these batch size choices are common in the literature, as the above guarantees are the 390 same as those achieved in standard SGD-based bilevel algorithms (e.g., (Arbel & Mairal, 2022; Ji 391 et al., 2021; Huang et al., 2022)) that require the computation of (multiple) Hessian/HVPs in each 392 iteration. It is also worth noting that the large batch sizes are required only for theoretical analysis and can be eliminated by using a third-order Lipschitz assumption, as done by SOBA (Dagréou et al., 393 2022). In our experiments, we use a small batch size instead, and our algorithm still outperforms the 394 baseline algorithms. 395

Given that LazyBLO can converge and even matches the performance of SOTA non-lazy bilevel methods, another question also arises: under which settings could LazyBLO theoretically outperform current bilevel approaches? The next result shows that if the LL problem is deterministic, we can, in fact, improve upon the current approaches and reduce the HVP evaluations from  $O(\epsilon^{-2})$  to  $O(\epsilon^{-1})$ .

400 **Corollary 5.5** (Computation Complexity for Deterministic LL Problems). Suppose the lower-level 401 problem is deterministic. Under the condition of Theorem 5.3, LazyBLO requires  $\mathcal{O}(\epsilon^{-1})$  for HVP 402 evaluations to achieve an  $\epsilon$ -stationary point.

Corollary 5.5 suggests that LazyBLO significantly reduces the HVP evaluations. In contrast, for standard bilevel optimization algorithms, the HVPs stay the same as the total number of rounds required by an algorithm to achieve the  $\epsilon$ -stationary solution. For example, under the same deterministic setting, the baseline methods BSA (Ghadimi & Wang, 2018), stocBiO (Ji et al., 2021) and AmIGO (Arbel & Mairal, 2022) require  $\mathcal{O}(\epsilon^{-2})$  gradient computations and TTSA (Hong et al., 2020) requires  $\mathcal{O}(\epsilon^{-2.5})$  gradient computations, which is equivalent to the number of outer function's gradients evaluated during the execution of the algorithm.

So far, our results characterize the performance of LazyBLO in the non-convex settings. However, for some problems (e.g., quadratic UL and LL problems), the implicit function may have additional structures that might lead to better convergence of LazyBLO. Next, we characterize the performance of LazyBLO when the implicit function is strongly convex, which is often of interest for applications in robust and inverse optimization, optimal control in robotics and aerospace with quadratic cost, etc.

(2) The Strongly Convex  $\ell(\mathbf{x})$  Setting: Under the setting where  $\ell(\cdot)$  is  $\mu_f$ -strongly convex, we provide a stronger performance guarantee for the convergence of LazyBLO, which is stated as follows:

**Theorem 5.6** (Strongly Convex  $\ell(\mathbf{x})$ ). Suppose the upper-level function  $\ell(\mathbf{x})$  is  $\mu_f$ -strongly-convex. Under Assumptions 3.1–3.4, choose the step-sizes  $\alpha_t = \alpha = \mathcal{O}\left(\frac{1}{N}\right)$ ,  $\beta_t \triangleq \hat{c}_{\beta}\alpha = \mathcal{O}\left(\frac{1}{N}\right)$  and  $\gamma_t \triangleq$   $\hat{c}_{\gamma}\alpha = \mathcal{O}\left(\frac{1}{N^2}\right)$  for all  $t \in \{0, 1, \dots, T-1\}$ , where  $\hat{c}_{\beta}$  and  $\hat{c}_{\gamma}$  are defined in Eq. 34 in Appendix D. Then, the iterates generated by LazyBLO satisfy:

$$\frac{1}{N}\sum_{n=0}^{N-1}\mathbb{E}\bigg[\ell\left(\mathbf{x}_{t}^{n}\right)-\ell^{*}\bigg] \leq (1-\mu_{f}\alpha)^{t}\hat{\Delta}_{0}+\mathcal{O}\bigg(\sigma_{g_{xy}}^{2}+\sigma_{f_{x}}^{2}+\frac{1}{N^{4}}\sigma_{g_{yy}}^{2}+\frac{1}{N^{4}}\sigma_{f_{y}}^{2}+\frac{1}{N}\sigma_{g_{y}}^{2}\bigg),$$

for any  $t \ge 1$ , where  $\hat{\Delta}_0 = \frac{1}{N} \sum_{n=0}^{N-1} (\ell(\mathbf{x}_0^n) - \ell^*) + \frac{1}{N} \sum_{n=0}^{N-1} \|\mathbf{y}_0^n - \mathbf{y}^*(\mathbf{x}_0^n)\|^2 + \frac{1}{N} \|\mathbf{z}_0 - \mathbf{z}_0^*\|^2$ .

425 426 427

428

429

422 423 424

The detailed proof of Theorem 5.6 is provided in Appendix D due to space limitations. Theorem 5.6 demonstrates that, under the strongly convex setting, LazyBLO achieves a much faster linear convergence rate. Theorem 5.6 also immediately implies the following computation complexity:

430 431 **Corollary 5.7** (Computation Complexity). Under the setting of Theorem 5.6, choosing  $|\mathcal{D}^{f_x}| = |\mathcal{D}^{g_{xy}}| = \Theta(\epsilon^{-1}), |\mathcal{D}^{g_y}| = \Theta(N^{-1}\epsilon^{-1}), and |\mathcal{D}^{f_y}| = |\mathcal{D}^{g_{yy}}| = \Theta(N^{-4}\epsilon^{-1}), LazyBLO requires$ 



Figure 1: Comparison for data hyper-cleaning on Figure 2: Training loss for deep hyper-MNIST (corruption rate p = 0.1, 10 repetitions). representation on CIFAR-10 (10 repetitions).

 $\mathcal{O}(\epsilon^{-1}\log\epsilon^{-1})$  partial gradient evaluations and  $\mathcal{O}(N^{-4}\epsilon^{-1}\log\epsilon^{-1})$  HVP evaluations to reach an  $\epsilon$ -optimal point.

Corollary 5.7 shows that LazyBLO significantly reduces the number of HVP evaluations. Again, note that the complexity of partial gradient evaluations in Corollary 5.7 matches the same guarantee achieved in (Arbel & Mairal, 2022), which is obtained by multiple Hessian evaluations per iteration.

#### 6 NUMERICAL RESULTS

440

441

442

443 444

445

446

447 448

449

In this section, we verify the theoretical performance of LazyBLO on different optimization tasks
 and with two different datasets: 1) data hyper-cleaning on the MNIST dataset, and 2) deep hyper representation with the ResNet network on the CIFAR-10 dataset. Due to space limitations, additional
 experimental details and results are included in Appendix B.

Task 1) Data Hyper-Cleaning on the MNIST Dataset: We conduct experiments on a data hyper-cleaning task with MNIST dataset (LeCun et al., 1998). Data hyper-cleaning aims to train a classifier on a corrupted dataset. We compare LazyBLO with stochastic bilevel algorithms AmIGO (Arbel & Mairal, 2022), stocBiO (Ji et al., 2021), BSA (Ghadimi & Wang, 2018), and MRBO (Yang et al., 2021) as baselines. We also perform data hyper-cleaning with fully single-loop bilevel algorithms TTSA (Hong et al., 2020), SOBA (Dagréou et al., 2022), and MA-SOBA (Chen et al., 2024).

Table 2 shows that TTSA, SOBA, and MA-SOBA all need an exceedingly long time to converge.
Specifically, the convergence of TTSA, SOBA, and MA-SOBA are 74, 73, and 82× slower, respectively, than LazyBLO. In addition, TTSA, SOBA, and MA-SOBA require 390, 126, and 150× more Hessian computations, respectively, compared to LazyBLO. Given the significantly slow convergence of these fully single-loop bilevel algorithms, we exclude them from the following comparison.

465 From Fig. 1a, we can see that AmIGO 466 and stocBiO have 467 similar conver-468 gence performance. 469 LazyBLO outperforms 470 all baseline methods 471 in terms of wall-clock 472 time, which shows the

Table 2: Convergence performance of TTSA, SOBA, and MA-SOBA compared with our LazyBLO on data hyper-cleaning on MNIST (corruption rate p = 0.1, average over 10 repetitions).

| Algorithm | WALL-CLOCK TIME | # OF HESSIAN | TRAINING LOSS |
|-----------|-----------------|--------------|---------------|
| TTSA      | 4290 s          | 1950         | 3.95          |
| SOBA      | 4210 s          | 630          | 3.28          |
| MA-SOBA   | 4740 \$         | 750          | 3.05          |
| LazyBLO   | 58 s            | 5            | 2.35          |

473 advantages of LazyBLO. Specifically, it only takes LazyBLO approximately 60 seconds to converge, 474 while AmIGO and stocBiO converge in around 100 seconds. This much-improved wall-clock 475 time is due to the fact that LazyBLO uses stale Hessian information and saves a lot of Hessian 476 computation time. It is worth pointing out that the comparison with MRBO is not entirely fair since 477 MRBO is equipped with more sophisticated momentum techniques to accelerate convergence, while LazyBLO only uses vanilla-SGD updates. LazyBLO can also be equipped with momentum-based 478 SGD updates to further accelerate the convergence. Furthermore, the training loss of LazyBLO is 479 similar to those of AmIGO, stocBiO, and BSA, which use up-to-date Hessian information during 480 the training. This result is surprising because LazyBLO with stale Hessian information can still 481 match the methods with non-lazy Hessian updates. This implies that the Hessian information evolves 482 gradually during the training, and one may use stale Hessians to construct good approximations of 483 the hypergradient in bilevel optimization. 484

485 It can be seen in Fig. 1b that the convergence speed with respect to the number of Hessian evaluations for LazyBLO is much faster compared with all the baseline algorithms (see the zoomed-in area in

Fig. 1b). Table 3 also demonstrates that, to achieve the same convergence training loss, AmIGO, stocBiO and BSA all need 252 Hessian computations, while LazyBLO only needs 5 Hessian computations (i.e., 50× faster). Note that we do not include MRBO in this table since it has a higher error floor compared to other algorithms. As a consequence, it can not reach the same training loss as the other algorithms.

491 Fig. 3 captures the effect of N on the perfor-492 mance of LazyBLO. Specifically, we observe 493 that as we increase the value of N, the execu-494 tion of the algorithm becomes faster and faster. 495 However, we note that increasing N beyond 496 a certain threshold may not yield additional benefits and could even lead to performance 497 degradation. This is because, as N increases, 498 the difference between stale and fresh Hes-499 sian information becomes larger, potentially 500 causing the hypergradient  $h_{t,n}^f$  to become less 501





Figure 3: Comparison of LazyBLO on data hypercleaning on MNIST at a different # of x-updates (N).

accurate and adversely affecting the training loss of LazyBLO.

Task 2) Deep Hyper-Representation with ResNet-20 on the CIFAR-10 Dataset: To demonstrate 504 the effectiveness of LazyBLO in training neural networks, we conduct experiments on a deep hyperrepresentation task (Yang et al., 2023; Sow et al., 2022) with the ResNet-20 model (He et al., 2016) on 505 CIFAR-10 dataset (Krizhevsky et al., 2009), which aims to classify CIFAR-10 images. We compare 506 LazyBLO with a standard stochastic bilevel algorithm AmIGO (Arbel & Mairal, 2022), and two 507 fully first-order (Hessian/Jacobian-free) stochastic bilevel algorithms  $F^2SA$  (Kwon et al., 2023) and 508  $F^3SA$  (Kwon et al., 2023) as baselines. We do not compare LazyBLO with other baselines from the 509 previous data hyper-cleaning experiments since stocBiO performs almost identically to AmIGO, and 510 they both outperform MRBO in terms of training loss and BSA in terms of wall-clock time. 511

As shown in Fig. 2a, LazyBLO converges 512 faster in terms of wall-clock time com-513 pared to AmIGO, F<sup>2</sup>SA and F<sup>3</sup>SA. In ad-514 dition, Fig. 2a indicates that the training 515 loss of LazyBLO is smaller than those 516 of F<sup>2</sup>SA and F<sup>3</sup>SA. The superior per-517 formance of LazyBLO in comparison to 518  $F^2SA$  and  $F^3SA$  establishes the neces-519 sity of Hessian/Jacobian evaluations in 520 stochastic bilevel optimization. Without 521 them, both the convergence speed and the training loss would degrade as demon-522 strated by the experiments. Fig. 2b il-523 lustrates the convergence performance of 524

Table 3: The number of hypergradient computations and Hessian computations required by various algorithms to achieve the same training loss in data hyper-cleaning experiments (Task 1) and hyper-representation experiments (Task 2) (average over 10 repetitions).

|                               | Algorithm  | # OF HGC   | # OF HESSIAN |
|-------------------------------|--|--|--------------|
|                               | ALGORITHM<br>AMIGO<br>STOCBIO<br>BSA<br>LazyBLO<br>2 AMIGO | 42   | 252          |
| TACK 1                        | STOCBIO  | DRITHM         # OF HGC         # OF HESSIAN           IIGO         42         252           CBIO         42         252           SA         21         252           yBLO         40         5           IIGO         361         722           yBLO         640         320 | 252          |
| IASK I                        | BSA  |  | 252          |
|                               | LazyBLO  |  | 5            |
| TAOK 2                        | АмIGO  | 361  | 722          |
| TASK Z                        | LazyBLO  | 640  | 320          |
| HGC HYDERGRADIENT COMPUTATION |  |  |              |

HGC: HYPERGRADIENT COMPUTATION

LazyBLO compared to AmIGO in terms of the number of Hessian computations. Note that we do not include F<sup>2</sup>SA and F<sup>3</sup>SA in this figure since they are Hessian-free. Fig. 2b demonstrates that with the same number of Hessian evaluations, LazyBLO has a lower training loss compared to AmIGO. Furthermore, as shown in Table 3, to reach the same training loss, LazyBLO uses 320 Hessian computations, while AmIGO uses 722 Hessian computations. This significantly reduces computational costs, especially for large-scale problems.

530 531 532

### 7 CONCLUSION

In this paper, we proposed the LazyBLO algorithmic framework for solving bilevel optimization problems. Compared to existing works, LazyBLO reduces the Hessian-vector product (HVP) evaluations by updating them periodically and less frequently. Although LazyBLO uses stale HVP evaluations that introduce additional errors, our theoretical analysis demonstrated that LazyBLO not only surprisingly enjoys the same convergence rate guarantee, but also achieves a much faster wall-clock time performance. Specifically, to reach an  $\epsilon$ -stationary point, LazyBLO requires  $\mathcal{O}(\epsilon^{-2})$  for both partial gradient evaluations and HVP evaluations, which matches the SOTA non-lazy methods. We conducted experiments on multi-hyperparameter optimization tasks to verify our theoretical findings.

## 540 REFERENCES

578

579

| 542 | Ilan Adler, Zhiyue T Hu, and Tianyi Lin. New proximal newton-type methods for convex optimization. |
|-----|--|
| 543 | In 2020 59th IEEE Conference on Decision and Control (CDC), pp. 4828–4835. IEEE, 2020.             |

- Gemayqzel Bouza Allende and Georg Still. Solving bilevel programs with the kkt-approach. *Mathematical programming*, 138(1):309–332, 2013.
- G Anandalingam and DJ White. A solution method for the linear static stackelberg problem using penalty functions. *IEEE Transactions on automatic control*, 35(10):1170–1173, 1990.
- Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *arXiv preprint arXiv:1810.09502*, 2018.
- Michael Arbel and Julien Mairal. Amortized implicit differentiation for stochastic bilevel optimization.
   In International Conference on Learning Representations, 2022.
- Fan Bao, Guoqiang Wu, Chongxuan Li, Jun Zhu, and Bo Zhang. Stability and generalization of
   bilevel programming in hyperparameter optimization. *Advances in neural information processing systems*, 34:4529–4541, 2021.
- Jerome Bracken and James T McGill. Mathematical programs with optimization problems in the constraints. *Operations Research*, 21(1):37–44, 1973.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
   Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and
   Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL
   http://github.com/google/jax.
- Lesi Chen, Yaohua Ma, and Jingzhao Zhang. Near-optimal fully first-order algorithms for finding stationary points in bilevel optimization. *arXiv preprint arXiv:2306.14853*, 2023a.
- Lesi Chen, Jing Xu, and Jingzhao Zhang. On bilevel optimization without lower-level strong convexity. *arXiv preprint arXiv:2301.00712*, 2023b.
- Tianyi Chen, Yuejiao Sun, and Wotao Yin. Tighter analysis of alternating stochastic gradient method for stochastic nested problems. *arXiv preprint arXiv:2106.13781*, 2021.
- Xuxing Chen, Tesi Xiao, and Krishnakumar Balasubramanian. Optimal algorithms for stochastic
   bilevel optimization under relaxed smoothness conditions. *Journal of Machine Learning Research*, 25(151):1–51, 2024.
- 575 Mathieu Dagréou, Pierre Ablin, Samuel Vaiter, and Thomas Moreau. A framework for bilevel
  576 optimization that enables stochastic and global variance reduction algorithms. *arXiv preprint*577 *arXiv:2201.13409*, 2022.
  - Stephan Dempe and Jonathan F Bard. Bundle trust-region algorithm for bilinear bilevel programming. *Journal of Optimization Theory and Applications*, 110(2):265–288, 2001.
- Nikita Doikov, El Mahdi Chayti, and Martin Jaggi. Second-order optimization with lazy hessians.
   *arXiv preprint arXiv:2212.00781*, 2022.
- Justin Domke. Generic methods for optimization-based modeling. In *Artificial Intelligence and Statistics*, pp. 318–326. PMLR, 2012.
- Bothina El-Sobky and Yousria Abo-Elnaga. A penalty method with trust-region mechanism for
   nonlinear bilevel optimization problem. *Journal of Computational and Applied Mathematics*, 340:
   360–374, 2018.
- James E Falk and Jiming Liu. On bilevel programming, part i: general nonlinear cases. *Mathematical Programming*, 70(1):47–72, 1995.
- Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. On the convergence theory of gradient-based
   model-agnostic meta-learning algorithms. In *International Conference on Artificial Intelligence* and Statistics, pp. 1082–1092. PMLR, 2020.

594 Jinyan Fan. A shamanskii-like levenberg-marquardt method for nonlinear equations. Computational 595 Optimization and Applications, 56(1):63–80, 2013. 596 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of 597 deep networks. In International conference on machine learning, pp. 1126–1135. PMLR, 2017. 598 Luca Franceschi, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. Forward and reverse 600 gradient-based hyperparameter optimization. In International Conference on Machine Learning, 601 pp. 1165-1173. PMLR, 2017. 602 Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazzi, and Massimiliano Pontil. Bilevel 603 programming for hyperparameter optimization and meta-learning. In International Conference on 604 Machine Learning, pp. 1568–1577. PMLR, 2018. 605 Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. arXiv preprint 607 arXiv:1802.02246, 2018. 608 Stephen Gould, Basura Fernando, Anoop Cherian, Peter Anderson, Rodrigo Santa Cruz, and Edison 609 Guo. On differentiating parameterized argmin and argmax problems with application to bi-level 610 optimization. arXiv preprint arXiv:1607.05447, 2016. 611 612 Riccardo Grazzi, Luca Franceschi, Massimiliano Pontil, and Saverio Salzo. On the iteration com-613 plexity of hypergradient computation. In International Conference on Machine Learning, pp. 614 3748-3758. PMLR, 2020. 615 Bin Gu, Guodong Liu, Yanfu Zhang, Xiang Geng, and Heng Huang. Optimizing large-scale 616 hyperparameters via automated learning algorithm. arXiv preprint arXiv:2102.09026, 2021. 617 618 Zhishuai Guo, Quanqi Hu, Lijun Zhang, and Tianbao Yang. Randomized stochastic variance-reduced 619 methods for multi-task stochastic bilevel optimization. arXiv preprint arXiv:2105.02266, 2021. 620 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image 621 recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, 622 pp. 770-778, 2016. 623 624 Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A two-timescale framework for bilevel 625 optimization: Complexity analysis and application to actor-critic. arXiv preprint arXiv:2007.05170, 626 2020. 627 Yibo Hu, Xiang Wu, and Ran He. Tf-nas: Rethinking three search freedoms of latency-constrained 628 differentiable neural architecture search. In European Conference on Computer Vision, pp. 123–139. 629 Springer, 2020. 630 631 Minhui Huang, Xuxing Chen, Kaiyi Ji, Shiqian Ma, and Lifeng Lai. Efficiently escaping saddle 632 points in bilevel optimization. arXiv preprint arXiv:2202.03684, 2022. 633 Kaiyi Ji and Yingbin Liang. Lower bounds and accelerated algorithms for bilevel optimization. arXiv 634 preprint arXiv:2102.03926, 2021. 635 636 Kaiyi Ji, Jason D Lee, Yingbin Liang, and H Vincent Poor. Convergence of meta-learning with 637 task-specific adaptation over partial parameters. Advances in Neural Information Processing 638 Systems, 33:11490-11500, 2020. 639 Kaiyi Ji, Junjie Yang, and Yingbin Liang. Bilevel optimization: Convergence analysis and enhanced 640 design. In International Conference on Machine Learning, pp. 4882-4892. PMLR, 2021. 641 642 Prashant Khanduri, Siliang Zeng, Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A 643 momentum-assisted single-timescale stochastic approximation algorithm for bilevel optimization. 644 *arXiv preprint arXiv:2102.07367v1*, 2021a. 645 Prashant Khanduri, Siliang Zeng, Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A 646 near-optimal algorithm for stochastic bilevel optimization via double-momentum. Advances in 647 Neural Information Processing Systems, 34:30271-30283, 2021b.

682

684

689

690

- 648 Vijay Konda and John Tsitsiklis. Actor-critic algorithms. Advances in neural information processing 649 systems, 12, 1999. 650
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 651
- 652 Jeongyeol Kwon, Dohyun Kwon, Stephen Wright, and Robert D Nowak. A fully first-order method 653 for stochastic bilevel optimization. In International Conference on Machine Learning, pp. 18083-654 18113. PMLR, 2023. 655
- Francesco Lampariello and Marco Sciandrone. Global convergence technique for the newton method 656 with periodic hessian evaluation. Journal of optimization theory and applications, 111:341–358, 657 2001. 658
- 659 Bruno Lecouat, Jean Ponce, and Julien Mairal. Designing and learning trainable priors with non-660 cooperative games. 2020a. 661
- 662 Bruno Lecouat, Jean Ponce, and Julien Mairal. A flexible framework for designing trainable priors with adaptive smoothing and game encoding. Advances in Neural Information Processing Systems, 663 33:15664-15675, 2020b. 664
- 665 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to 666 document recognition. Proceedings of the IEEE, 86(11):2278-2324, 1998. 667
- Jiaxiang Li, Siliang Zeng, Hoi To Wai, Chenliang Li, Alfredo Garcia, and Mingyi Hong. Getting 668 more juice out of the SFT data: Reward learning from human demonstration improves SFT for 669 LLM alignment. In ICML 2024 Workshop on Theoretical Foundations of Foundation Models, 670 2024. 671
- 672 Junyi Li, Bin Gu, and Heng Huang. A fully single loop algorithm for bilevel optimization without 673 hessian inverse. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 36, pp. 674 7426-7434, 2022.
- Yi Li, Lingxiao Song, Xiang Wu, Ran He, and Tieniu Tan. Learning a bi-level adversarial network 676 with global and local perception for makeup-invariant face verification. Pattern Recognition, 90: 677 99-108, 2019. 678
- 679 Dongze Lian, Yin Zheng, Yintao Xu, Yanxiong Lu, Leyu Lin, Peilin Zhao, Junzhou Huang, and 680 Shenghua Gao. Towards fast adaptation of neural architectures with meta learning. In International 681 Conference on Learning Representations, 2019.
- Renjie Liao, Yuwen Xiong, Ethan Fetaya, Lisa Zhang, KiJung Yoon, Xaq Pitkow, Raquel Urtasun, and 683 Richard Zemel. Reviving and improving recurrent back-propagation. In International Conference on Machine Learning, pp. 3082–3091. PMLR, 2018. 685
- 686 Bo Liu, Mao Ye, Stephen Wright, Peter Stone, and Qiang Liu. Bome! bilevel optimization made 687 easy: A simple first-order approach. Advances in Neural Information Processing Systems, 35: 17248–17262, 2022a. 688
  - Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. arXiv preprint arXiv:1806.09055, 2018.
- 692 Risheng Liu, Xuan Liu, Xiaoming Yuan, Shangzhi Zeng, and Jin Zhang. A value-function-based 693 interior-point method for non-convex bi-level optimization. In International Conference on Machine Learning, pp. 6882–6892. PMLR, 2021. 694
- Zhuqing Liu, Xin Zhang, Prashant Khanduri, Songtao Lu, and Jia Liu. Interact: achieving low sample 696 and communication complexities in decentralized bilevel learning over networks. In Proceedings 697 of the Twenty-Third International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing, pp. 61–70, 2022b. 699
- Jonathan Lorraine, Paul Vicol, and David Duvenaud. Optimizing millions of hyperparameters by 700 implicit differentiation. In International Conference on Artificial Intelligence and Statistics, pp. 1540-1552. PMLR, 2020.

| 702<br>703<br>704<br>705 | Matthew MacKay, Paul Vicol, Jon Lorraine, David Duvenaud, and Roger Grosse. Self-tuning networks: Bilevel optimization of hyperparameters using structured best-response functions. <i>arXiv</i> preprint arXiv:1903.03088, 2019.  |
|--------------------------|--|
| 706<br>707<br>708        | Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimization through reversible learning. In <i>International conference on machine learning</i> , pp. 2113–2122. PMLR, 2015.  |
| 709<br>710<br>711        | Julien Mairal, Francis Bach, and Jean Ponce. Task-driven dictionary learning. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 34(4):791–804, 2011.   |
| 712<br>713               | Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. <i>arXiv</i> preprint arXiv:1803.02999, 2018.  |
| 714<br>715<br>716        | Rui Pan, Jipeng Zhang, Xingyuan Pan, Renjie Pi, Xiaoyu Wang, and Tong Zhang. Scalebio: Scalable bilevel optimization for llm data reweighting. <i>arXiv preprint arXiv:2406.19976</i> , 2024.  |
| 717<br>718               | Barak A Pearlmutter. Fast exact multiplication by the hessian. <i>Neural computation</i> , 6(1):147–160, 1994.   |
| 719<br>720<br>721        | Fabian Pedregosa. Hyperparameter optimization with approximate gradient. In <i>International conference on machine learning</i> , pp. 737–746. PMLR, 2016.   |
| 722<br>723<br>724        | Peiwen Qiu, Yining Li, Zhuqing Liu, Prashant Khanduri, Jia Liu, Ness B Shroff, Elizabeth Ser-<br>ena Bentley, and Kurt Turck. Diamond: Taming sample and communication complexities in<br>decentralized bilevel optimization. <i>arXiv preprint arXiv:2212.02376</i> , 2022. |
| 725<br>726               | Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. <i>Advances in neural information processing systems</i> , 32, 2019.  |
| 728                      | Walter Rudin et al. Principles of mathematical analysis, volume 3. McGraw-hill New York, 1976.   |
| 729<br>730<br>731        | Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. Truncated back-propagation for bilevel optimization. In <i>The 22nd International Conference on Artificial Intelligence and Statistics</i> , pp. 1723–1732. PMLR, 2019.                                      |
| 732<br>733<br>734        | VE Shamanskii. A modification of newton's method. Ukrainian Mathematical Journal, 19(1): 118–122, 1967.  |
| 735<br>736<br>737        | Han Shen and Tianyi Chen. On penalty-based bilevel gradient descent method. <i>arXiv preprint arXiv:2302.05185</i> , 2023.   |
| 738<br>739<br>740        | Ankur Sinha, Samish Bedi, and Kalyanmoy Deb. Bilevel optimization based on kriging approxima-<br>tions of lower level optimal value function. In 2018 IEEE congress on evolutionary computation<br>(CEC), pp. 1–8. IEEE, 2018.   |
| 741<br>742               | Ankur Sinha, Tharo Soun, and Kalyanmoy Deb. Using karush-kuhn-tucker proximity measure for solving bilevel optimization problems. <i>Swarm and evolutionary computation</i> , 44:496–510, 2019.  |
| 744<br>745               | Xingyou Song, Wenbo Gao, Yuxiang Yang, Krzysztof Choromanski, Aldo Pacchiano, and Yunhao Tang. Es-maml: Simple hessian-free meta learning. <i>arXiv preprint arXiv:1910.01215</i> , 2019.  |
| 746<br>747<br>748        | Daouda Sow, Kaiyi Ji, and Yingbin Liang. On the convergence theory for hessian-free bilevel algorithms. <i>Advances in Neural Information Processing Systems</i> , 35:4136–4149, 2022.   |
| 749<br>750<br>751        | Yuesong Tian, Li Shen, Guinan Su, Zhifeng Li, and Wei Liu. Alphagan: Fully differentiable architecture search for generative adversarial networks. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 44(10):6752–6766, 2021.                           |
| 752<br>753               | Luis Vicente, Gilles Savard, and Joaquim Júdice. Descent approaches for quadratic bilevel program-<br>ming. <i>Journal of Optimization theory and applications</i> , 81(2):379–399, 1994.  |
| 755                      | Zhongping Wan, Lijun Mao, and Guangmin Wang. Estimation of distribution algorithm for a class of nonlinear bilevel programming problems. <i>Information Sciences</i> , 256:184–196, 2014.  |

| 756<br>757<br>759 | Chang-yu Wang, Yuan-yuan Chen, and Shou-qiang Du. Further insight into the shamanskii modifica-<br>tion of newton method. <i>Applied mathematics and computation</i> , 180(1):46–52, 2006.  |
|-------------------|---|
| 759<br>760        | Douglas J White and G Anandalingam. A penalty function approach for solving bi-level linear programs. <i>Journal of Global Optimization</i> , 3(4):397–419, 1993.   |
| 761<br>762<br>763 | Feijie Wu, Zitao Li, Yaliang Li, Bolin Ding, and Jing Gao. Fedbiot: Llm local fine-tuning in federated learning without full model. In <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pp. 3345–3355, 2024.                         |
| 765<br>766        | Junjie Yang, Kaiyi Ji, and Yingbin Liang. Provably faster algorithms for bilevel optimization.<br>Advances in Neural Information Processing Systems, 34:13670–13682, 2021.  |
| 767<br>768        | Yifan Yang, Peiyao Xiao, and Kaiyi Ji. Achieving $\mathcal{O}(\epsilon^{-1.5})$ complexity in hessian/jacobian-free stochastic bilevel optimization. Advances in Neural Information Processing Systems, 2023.   |
| 770<br>771<br>772 | Alain B Zemkoho and Shenglong Zhou. Theoretical and numerical comparison of the karush–kuhn–<br>tucker and value function reformulations in bilevel optimization. <i>Computational Optimization and</i><br><i>Applications</i> , 78(2):625–674, 2021.                             |
| 773<br>774<br>775 | Yihua Zhang, Guanhua Zhang, Prashant Khanduri, Mingyi Hong, Shiyu Chang, and Sijia Liu.<br>Revisiting and advancing fast adversarial training through the lens of bi-level optimization. In<br><i>International Conference on Machine Learning</i> , pp. 26693–26712. PMLR, 2022. |
| 776               |   |
| 777               |   |
| 778               |   |
| 779               |   |
| 700               |   |
| 782               |   |
| 783               |   |
| 784               |   |
| 785               |   |
| 786               |   |
| 787               |   |
| 788               |   |
| 789               |   |
| 790               |   |
| 791               |   |
| 792               |   |
| 793               |   |
| 794               |   |
| 795               |   |
| 796               |   |
| 797               |   |
| 798               |   |
| 799               |   |
| 008               |   |
| 900               |   |
| 802               |   |
| 804               |   |
| 805               |   |
| 806               |   |
| 807               |   |
| 808               |   |
| 809               |   |

# 810 A ADDITIONAL RELATED WORK

812 Bilevel Optimization: The history of bilevel optimization dates back to 1973 (Bracken & McGill, 813 1973). Some early attempts for solving bilevel problems include: value function (Liu et al., 2021; 814 Sinha et al., 2018; Zemkoho & Zhou, 2021), Karush-Kuhn-Tucker conditions based reformula-815 tions (Allende & Still, 2013; Sinha et al., 2019; Zemkoho & Zhou, 2021), penalty function (White & Anandalingam, 1993; Anandalingam & White, 1990; Wan et al., 2014), approximate descent (Falk 816 & Liu, 1995; Vicente et al., 1994), and trust region methods (Dempe & Bard, 2001; El-Sobky & 817 Abo-Elnaga, 2018). Among these approaches, approximate descent methods have gained promi-818 nence recently because of their ease of implementation as well as strong theoretical and empirical 819 performance in many machine learning applications. Two standard descent-based approaches to 820 tackle problems of form (1) are iterative differentiation (ITD) (Domke, 2012; Maclaurin et al., 2015; 821 Franceschi et al., 2017; 2018; Shaban et al., 2019; Grazzi et al., 2020; MacKay et al., 2019) and ap-822 proximate implicit differentiation (AID) (Domke, 2012; Pedregosa, 2016; Liao et al., 2018; Ghadimi 823 & Wang, 2018; Grazzi et al., 2020; Lorraine et al., 2020; Gould et al., 2016; Ji & Liang, 2021; 824 MacKay et al., 2019; Khanduri et al., 2021a; Hong et al., 2020). The basic idea of ITD is to obtain 825 an approximate hypergradient of the loss function  $\ell(\mathbf{x})$  in Eq. (1) by differentiating the unrolled 826 iterates of the LL problem. Consequently, ITD-based approaches need to store all the LL iterates in 827 the memory (Shaban et al., 2019). On the other hand, AID relies on the implicit function theorem to compute the implicit gradient of  $\ell(\mathbf{x})$  without the need to maintain the sequence of LL iterates. 828 Instead of differentiating the iterates of the LL problem, AID computes the implicit gradient by 829 approximately solving a linear system of equations using HVPs. In this work, we focus on AID-based 830 approaches for solving stochastic bilevel problems. 831

## **B** ADDITIONAL EXPERIMENTAL DETAILS AND RESULTS

#### B.1 SPECIFICATIONS OF THE BASELINE ALGORITHMS IN SECTION 6

In this section, we provide more description of the baseline algorithms used in our experiments, as follows:

- AmIGO (Arbel & Mairal, 2022): a double-loop stochastic AID-based bilevel algorithm that uses conjugate gradient to estimate the Hessian inverse.
- stocBiO (Ji et al., 2021): a two timescale stochastic AID-based bilevel approach that uses Neumann Series to estimate the Hessian inverse. The repository of stocBiO is available at https://github.com/JunjieYang97/StocBio.
- BSA (Ghadimi & Wang, 2018): an AID-based bilevel method that uses single-sample sampling.
- MRBO (Yang et al., 2021): a single-loop AID-based stochastic bilevel algorithm that uses momentum-based SGD to accelerate convergence. The implementation of MRBO is available at https://github.com/JunjieYang97/MRVRBO.
- F<sup>2</sup>SA (Kwon et al., 2023): a fully first-order (Hessian/Jacobian-free) stochastic bilevel method, which is double-loop.
- F<sup>3</sup>SA (Kwon et al., 2023): a fully first-order stochastic bilevel approach that uses momentum-based SGD to accelerate convergence and is single timescale.

#### B.2 EXPERIMENTAL DETAILS FOR DATA HYPER-CLEANING

In this section, we describe the details of the experiments on data hyper-cleaning. The goal of data hyper-cleaning is to train a classifier on a potentially corrupt dataset. To make fair comparison, we follow the same implementation as in (Ji et al., 2021; Yang et al., 2021) and apply it to other algorithms. The objective function can be written as follows:

$$\min_{\lambda} \mathcal{L}_{\mathcal{D}_{val}}\left(\lambda, w^*\right) = \frac{1}{|\mathcal{D}_{val}|} \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_{val}} \mathcal{L}\left(w^* \mathbf{x}_i, \mathbf{y}_i\right)$$

861 862 863

832 833

834 835

836

837

838 839

840

841

842

843

844

845

846

847

848

849

850

851

852 853

854 855

856

858

s.t. 
$$w^* = \operatorname*{arg\,min}_{w} \left( \frac{1}{|\mathcal{D}_{tr}|} \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_{tr}} \sigma\left(\lambda_i\right) \mathcal{L}\left(w\mathbf{x}_i, \mathbf{y}_i\right) + C_r \|w\|^2 \right),$$

where  $(\mathbf{x}_i, \mathbf{y}_i)$  represents the data samples,  $\mathcal{D}_{val}$  and  $\mathcal{D}_{tr}$  correspond to the validation data and the training data,  $\mathcal{L}$  denotes the cross-entropy loss,  $\sigma$  represents the sigmoid function, and  $C_r$  is the regularization parameter. Note that the training loss corresponds to the upper-level loss. We choose  $C_r = 0.001$  in our experiments, which is the same as (Shaban et al., 2019; Ji et al., 2021). We conduct experiments on the MNIST dataset (LeCun et al., 1998), which is corrupted by replacing the training data label with a uniformly random one. Such replacement has a probability p, referred to as the corruption rate. We run the experiments with corruption rates of  $p = \{0.1, 0.15, 0.2, 0.25, 0.3\}$ .

871 We compare the performance of LazyBLO with AmIGO (Arbel & Mairal, 2022), stocBiO (Ji et al., 872 2021), BSA (Ghadimi & Wang, 2018), and MRBO (Yang et al., 2021). For all algorithms, we tune the 873 parameters using grid search to achieve the best convergence performance based on the training loss 874 as the metric. As a result, we set the batch size to 1000 for AmIGO, stocBiO, MRBO and LazyBLO. We set both the outer stepsize  $\alpha$  and the inner stepsize  $\beta$  as 0.1, and the Hessian update stepsize  $\gamma$  as 875 0.5 for AmIGO, stocBiO and MRBO. We choose both the outer stepsize  $\alpha$  and the inner stepsize  $\beta$  to 876 be 0.01, and the Hessian update stepsize  $\gamma$  to be 0.1 for BSA. For LazyBLO, We set 0.5 as the inner 877 stepsize  $\beta$ , and 0.1 as both the outer stepsize  $\alpha$  and the Hessian update stepsize  $\gamma$ . We set the number 878 of inner-loop iterations for y-update to 64 for AmIGO, stocBiO and BSA. We choose the number of 879 iterations for Hessian inverse evaluations to be 6 for AmIGO and stocBiO, and 12 for MRBO and 880 BSA. For LazyBLO, we set 8 as the inner-loop iteration number N for x- and y-update. We conduct 10 repetitions for the experiments using different random seeds. The solid line shows the average 882 training loss, and the shaded area represents the variance containing the maximum and the minimum 883 values. We run the data hyper-cleaning experiments using NVIDIA GeForce RTX 3060 GPU. 884

#### B.3 EXPERIMENTAL DETAILS FOR DEEP HYPER-REPRESENTATION

In this section, we show the details of the experiments on deep hyper-representation, which aims to classify the images. The objective function is given by:

$$\min_{\lambda} \mathcal{L}_{\mathcal{D}_{val}}(\lambda, w^*) = \frac{1}{|\mathcal{D}_{val}|} \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_{val}} \mathcal{L}\left(w^* f\left(\lambda; \mathbf{x}_i\right), \mathbf{y}_i\right)$$
  
s.t.  $w^* = \operatorname*{arg\,min}_{w} \frac{1}{|\mathcal{D}_{tr}|} \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_{tr}} \mathcal{L}\left(wf\left(\lambda, \mathbf{x}_i\right), \mathbf{y}_i\right),$ 

885

886 887

889

890

893

where  $(\mathbf{x}_i, \mathbf{y}_i)$  denotes the data samples,  $\mathcal{D}_{val}$  and  $\mathcal{D}_{tr}$  are the validation data and the training data,  $\mathcal{L}$ corresponds to the cross-entropy loss,  $f(\lambda; \mathbf{x}_i)$  represents the features extracted from the data sample. We run the experiments with ResNet-20 network (He et al., 2016) on CIFAT-10 dataset (Krizhevsky et al., 2009) using a batch size of 128. We treat the last two layers in ResNet-20 as the LL parameters w with a dimension of 5, 130, and all remaining layers as the UL parameters  $\lambda$  with a dimension of 11, 168, 832.

900 We compare LazyBLO with AmIGO (Arbel & Mairal, 2022), F<sup>2</sup>SA (Kwon et al., 2023) and 901  $F^{3}SA$  (Kwon et al., 2023). To ensure the best performance of all the algorithms, we fine tune 902 the parameters using grid search with the goal of finding the lowest training loss. Consequently, for 903 AmIGO, we set all the stepsize for updating x, y and z to 0.01. We choose the number of y-update 904 iterations to be 8 and the number of z-update iterations to be 2. For LazyBLO, we choose the stepsize 905  $\alpha$  and  $\gamma$  to be 0.01, and  $\beta$  to be 0.05. We set 2 as the inner-loop iteration number N. Following the same notations as in (Kwon et al., 2023), for F<sup>2</sup>SA, we choose the stepsize  $\alpha$  as 0.1 and  $\gamma$  as 906 0.01. We set both the step-size ratio  $\xi$  and the Lagrangian multiplier  $\lambda$  to 0.5. We choose the number 907 of inner-loop iterations to be 1. For F<sup>3</sup>SA, we set 0.05 as  $\alpha$ , 0.01 as  $\gamma$ , 0.1 as  $\xi$ , 0.5 as  $\lambda$ , and 0.9 908 as momentum-weight  $\eta$ . We repeat the experiments 10 times with different random seeds, where 909 the solid line represents the average training loss or test accuracy, and the shaded area shows the 910 variance containing the maximum and the minimum values. We run the deep hyper-representation 911 experiments using NVIDIA Tesla V100 GPU. 912

913

## 914 B.4 ADDITIONAL EXPERIMENT RESULTS

#### 915 B.4.1 DATA HYPER-CLEANING 916

917 We can see in Table 4 that the test accuracy of LazyBLO is *comparable* to the SOTA baseline algorithms although LazyBLO uses stale Hessian information. In addition, the number of Hessian

| 920 |                               |                |             |               |               |           |
|-----|-------------------------------|----------------|-------------|---------------|---------------|-----------|
| 921 |                               | LAZYBLO        | АмIGO       | STOCBIO       | MRBO          | BSA       |
| 922 | TEST ACCURACY (%)             | 72.31          | 72.12       | 72.75         | 69.46         | 72.92     |
| 923 | # OF HESSIAN                  | 6              | 60          | 60            | 1440          | 720       |
| 924 |                               |                |             |               |               |           |
| 925 |                               | 20.0           | — LazyBL    | O p=0.10      |               |           |
| 926 |                               | ر 17.5 ·       | — LazyBL    | O p=0.15      |               |           |
| 927 |                               | 8<br>15.0      | LazyBL      | O p=0.20      |               |           |
| 928 |                               | b 12.5         | — LazyBL    | O p=0.23      |               |           |
| 929 |                               | LE 10.01       |             |               |               |           |
| 930 |                               | 5.0            |             | nam           |               |           |
| 931 |                               | 2.5            |             |               |               |           |
| 932 |                               | 0 20 40        | 0 60 80 10  | $0\ 120\ 140$ |               |           |
| 933 |                               | ĸ              | unning time | (5)           |               |           |
| 934 | Eisen 4. Companies of Long PL | O on data huma | m alaamina. | MNIGT 4       | ata a at with | different |

Table 4: Convergence performance of different bilevel algorithms on data hyper-cleaning on MNIST dataset (corruption rate p = 0.1, average over 10 repetitions).

Figure 4: Comparison of LazyBLO on data hyper-cleaning on MNIST dataset with different corruption rates (p).

computations required for LazyBLO to converge is significantly reduced, which is ten times fewer
 than AmIGO and stocBiO, 240 times fewer than MRBO, and 120 times fewer than BSA.

Figure 4 illustrates the robustness of LazyBLO against corrupted datasets. We can see from Figure 4
that when the corruption rate *p* (the probability that a training data label is replaced by a uniformly
random one) is larger, the training loss becomes higher, which is natural since with larger corruption
rate the classification problem becomes challenging. However, the convergence speed of LazyBLO is
similar regardless of the corruption rate *p*.

#### B.4.2 DEEP HYPER-REPRESENTATION



Figure 5: Test accuracy of different bilevel algorithms on deep hyper-representation on CIFAR-10 dataset (10 repetitions).

Figure 5a illustrates the test accuracy of LazyBLO compared with the baseline algorithms, and it demonstrates that LazyBLO converges faster in terms of wall-clock time compared to both  $F^2SA$ and  $F^3SA$ . Figure 5b shows the test accuracy of LazyBLO compared to AmIGO in terms of the number of Hessian computations, and it indicates that with the same number of Hessian evaluations, LazyBLO has a higher test accuracy compared to AmIGO.

## C PROOF OF THEOREM 5.3: NON-CONVEX $\ell(\mathbf{x})$

- - C.1 PROOF SKETCHES
- 971 Here, we provide a detailed proof sketch of Theorem 5.3. The detailed proof is provided in Appendix C.2. The proof is organized into five key steps:

Step 1) Descent in the upper-level objective function: First, we show the bound for the per-iterate descent of the UL problem as follows: 

**Lemma C.1.** Under Assumptions 3.1–3.4, the following inequality holds for successive iterations of Algorithm 1: 

$$\mathbb{E}\left[\ell\left(\mathbf{x}_{t}^{n+1}\right)-\ell\left(\mathbf{x}_{t}^{n}\right)\right] \leq -\frac{\alpha_{t}}{2}\mathbb{E}\left[\left\|\nabla\ell\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right]-\left(\frac{\alpha_{t}}{2}-\frac{\alpha_{t}^{2}L_{l}}{2}\right)\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right]+4\tilde{\sigma}_{g_{xy}}^{2}\frac{B_{fy}^{2}}{\mu_{g}^{2}}\alpha_{t}$$

980 
$$+8L_{f}^{2}L_{z}^{2}\alpha_{t}^{3}N\sum_{i=0}\mathbb{E}\left[\left\|h_{t,i}^{f}\right\|^{2}\right]+2\alpha_{t}L_{f}^{2}\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right]+2\sigma_{fx}^{2}\alpha_{t}+16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N^{2}\sigma_{gy}^{2}\alpha_{t}$$
982 
$$N-1$$

$$+ \left(4\sigma_{g_{xy}}^{2}\alpha_{t} + 4L_{f}^{2}\alpha_{t}\right)\mathbb{E}\left[\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2}\right] + 16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\alpha_{t}\sum_{i=0}^{N-1}\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2}\right]$$

for all  $t \in \{0, 1, \dots, T-1\}$  and  $n \in \{0, 1, \dots, N-1\}$ , where the expectation is taken over the stochasticity of the algorithm.

Lemma C.1 indicates that the descent in the upper-level objective function depends on i) the stochastic gradient estimator  $\mathbb{E}[\|h_{t,n}^f\|^2]$ , ii) the full gradient  $\mathbb{E}[\|\nabla_{\mathbf{y}}g(\mathbf{x}_t^n,\mathbf{y}_t^n)\|^2]$ , iii) the approximation error of  $\mathbf{y}^*(\mathbf{x})$ , which is  $\mathbb{E}[\|\mathbf{y}_t^n - \mathbf{y}^*(\mathbf{x}_t^n)\|^2]$  and will be bounded in Step 2), and iv) the approximation gap of  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$ , which is  $\mathbb{E}[\|\mathbf{z}_t - \mathbf{z}_t^*\|^2]$  and will be bounded in Step 3).

Step 2) Descent in the error of  $y^*(x)$ : We bound the approximation error of  $y^*(x)$  as follows:

**Lemma C.2.** Under Assumptions 3.2–3.4, the approximation error of  $\mathbf{y}^*(\mathbf{x})$  for Algorithm 1 satisfies the following inequality:

$$\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n+1} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n+1}\right)\right\|^{2}\right] \leq (1+c_{1})\left(1+c_{2}\right)\left(1 - \frac{2\beta_{t}\mu_{g}L_{g}}{\mu_{g}+L_{g}}\right)\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] \\ + \left(1 + \frac{1}{c_{1}}\right)L_{y}^{2}\alpha_{t}^{2}\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] + (1+c_{1})\left(1+c_{2}\right)\left(\beta_{t}^{2} - \frac{2\beta_{t}}{\mu_{g}+L_{g}}\right)\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)\right\|^{2}\right]$$

 $+\left(1+c_1\right)\left(1+\frac{1}{c_2}\right)\beta_t^2\sigma_{g_y}^2,$ 

for all  $t \in \{0, 1, ..., T-1\}$  and  $n \in \{0, 1, ..., N-1\}$  with some constants  $c_1, c_2 > 0$ , where the expectation is taken over the randomness of the algorithm. 

Lemma C.2 shows that the approximation error of  $\mathbf{y}^{*}(\mathbf{x})$  is affected by the full gradient  $\mathbb{E}[\|\nabla_{\mathbf{y}}g(\mathbf{x}_t^n,\mathbf{y}_t^n)\|^2]$ , and the stochastic gradient estimator  $\mathbb{E}[\|h_{t,n}^f\|^2]$ , which is due to the cou-pled structure of the bilevel optimization problem. 

Step 3) Descent in the error of  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$ : Next, we demonstrate that the approximation error of  $\mathbf{z}^{*}(\mathbf{x}, \mathbf{y})$  can be bounded as follows:

**Lemma C.3.** Under Assumptions 3.1–3.4, the following inequality of the approximation error of  $\mathbf{z}^{*}(\mathbf{x}, \mathbf{y})$  holds for Algorithm 1: 

$$\mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t+1}^{*}\right\|^{2}\right] \leq (1 + c_{3})\left(1 + c_{4}\right)\left(\gamma_{t}^{2} - \frac{2\gamma_{t}}{\mu_{g} + L_{q}}\right)\mathbb{E}\left[\left\|\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right)\right\|^{2}\right]$$

$$+ \left( (1+c_3) \left(1+c_4\right) \left(1 - \frac{2\gamma_t \mu_g L_q}{\mu_g + L_q}\right) + 4\sigma_{g_{yy}}^2 \gamma_t^2 \left(1+c_3\right) \left(1 + \frac{1}{c_4}\right) \right) \mathbb{E} \left[ \|\mathbf{z}_t - \mathbf{z}_t^*\|^2 \right]$$

$$+ \left( (1+c_3) \left(1+c_4\right) \left(1 - \frac{2\gamma_t \mu_g L_q}{\mu_g + L_q}\right) + 4\sigma_{g_{yy}}^2 \gamma_t^2 \left(1+c_3\right) \left(1 + \frac{1}{c_4}\right) \right) \mathbb{E} \left[ \|\mathbf{z}_t - \mathbf{z}_t^*\|^2 \right]$$

$$+ 2\left(1 + \frac{1}{c_3}\right)L_z^2\alpha_t^2N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|h_{t,n}^f\right\|^2\right] + 4\left(1 + \frac{1}{c_3}\right)L_z^2\beta_t^2N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_t^n, \mathbf{y}_t^n\right)\right\|^2\right]$$

 $+2\sigma_{f_y}^2\left(1+c_3\right)\left(1+\frac{1}{c_4}\right)\gamma_t^2+4\sigma_{g_{yy}}^2\frac{B_{f_y}^2}{\mu_a^2}\left(1+c_3\right)\left(1+\frac{1}{c_4}\right)\gamma_t^2+4\left(1+\frac{1}{c_3}\right)L_z^2\beta_t^2N^2\sigma_{g_y}^2,$ 

for all  $t \in \{0, 1, ..., T-1\}$  and  $n \in \{0, 1, ..., N-1\}$  with some constants  $c_3, c_4 > 0$ , where  $\mathbf{z}_t = \mathbf{z} \left( \mathbf{x}_t^0, \mathbf{y}_t^0 \right)$  and  $\mathbf{z}_t^* = \mathbf{z}^* \left( \mathbf{x}_t^0, \mathbf{y}_t^0 \right)$ . The expectation is taken over the stochasticity of the algorithm.

Lemma C.3 shows that the approximation error of  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$  is influenced by the full gradients  $\mathbb{E}[\|\nabla_{\mathbf{y}}g(\mathbf{x}_t^n, \mathbf{y}_t^n)\|^2]$  and  $\mathbb{E}[\|\nabla_{\mathbf{z}}q(\mathbf{x}_t^N, \mathbf{y}_t^N, \mathbf{z}_t)\|^2]$ , and the stochastic gradient estimator  $\mathbb{E}[\|h_{t,n}^f\|^2]$ , which is due to the coupled structure of the quadratic problem in (3).

Step 4) Descent in the potential function: We define the potential function  $W_t$  as follows:

$$W_{t} = \ell \left( \mathbf{x}_{t}^{0} \right) + K_{y} \left\| \mathbf{y}_{t}^{0} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{0} \right) \right\|^{2} + K_{z} \left\| \mathbf{z}_{t} \left( \mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0} \right) - \mathbf{z}^{*} \left( \mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0} \right) \right\|^{2}$$

<sup>1034</sup> To demonstrate the descent in the potential function, we prove the following lemma.

**Lemma C.4.** Set  $c_1 = \frac{\beta_t L_{\mu_g}}{2(1-\beta_t L_{\mu_g})}$ ,  $c_2 = \frac{\beta_t L_{\mu_g}}{1-2\beta_t L_{\mu_g}}$ ,  $c_3 = \frac{\gamma_t L_{\mu_q}}{2(1-\gamma_t L_{\mu_q})}$ , and  $c_4 = \frac{\gamma_t L_{\mu_q}}{1-2\gamma_t L_{\mu_q}}$ . Under the same conditions as described in Theorem 5.3 and using Lemmas C.1-C.3, the iterates generated by Algorithm 1 satisfies: for all  $t \in \{0, 1, ..., T-1\}$ ,

$$\mathbb{E}\left[W_{t+1} - W_t\right] \leq -\frac{\alpha_t}{2} \sum_{n=0}^{N-1} \mathbb{E}\left[\left\|\nabla \ell\left(\mathbf{x}_t^n\right)\right\|^2\right] + \sigma_{g_{xy}}^2 \alpha_t C_{g_{xy}} N + \sigma_{f_x}^2 \alpha_t C_{f_x} N + \sigma_{g_{yy}}^2 \alpha_t C_{g_{yy}} N + \sigma_{f_y}^2 \alpha_t C_{f_y} N + \sigma_{g_y}^2 \alpha_t \left(C_{g_1} N + C_{g_2} \frac{1}{N}\right)$$

1041 1042

1039 1040

1030

1031 1032 1033

1043 1044

1047

where the constant values  $C_{g_{xy}}$ ,  $C_{f_x}$ ,  $C_{g_1}$ ,  $C_{g_2}$ ,  $C_{g_{yy}}$  and  $C_{f_y}$ , which are independent of N, are defined in (18) of Appendix C.

With the proper parameter choices, the coefficients of  $\mathbb{E}[\|\mathbf{y}_t^n - \mathbf{y}^*(\mathbf{x}_t^n)\|^2]$ ,  $\mathbb{E}[\|\mathbf{z}_t - \mathbf{z}_t^*\|^2]$ ,  $\mathbb{E}[\|\nabla_{\mathbf{y}}g(\mathbf{x}_t^n, \mathbf{y}_t^n)\|^2]$ ,  $\mathbb{E}[\|h_{t,n}^f\|^2]$  and  $\mathbb{E}[\|\nabla_{\mathbf{z}}q(\mathbf{x}_t^N, \mathbf{y}_t^N, \mathbf{z}_t)\|^2]$  are made to be non-positive within the ranges of  $\alpha_t$ ,  $\beta_t$  and  $\gamma_t$ .

1052 Step 5) Proof of Theorem 5.3: Choose a constant step-size  $\alpha_t = \alpha$ . Under the same conditions as 1053 described in Theorem 5.3, telescoping the result in Lemma C.4 from 0 to T - 1 yields:

$$\frac{1}{TN} \sum_{t=0}^{T-1} \sum_{n=0}^{N-1} \mathbb{E}\left[ \left\| \nabla \ell \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} \right] \leq \frac{2 \left( W_{0} - \ell^{*} \right)}{\alpha NT} + 2 \left( \sigma_{f_{x}}^{2} C_{f_{x}} + \sigma_{f_{y}}^{2} C_{f_{y}} + \sigma_{g_{y}}^{2} \left( C_{g_{1}} + C_{g_{2}} \frac{1}{N^{2}} \right) + \sigma_{g_{yy}}^{2} C_{g_{yy}} + \sigma_{g_{xy}}^{2} C_{g_{xy}} \right),$$

1058 1059

1061

1062

1055 1056 1057

where  $W_0 = \ell \left( \mathbf{x}_0^0 \right) + K_y \left\| \mathbf{y}_0^0 - \mathbf{y}^* \left( x_0^0 \right) \right\|^2 + K_z \left\| \mathbf{z}_0 - \mathbf{z}^* \left( \mathbf{x}_0^0, \mathbf{y}_0^0 \right) \right\|^2$ . The proof of Theorem 5.3 is completed.

C.2 DETAILED PROOF

1065 C.2.1 DESCENT IN THE UPPER-LEVEL OBJECTIVE FUNCTION

**Lemma C.5.** Under Assumptions 3.1–3.4, the following inequality holds for successive iterations of Algorithm 1:

$$\mathbb{E}\left[\ell\left(\mathbf{x}_{t}^{n+1}
ight)-\ell\left(\mathbf{x}_{t}^{n}
ight)
ight]$$

$$\leq -\frac{\alpha_{t}}{2}\mathbb{E}\left[\left\|\nabla \ell\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] - \left(\frac{\alpha_{t}}{2} - \frac{\alpha_{t}^{2}L_{l}}{2}\right)\mathbb{E}\left[\left\|\boldsymbol{h}_{t,n}^{f}\right\|^{2}\right] + 8L_{f}^{2}L_{z}^{2}\alpha_{t}^{3}N\sum_{i=0}^{N-1}\mathbb{E}\left[\left\|\boldsymbol{h}_{t,i}^{f}\right\|^{2}\right]$$

$$+ 2\alpha_{t}L_{f}^{2}\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] + \left(4\sigma_{g_{xy}}^{2}\alpha_{t}+4L_{f}^{2}\alpha_{t}\right)\mathbb{E}\left[\left\|\mathbf{z}_{t}-\mathbf{z}_{t}^{*}\right\|^{2}\right] + 2\sigma_{f_{x}}^{2}\alpha_{t}$$

1075

1069

1070 1071

$$+16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\alpha_{t}\sum_{i=0}^{N-1}\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i},\mathbf{y}_{t}^{i}\right)\right\|^{2}\right]+16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N^{2}\sigma_{g_{y}}^{2}\alpha_{t}+4\sigma_{g_{xy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\alpha_{t}$$

1077 1078

for all  $t \in \{0, 1, ..., T-1\}$  and  $n \in \{0, 1, ..., N-1\}$ , where the expectation is taken over the stochasticity of the algorithm.

 $\mathbb{E}\left[\ell\left(\mathbf{x}_{t}^{n+1}\right)-\ell\left(\mathbf{x}_{t}^{n}\right)\right]$ 

 Proof. We have

$$\overset{(a)}{\leq} \mathbb{E} \left[ \left\langle \nabla \ell \left( \mathbf{x}_{t}^{n} \right), \mathbf{x}_{t}^{n+1} - \mathbf{x}_{t}^{n} \right\rangle + \frac{L_{l}}{2} \left\| \mathbf{x}_{t}^{n+1} - \mathbf{x}_{t}^{n} \right\|^{2} \right]$$

$$\overset{(b)}{=} \mathbb{E} \left[ -\alpha_{t} \left\langle \nabla \ell \left( \mathbf{x}_{t}^{n} \right), h_{t,n}^{f} \right\rangle + \frac{\alpha_{t}^{2} L_{l}}{2} \left\| h_{t,n}^{f} \right\|^{2} \right]$$

$$\overset{(c)}{=} \mathbb{E} \left[ -\frac{\alpha_{t}}{2} \left\| \nabla \ell \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} - \frac{\alpha_{t}}{2} \left\| h_{t,n}^{f} \right\|^{2} + \frac{\alpha_{t}}{2} \left\| \nabla \ell \left( \mathbf{x}_{t}^{n} \right) - h_{t,n}^{f} \right\|^{2} + \frac{\alpha_{t}^{2} L_{l}}{2} \left\| h_{t,n}^{f} \right\|^{2} \right], \quad (8)$$

where (a) uses the Lipschitz continuous gradients of  $\ell$  (see Lemma 5.1). (b) follows from the update rule of Algorithm 1. (c) is because of  $\langle x, y \rangle = \frac{1}{2} ||x||^2 + \frac{1}{2} ||y||^2 - \frac{1}{2} ||x - y||^2$ .

Next, we bound the third term on the right in (8) above. Before that, we bound  $\|\mathbf{x}_t^n - \mathbf{x}_t^0\|^2$  and  $\|\mathbf{y}_t^n - \mathbf{y}_t^0\|^2$ .

$$\left\|\mathbf{x}_{t}^{n} - \mathbf{x}_{t}^{0}\right\|^{2} \stackrel{(a)}{=} \alpha_{t}^{2} \left\|\sum_{i=0}^{n-1} h_{t,i}^{f}\right\|^{2} \stackrel{(b)}{\leq} \alpha_{t}^{2} n \sum_{i=0}^{n-1} \left\|h_{t,i}^{f}\right\|^{2} \leq \alpha_{t}^{2} N \sum_{i=0}^{N-1} \left\|h_{t,i}^{f}\right\|^{2},\tag{9}$$

where (a) is because of the update rule of Algorithm 1. (b) is due to  $||z_1 + \cdots + z_k||^2 \le k ||z_1||^2 + \cdots + k ||z_k||^2$ .

1102 Similarly,

$$\left\|\mathbf{y}_{t}^{n} - \mathbf{y}_{t}^{0}\right\|^{2} \leq \beta_{t}^{2} N \sum_{i=0}^{N-1} \left\|h_{t,i}^{g}\right\|^{2}.$$
 (10)

$$\begin{aligned} & 1134 \\ & +16L_{f}^{2}L_{z}^{2}\alpha_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2} + 16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{g}\right\|^{2} \\ & 1137 \\ & \leq \mathbb{E}\left[2\left\|h_{t,n}^{f} - \nabla f\left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathbf{z}_{t}\right)\right\|^{2} + 4L_{f}^{2}\left\|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2} + 8L_{f}^{2}\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2} \\ & + 16L_{f}^{2}L_{z}^{2}\alpha_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2} + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{g} - \nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \\ & + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \\ & 1142 \\ & + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \\ & 1144 \\ & + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \\ & + 16L_{f}^{2}L_{z}^{2}\alpha_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2} + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \right] \\ & + 16L_{f}^{2}L_{z}^{2}\alpha_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2} + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \right] \\ & + 16L_{f}^{2}L_{z}^{2}\alpha_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2} + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \right] \\ & + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2} + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \right] \\ & + 16L_{f}^{2}L_{z}^{2}\alpha_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2} + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \right] \\ & + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2} + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \right] \\ & + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2} + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{i}, \mathbf{y}_{t}^{i}\right)\right\|^{2} \right] \\ & + 32L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\sum_{i=0}^{N-1}\left\|h_{t,i}^{f}\right\|^{2}$$

where  $\mathbf{z}_t = \mathbf{z} \left( \mathbf{x}_t^{\mathrm{o}}, \mathbf{y}_t^{\mathrm{o}} \right)$  and  $\mathbf{z}_t^* = \mathbf{z}^* \left( \mathbf{x}_t^{\mathrm{o}}, \mathbf{y}_t^{\mathrm{o}} \right)$ . (a), (c) and (f) follow from  $||x + y||^2 \le 2 ||x||^2 + 2 ||y||^2$  and  $\left\| \nabla_{\mathbf{x}\mathbf{y}}^2 g(\mathbf{x}, \mathbf{y}) \right\| \le B_{g_{xy}}$ . (b) utilizes the Lipschitzness of  $\nabla f(\mathbf{x}, \mathbf{y}, \mathbf{z})$  (see Lemma 5.1), and (d) is due to the Lipschitzness of  $z^*(x, y)$  (see Lemma 5.2). (e) uses equations (9) and (10). (g) is because of the bounded variance in Assumption 3.4. 

Then, we bound the term 
$$\mathbb{E} \left[ \left\| h_{t,n}^{f} - \nabla f \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathbf{z}_{t} \right) \right\|^{2} \right].$$
  
Then, we bound the term  $\mathbb{E} \left[ \left\| h_{t,n}^{f} - \nabla f \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathbf{z}_{t} \right) \right\|^{2} \right]$   
 $\mathbb{E} \left[ \left\| h_{t,n}^{f} - \nabla f \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathbf{z}_{t} \right) \right\|^{2} \right]$   
 $\mathbb{E} \left[ \left\| \nabla_{\mathbf{x}} f \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathcal{D}_{t,n}^{f_{x}} \right) + \nabla_{\mathbf{xy}}^{2} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathcal{D}_{t,n}^{g_{xy}} \right) \mathbf{z}_{t} - \nabla_{\mathbf{x}} f \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) - \nabla_{\mathbf{x}}^{2} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathcal{D}_{t,n}^{g_{xy}} \right) \mathbf{z}_{t} - \nabla_{\mathbf{x}} f \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathbf{z}_{t} \right) \right\|^{2} \right]$   
 $\mathbb{E} \left[ \left\| \nabla_{\mathbf{x}} f \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathcal{D}_{t,n}^{f_{x}} \right) - \nabla_{\mathbf{x}} f \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} + \left\| \mathbf{z}_{t} \right\|^{2} \left\| \nabla_{\mathbf{xy}}^{2} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}, \mathcal{D}_{t,n}^{g_{xy}} \right) - \nabla_{\mathbf{x}}^{2} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} \right]$   
 $\mathbb{E} \left[ \mathbb{E} \left[ 2\sigma_{g_{xy}}^{2} \| \mathbf{z}_{t} \|^{2} + 2\sigma_{f_{x}}^{2} \right]$   
 $\mathbb{E} \left[ \left\{ 4\sigma_{g_{xy}}^{2} \| \mathbf{z}_{t} - \mathbf{z}_{t}^{*} \|^{2} + 4\sigma_{g_{xy}}^{2} \| \frac{B_{f_{y}}^{2}}{\mu_{g}^{2}} + 2\sigma_{f_{x}}^{2} \right]$   
 $\mathbb{E} \left[ \left\{ 4\sigma_{g_{xy}}^{2} \| \mathbf{z}_{t} - \mathbf{z}_{t}^{*} \|^{2} + 4\sigma_{g_{xy}}^{2} \frac{B_{f_{y}}^{2}}{\mu_{g}^{2}} + 2\sigma_{f_{x}}^{2} \right]$ , (12)

where (a) uses the definitions of  $h_{t,n}^f$  and  $\nabla f(\mathbf{x}_t^n, \mathbf{y}_t^n, \mathbf{z}_t)$ . (b) utilizes the bounded variance in Assumption 3.4. (c) uses  $||x + y||^2 \le 2 ||x||^2 + 2 ||y||^2$ , and (d) is due to the bound of  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$  in Lemma 5.2. 

Combining (8), (11) and (12) completes the proof of the lemma. 

#### C.2.2 DESCENT IN THE ERROR OF $\mathbf{y}^{*}(\mathbf{x})$

**Lemma C.6.** Under Assumptions 3.2–3.4, the approximation error of  $\mathbf{y}^*(\mathbf{x})$  of Algorithm 1 satisfies the following inequality:

$$\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n+1} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n+1}\right)\right\|^{2}\right] \leq (1 + c_{1})\left(1 + c_{2}\right)\left(1 - 2\beta_{t}\frac{\mu_{g}L_{g}}{\mu_{g} + L_{g}}\right)\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] + (1 + c_{1})\left(1 + \frac{1}{c_{2}}\right)\beta_{t}^{2}\sigma_{g_{y}}^{2}$$

1186  
1187 
$$+ (1+c_1)(1+c_2)\left(\beta_t^2 - \frac{2\beta_t}{\mu_g + L_g}\right) \mathbb{E}\left[\|\nabla_{\mathbf{y}}g(\mathbf{x}_t^n, \mathbf{y}_t^n)\|^2\right] + \left(1 + \frac{1}{c_1}\right)L_y^2\alpha_t^2\mathbb{E}\left[\left\|h_{t,n}^f\right\|^2\right],$$

for all  $t \in \{0, 1, ..., T-1\}$  and  $n \in \{0, 1, ..., N-1\}$  with some constants  $c_1, c_2 > 0$ , where the expectation is taken over the stochasticity of the algorithm.

<sup>1191</sup> *Proof.* We have

$$\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n+1}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n+1}\right)\right\|^{2}\right] \\
\stackrel{(a)}{\leq} \mathbb{E}\left[\left(1+c_{1}\right)\left\|\mathbf{y}_{t}^{n+1}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}+\left(1+\frac{1}{c_{1}}\right)\left\|\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n+1}\right)\right\|^{2}\right] \\
\stackrel{(b)}{\leq} \mathbb{E}\left[\left(1+c_{1}\right)\left\|\mathbf{y}_{t}^{n}-\beta_{t}h_{t,n}^{g}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}+\left(1+\frac{1}{c_{1}}\right)L_{y}^{2}\left\|\mathbf{x}_{t}^{n+1}-\mathbf{x}_{t}^{n}\right\|^{2}\right] \\
\stackrel{(c)}{\leq}\left(1+c_{1}\right)\left(1+c_{2}\right)\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n}-\beta_{t}\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] \\
+\left(1+c_{1}\right)\left(1+\frac{1}{c_{2}}\right)\beta_{t}^{2}\mathbb{E}\left[\left\|h_{t,n}^{g}-\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)\right\|^{2}\right]+\left(1+\frac{1}{c_{1}}\right)L_{y}^{2}\alpha_{t}^{2}\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] \\
\stackrel{(d)}{\leq}\left(1+c_{1}\right)\left(1+c_{2}\right)\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n}-\beta_{t}\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right]+\left(1+\frac{1}{c_{1}}\right)L_{y}^{2}\alpha_{t}^{2}\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] \\
+\left(1+c_{1}\right)\left(1+\frac{1}{c_{2}}\right)\beta_{t}^{2}\sigma_{g_{y}}^{2},$$
(13)

where (a) results from Young's inequality. (b) is because of the update rule of Algorithm 1 and the Lipschitzness of  $y^*(\cdot)$  (see Lemma 5.1). (c) follows from Young's inequality and the update rule of Algorithm 1. (d) uses the bounded variance in Assumption 3.4.

 $\mathbf{v}^* (\mathbf{x}^n_t) \|^2$ 

<sup>1212</sup> To bound the first term on the right, we have

1214 
$$\|\mathbf{y}_{t}^{n} - \beta_{t} \nabla_{\mathbf{y}} g\left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}\right) -$$

$$= \|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}(\mathbf{x}_{t}^{n})\|^{2} + \beta_{t}^{2} \|\nabla_{\mathbf{y}}g(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n})\|^{2} - 2\beta_{t} \langle \nabla_{\mathbf{y}}g(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}), \mathbf{y}_{t}^{n} - \mathbf{y}^{*}(\mathbf{x}_{t}^{n}) \rangle$$

$$\stackrel{(a)}{\leq} \left(1 - 2\beta_{t} \frac{\mu_{g}L_{g}}{\mu_{g} + L_{g}}\right) \|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}(\mathbf{x}_{t}^{n})\|^{2} + \left(\beta_{t}^{2} - \frac{2\beta_{t}}{\mu_{g} + L_{g}}\right) \|\nabla_{\mathbf{y}}g(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n})\|^{2}, \quad (14)$$

where (a) is due to  $\mu_g$ -strongly convexity and  $L_g$ -smoothness of the lower-level function  $g(\mathbf{x}, \mathbf{y})$  (see Assumption 3.2), which implies

$$\left\langle \nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right),\mathbf{y}_{t}^{n}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\rangle \geq\frac{\mu_{g}L_{g}}{\mu_{g}+L_{g}}\left\|\mathbf{y}_{t}^{n}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}+\frac{1}{\mu_{g}+L_{g}}\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)\right\|^{2}.$$

1225 The Lemma is proved by substituting (14) in (13).

1227 C.2.3 Descent in the error of  $\mathbf{z}^{*}(\mathbf{x}, \mathbf{y})$ 

**Lemma C.7.** Under Assumptions 3.1–3.4, the following inequality of the approximation error of  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$  holds for Algorithm 1:

$$\begin{split} & \mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t+1}^{*}\right\|^{2}\right] \\ & \leq \left(\left(1 + c_{3}\right)\left(1 + c_{4}\right)\left(1 - \frac{2\gamma_{t}\mu_{g}L_{q}}{\mu_{g} + L_{q}}\right) + 4\sigma_{g_{yy}}^{2}\gamma_{t}^{2}\left(1 + c_{3}\right)\left(1 + \frac{1}{c_{4}}\right)\right)\mathbb{E}\left[\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2}\right] \\ & + \left(1 + c_{3}\right)\left(1 + c_{4}\right)\left(\gamma_{t}^{2} - \frac{2\gamma_{t}}{\mu_{g} + L_{q}}\right)\mathbb{E}\left[\left\|\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right)\right\|^{2}\right] + 4\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\beta_{t}^{2}N^{2}\sigma_{g_{y}}^{2} \\ & + 2\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\alpha_{t}^{2}N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] + 4\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\beta_{t}^{2}N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}\right)\right\|^{2}\right] \\ & + 4\sigma_{g_{yy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\left(1 + c_{3}\right)\left(1 + \frac{1}{c_{4}}\right)\gamma_{t}^{2} + 2\sigma_{f_{y}}^{2}\left(1 + c_{3}\right)\left(1 + \frac{1}{c_{4}}\right)\gamma_{t}^{2}, \end{split}$$

for all  $t \in \{0, 1, ..., T-1\}$  and  $n \in \{0, 1, ..., N-1\}$  with some constants  $c_3, c_4 > 0$ , where  $\mathbf{z}_t = \mathbf{z} \left( \mathbf{x}_t^0, \mathbf{y}_t^0 \right)$  and  $\mathbf{z}_t^* = \mathbf{z}^* \left( \mathbf{x}_t^0, \mathbf{y}_t^0 \right)$ . The expectation is taken over the stochasticity of the algorithm. Proof. We have  $\mathbb{E} \left\| \left\| \mathbf{z}_{t+1} - \mathbf{z}_{t+1}^* \right\|^2 \right\|$  $\stackrel{(a)}{\leq} \mathbb{E}\left[ (1+c_3) \|\mathbf{z}_{t+1} - \mathbf{z}_t^*\|^2 + \left(1 + \frac{1}{c_3}\right) \|\mathbf{z}^* \left(\mathbf{x}_{t+1}^0, \mathbf{y}_{t+1}^0\right) - \mathbf{z}^* \left(\mathbf{x}_t^0, \mathbf{y}_t^0\right) \|^2 \right]$  $\overset{(b)}{\leq} \mathbb{E}\left[ (1+c_3) \left\| \mathbf{z}_{t+1} - \mathbf{z}_t^* \right\|^2 + \left( 1 + \frac{1}{c_2} \right) L_z^2 \left( \left\| \mathbf{x}_{t+1}^0 - \mathbf{x}_t^0 \right\| + \left\| \mathbf{y}_{t+1}^0 - \mathbf{y}_t^0 \right\| \right)^2 \right]$  $\stackrel{(c)}{\leq} \mathbb{E}\left[ (1+c_3) \left\| \mathbf{z}_{t+1} - \mathbf{z}_t^* \right\|^2 + 2\left( 1 + \frac{1}{c_3} \right) L_z^2 \left\| \mathbf{x}_t^N - \mathbf{x}_t^0 \right\|^2 + 2\left( 1 + \frac{1}{c_3} \right) L_z^2 \left\| \mathbf{y}_t^N - \mathbf{y}_t^0 \right\|^2 \right]$  $\stackrel{(d)}{\leq} (1+c_3) \mathbb{E}\left[ \left\| \mathbf{z}_{t+1} - \mathbf{z}_t^* \right\|^2 \right] + 2\left( 1 + \frac{1}{c_3} \right) L_z^2 \alpha_t^2 N \sum_{s=1}^{N-1} \mathbb{E}\left[ \left\| h_{t,n}^f \right\|^2 \right]$  $+2\left(1+\frac{1}{c_3}\right)L_z^2\beta_t^2N\sum_{n=1}^{N-1}\mathbb{E}\left[\left\|h_{t,n}^g\right\|^2\right]$  $\stackrel{(e)}{\leq} (1+c_3) \mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_t^*\right\|^2\right] + 4\left(1 + \frac{1}{c_3}\right) L_z^2 \beta_t^2 N \sum_{i=1}^{N-1} \mathbb{E}\left[\left\|\nabla_{\mathbf{y}} g\left(\mathbf{x}_t^n, \mathbf{y}_t^n\right)\right\|^2\right]$  $+4\left(1+\frac{1}{c_{3}}\right)L_{z}^{2}\beta_{t}^{2}N\sum_{n=1}^{N-1}\mathbb{E}\left[\left\|h_{t,n}^{g}-\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)\right\|^{2}\right]+2\left(1+\frac{1}{c_{3}}\right)L_{z}^{2}\alpha_{t}^{2}N\sum_{n=1}^{N-1}\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right]$  $\stackrel{(f)}{\leq} (1+c_3) \mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_t^*\right\|^2\right] + 2\left(1 + \frac{1}{c_3}\right) L_z^2 \alpha_t^2 N \sum^{N-1} \mathbb{E}\left[\left\|h_{t,n}^f\right\|^2\right]$  $+4\left(1+\frac{1}{c_{3}}\right)L_{z}^{2}\beta_{t}^{2}N\sum^{N-1}\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)\right\|^{2}\right]+4\left(1+\frac{1}{c_{3}}\right)L_{z}^{2}\beta_{t}^{2}N^{2}\sigma_{g_{y}}^{2},$ (15)where (a) follows from Young's inequality. (b) is due to the Lipschitzness of  $\mathbf{z}^{*}(\cdot, \cdot)$  (see Lemma 5.2).

where (a) follows from Young's inequality. (b) is due to the Lipschitzness of  $z^*(\cdot, \cdot)$  (see Lemma 5.2). (c) and (e) result from  $||x + y||^2 \le 2 ||x||^2 + 2 ||y||^2$ . (d) is because of equations (9) and (10). (f) uses the bounded variance in Assumption 3.4.

1278 Next, we bound the first term on the right:

$$\mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t}^{*}\right\|^{2}\right] \stackrel{(a)}{=} \mathbb{E}\left[\left\|\mathbf{z}_{t} - \gamma_{t}h_{t}^{q} - \mathbf{z}_{t}^{*}\right\|^{2}\right] \\
\stackrel{(b)}{\leq} \mathbb{E}\left[\left(1 + c_{4}\right)\left\|\mathbf{z}_{t} - \gamma_{t}\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right) - \mathbf{z}_{t}^{*}\right\|^{2} + \left(1 + \frac{1}{c_{4}}\right)\gamma_{t}^{2}\left\|\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right) - h_{t}^{q}\right\|^{2}\right] \\
= \mathbb{E}\left[\left(1 + c_{4}\right)\left(\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2} + \gamma_{t}^{2}\left\|\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right)\right\|^{2} - 2\gamma_{t}\left\langle\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right), \mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\rangle\right)\right] \\
+ \left(1 + \frac{1}{c_{4}}\right)\gamma_{t}^{2}\mathbb{E}\left[\left\|\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right) - h_{t}^{q}\right\|^{2}\right] \\
\stackrel{(c)}{\leq}\left(1 + c_{4}\right)\left(1 - 2\gamma_{t}\frac{\mu_{g}L_{q}}{\mu_{g} + L_{q}}\right)\mathbb{E}\left[\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2}\right] + \left(1 + \frac{1}{c_{4}}\right)\gamma_{t}^{2}\mathbb{E}\left[\left\|\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right) - h_{t}^{q}\right\|^{2}\right] \\
+ \left(1 + c_{4}\right)\left(\gamma_{t}^{2} - \frac{2\gamma_{t}}{\mu_{g} + L_{q}}\right)\mathbb{E}\left[\left\|\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right)\right\|^{2}\right],$$
(16)

where (a) results from the update rule of Algorithm 1, and (b) uses Young's inequality. (c) follows from  $\mu_g$ -strongly convexity and  $L_q$ -smoothness of  $q(\mathbf{x}, \mathbf{y}, \mathbf{z})$ , which implies

$$\left\langle \nabla_{\mathbf{z}} q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right), \mathbf{z}_{t} - \mathbf{z}_{t}^{*} \right\rangle \geq \frac{\mu_{g} L_{q}}{\mu_{g} + L_{q}} \left\| \mathbf{z}_{t} - \mathbf{z}_{t}^{*} \right\|^{2} + \frac{1}{\mu_{g} + L_{q}} \left\| \nabla_{\mathbf{z}} q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right) \right\|^{2}$$

Then, we bound the second term on the right as follows:

1312 where (a) follows from the definitions of  $h_t^q$  and  $\nabla_{\mathbf{z}} q(\mathbf{x}, \mathbf{y}, \mathbf{z})$ . (b) and (d) are because of  $||x + y||^2 \le 2 ||x||^2 + 2 ||y||^2$ . (c) results from the bounded variances in Assumption 3.4. (e) utilizes the bound of 1314  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$  in Lemma 5.2.

Substituting (17) in (16) and then substituting the result in (15), the lemma is proved.

# 1317 C.2.4 DESCENT IN THE POTENTIAL FUNCTION

1319 We define the potential function  $W_t$  as follows:

$$W_{t} = \ell \left( \mathbf{x}_{t}^{0} \right) + K_{y} \left\| \mathbf{y}_{t}^{0} - \mathbf{y}^{*} \left( x_{t}^{0} \right) \right\|^{2} + K_{z} \left\| \mathbf{z}_{t} \left( \mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0} \right) - \mathbf{z}^{*} \left( \mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0} \right) \right\|^{2}$$

**Lemma C.8.** Set  $c_1 = \frac{\beta_t L_{\mu_g}}{2(1-\beta_t L_{\mu_g})}$ ,  $c_2 = \frac{\beta_t L_{\mu_g}}{1-2\beta_t L_{\mu_g}}$ ,  $c_3 = \frac{\gamma_t L_{\mu_q}}{2(1-\gamma_t L_{\mu_q})}$ , and  $c_4 = \frac{\gamma_t L_{\mu_q}}{1-2\gamma_t L_{\mu_q}}$ . Under the same conditions as described in Theorem C.9 and using Lemmas C.1-C.3, the iterates generated by Algorithm 1 satisfies: for all  $t \in \{0, 1, ..., T-1\}$ ,

 $\mathbb{E}\left[W_{t+1} - W_t\right] \leq -\frac{\alpha_t}{2} \sum_{n=0}^{N-1} \mathbb{E}\left[\left\|\nabla \ell\left(\mathbf{x}_t^n\right)\right\|^2\right] + \sigma_{g_{xy}}^2 \alpha_t C_{g_{xy}} N + \sigma_{f_x}^2 \alpha_t C_{f_x} N + \sigma_{g_{yy}}^2 \alpha_t C_{g_{yy}} N + \sigma_{f_y}^2 \alpha_t C_{f_y} N + \sigma_{g_y}^2 \alpha_t \left(C_{g_1} N + C_{g_2} \frac{1}{N}\right),$ 

where the constant values  $C_{g_{xy}}$ ,  $C_{f_x}$ ,  $C_{g_1}$ ,  $C_{g_2}$ ,  $C_{g_{yy}}$  and  $C_{f_y}$ , which are independent of N, are defined as:

$$C_{g_{xy}} = \frac{4B_{f_y}^2}{\mu_g^2}, \qquad C_{f_x} = 2, \qquad C_{g_1} = \frac{2c_\beta K_y}{L_{\mu_g}} + \frac{8L_z^2 c_\beta^2 N K_z}{c_\gamma L_{\mu_q}}, C_{g_2} = \frac{L_l^2 c_\beta^2}{16L_f^2 L_z^2}, \qquad C_{g_{yy}} = \frac{8B_{f_y}^2 c_\gamma K_z}{\mu_g^2 L_{\mu_q} N}, \qquad C_{f_y} = \frac{4c_\gamma K_z}{L_{\mu_q} N},$$
(18)

1335 1336 1337

1338

1340

1334

1320 1321

1326

1327 1328 1329

1330 1331

where 
$$K_y$$
 and  $K_z$  are defined in (23) of Theorem C.9.

1341 *Proof.* From Lemma C.1, we have

$$\begin{aligned} & \begin{array}{l} & \begin{array}{l} \mathbf{1342} \\ & \mathbf{1343} \\ & \mathbf{1343} \\ & \begin{array}{l} \mathbf{1344} \\ & \mathbf{1344} \\ & \begin{array}{l} \mathbf{1345} \\ & \mathbf{1345} \\ & \begin{array}{l} \mathbf{1346} \\ & \mathbf{1346} \\ & \begin{array}{l} \mathbf{1347} \\ & \end{array} \\ & \begin{array}{l} \mathbf{1346} \\ & \begin{array}{l} \mathbf{1347} \\ & \end{array} \\ & \begin{array}{l} \mathbf{1347} \\ & \end{array} \\ & \begin{array}{l} \mathbf{1348} \\ & \end{array} \\ & \begin{array}{l} \mathbf{1349} \end{array} \\ & \end{array} \\ & \begin{array}{l} \mathbf{1348} \\ & \end{array} \\ & \begin{array}{l} \mathbf{1349} \\ & \end{array} \\ & \begin{array}{l} \mathbf{1348} \\ \mathbf{1349} \end{array} \\ \\ & \begin{array}{l} \mathbf{1348} \\ \mathbf{1349} \end{array} \\ \\$$

$$\begin{array}{l} \mathbf{1350} \\ \mathbf{1351} \\ \mathbf{1352} \end{array} + 2\alpha_t L_f^2 \sum_{n=0}^{N-1} \mathbb{E}\left[ \left\| \mathbf{y}_t^n - \mathbf{y}^* \left( \mathbf{x}_t^n \right) \right\|^2 \right] + 4\tilde{\sigma}_{g_{xy}}^2 \left| \mathcal{D}^{g_{xy}} \right|^{-1} \frac{B_{f_y}^2}{\mu_g^2} N \alpha_t + 16L_f^2 L_z^2 \beta_t^2 N^3 \sigma_{g_y}^2 \alpha_t.$$

Choosing  $\alpha_t \leq \frac{L_l}{16L_f^2 L_z^2 N^2}$  and using the definition of  $\beta_t = c_\beta \alpha_t$ , we get

$$\begin{aligned} & 1356 \\ 1357 \\ & 1358 \\ 1358 \\ 1358 \\ 1359 \\ 1360 \\ & + \left(-\frac{\alpha_t}{2} + \alpha_t^2 L_l\right) \sum_{n=0}^{N-1} \mathbb{E}\left[ \left\| \nabla \ell \left( \mathbf{x}_t^n \right) \right\|^2 \right] + \left( 4\sigma_{g_{xy}}^2 N \alpha_t + 4L_f^2 N \alpha_t \right) \mathbb{E}\left[ \left\| \mathbf{z}_t - \mathbf{z}_t^* \right\|^2 \right] \\ & + \left( -\frac{\alpha_t}{2} + \alpha_t^2 L_l \right) \sum_{n=0}^{N-1} \mathbb{E}\left[ \left\| h_{t,n}^f \right\|^2 \right] + 2\alpha_t L_f^2 \sum_{n=0}^{N-1} \mathbb{E}\left[ \left\| \mathbf{y}_t^n - \mathbf{y}^* \left( \mathbf{x}_t^n \right) \right\|^2 \right] + 2\sigma_{f_x}^2 N \alpha_t \end{aligned}$$

$$+16L_{f}^{2}L_{z}^{2}c_{\beta}^{2}N^{2}\alpha_{t}^{3}\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)\right\|^{2}\right]+\frac{L_{l}^{2}c_{\beta}^{2}}{16L_{f}^{2}L_{z}^{2}N}\sigma_{g_{y}}^{2}\alpha_{t}+4\sigma_{g_{xy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}N\alpha_{t}.$$
(19)

With the result from Lemma C.2, we have

$$\begin{split} & \sum_{n=0}^{N-1} \mathbb{E} \left[ \left\| \mathbf{y}_{t}^{n+1} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n+1} \right) \right\|^{2} - \left\| \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} \right] = \mathbb{E} \left[ \left\| \mathbf{y}_{t+1}^{0} - \mathbf{y}^{*} \left( \mathbf{x}_{t+1}^{0} \right) \right\|^{2} - \left\| \mathbf{y}_{t}^{0} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{0} \right) \right\|^{2} \right] \\ & \left\{ \left( (1+c_{1}) \left( 1+c_{2} \right) \left( 1-2\beta_{t} \frac{\mu_{g}L_{g}}{\mu_{g}+L_{g}} \right) - 1 \right) \sum_{n=0}^{N-1} \mathbb{E} \left[ \left\| \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} \right] \\ & \left\{ \left( 1+c_{1} \right) \left( 1+c_{2} \right) \left( 1-2\beta_{t} \frac{\mu_{g}L_{g}}{\mu_{g}+L_{g}} \right) - 1 \right) \sum_{n=0}^{N-1} \mathbb{E} \left[ \left\| \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} \right] \\ & \left\{ \left( 1+c_{1} \right) \left( 1+c_{2} \right) \left( 1-2\beta_{t} \frac{\mu_{g}L_{g}}{\mu_{g}+L_{g}} \right) + \left( 1+c_{1} \right) \left( 1+c_{2} \right) \beta_{t}^{2} \tilde{\sigma}_{g}^{2} \left\| \mathcal{D}^{g_{y}} \right\|^{-1} N \\ & \left\{ \left( 1+c_{1} \right) \left( 1+c_{2} \right) \left( \beta_{t}^{2} - \frac{2\beta_{t}}{\mu_{g}+L_{g}} \right) \sum_{n=0}^{N-1} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} \right] . \\ & \left\{ 1+c_{1} \right) \left( 1+c_{2} \right) \left( \beta_{t}^{2} - \frac{2\beta_{t}}{\mu_{g}+L_{g}} \right) \sum_{n=0}^{N-1} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} \right] . \end{split} \right] \end{split}$$

Denote  $L_{\mu_g} = \frac{\mu_g L_g}{\mu_g + L_g}$ . Choose  $c_1$  and  $c_2$  such that  $(1+c_1)(1+c_2)(1-2\beta_t L_{\mu_g}) = 1 - \frac{\beta_t L_{\mu_g}}{2}.$ 

$$(1+c_2)\left(1-2\beta_t L_{\mu_g}\right) = 1-\beta_t L_{\mu_g} \implies c_2 = \frac{\beta_t L_{\mu_g}}{1-2\beta_t L_{\mu_g}} \& \beta_t \le \frac{1}{2L_{\mu_g}}.$$

Thus,

$$c_1 = \frac{\beta_t L_{\mu_g}}{2\left(1 - \beta_t L_{\mu_g}\right)}.$$

Moreover, this implies that

$$1 + \frac{1}{c_2} = 1 + \frac{1 - 2\beta_t L_{\mu_g}}{\beta_t L_{\mu_g}} \le \frac{1}{\beta_t L_{\mu_g}}, \qquad 1 + \frac{1}{c_1} = \frac{2\left(1 - \beta_t L_{\mu_g}\right)}{\beta_t L_{\mu_g}} \le \frac{2}{\beta_t L_{\mu_g}}.$$

Use the definition of  $\beta_t = c_\beta \alpha_t$ . Substituting  $c_1$  and  $c_2$  and choosing  $\beta_t \leq \frac{1}{\mu_g + L_g}$ , we have 

$$\mathbb{E}\left[\left\|\mathbf{y}_{t+1}^{0}-\mathbf{y}^{*}\left(\mathbf{x}_{t+1}^{0}\right)\right\|^{2}-\left\|\mathbf{y}_{t}^{0}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{0}\right)\right\|^{2}\right] \leq -\frac{c_{\beta}L_{\mu_{g}}}{2}\alpha_{t}\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right]$$

$$\frac{1402}{1403} + \frac{2L_{y}^{2}\alpha_{t}}{c_{\beta}L_{\mu_{g}}} \sum_{n=0}^{N-1} \mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] - \frac{c_{\beta}\alpha_{t}}{\mu_{g} + L_{g}} \sum_{n=0}^{N-1} \mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}\right)\right\|^{2}\right] + \frac{2}{L_{\mu_{g}}}c_{\beta}\alpha_{t}\sigma_{g_{y}}^{2}N.$$
(20)

1404  
1405  
1406 According to Lemma C.3, we have  
$$\mathbb{E} \left[ \| \mathbf{z}_{t+1} - \mathbf{z}_{t+1}^* \|^2 - \| \mathbf{z}_t - \mathbf{z}_t^* \|^2 \right]$$

$$\begin{split} & \mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t+1}^{*}\right\|^{2} - \|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\|^{2}\right] \\ & = \mathbb{E}\left[\left\|\mathbf{z}\left(\mathbf{x}_{t+1}^{0}, \mathbf{y}_{t+1}^{0}\right) - \mathbf{z}^{*}\left(\mathbf{x}_{t+1}^{0}, \mathbf{y}_{t+1}^{0}\right)\right\|^{2} - \left\|\mathbf{z}\left(\mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0}\right) - \mathbf{z}^{*}\left(\mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0}\right)\right\|^{2}\right] \\ & = \mathbb{E}\left[\left\|\mathbf{z}\left(\mathbf{x}_{t+1}^{0}, \mathbf{y}_{t+1}^{0}\right) - \mathbf{z}^{*}\left(\mathbf{x}_{t+1}^{0}, \mathbf{y}_{t+1}^{0}\right)\right\|^{2} - \left\|\mathbf{z}\left(\mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0}\right) - \mathbf{z}^{*}\left(\mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0}\right)\right\|^{2}\right] \\ & \leq \left(\left(1 + c_{3}\right)\left(1 + c_{4}\right)\left(1 - \frac{2\gamma_{t}\mu_{g}L_{q}}{\mu_{g} + L_{q}}\right) + 4\sigma_{g_{yy}}^{2}\gamma_{t}^{2}\left(1 + c_{3}\right)\left(1 + \frac{1}{c_{4}}\right) - 1\right)\mathbb{E}\left[\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2}\right] \\ & + \left(1 + c_{3}\right)\left(1 + c_{4}\right)\left(\gamma_{t}^{2} - \frac{2\gamma_{t}}{\mu_{g} + L_{q}}\right)\mathbb{E}\left[\left\|\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right)\right\|^{2}\right] + 4\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\beta_{t}^{2}N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\mathbf{x}_{z}^{2}\beta_{t}^{2}N^{2}\sigma_{g_{yy}}^{2}\right. \\ & + 2\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\alpha_{t}^{2}N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] + 4\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\beta_{t}^{2}N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}\right)\right\|^{2}\right] \\ & + 2\sigma_{f_{y}}^{2}\left(1 + c_{3}\right)\left(1 + \frac{1}{c_{4}}\right)\gamma_{t}^{2} + 4\sigma_{g_{yy}}^{2}\frac{B_{f_{yy}}^{2}}{\mu_{g}^{2}}\left(1 + c_{3}\right)\left(1 + \frac{1}{c_{4}}\right)\gamma_{t}^{2}. \end{split}$$

Similar as  $c_1$  and  $c_2$ , we choose

$$c_3 = \frac{\gamma_t L_{\mu_q}}{2(1 - \gamma_t L_{\mu_q})}, \qquad c_4 = \frac{\gamma_t L_{\mu_q}}{1 - 2\gamma_t L_{\mu_q}}$$

where  $\gamma_t \leq \frac{1}{2L_{\mu_q}}$  and we denote  $L_{\mu_q} = \frac{\mu_g L_q}{\mu_g + L_q}$ . This implies that 

$$1 + \frac{1}{c_4} \le \frac{1}{\gamma_t L_{\mu_q}}, \qquad 1 + \frac{1}{c_3} \le \frac{2}{\gamma_t L_{\mu_q}}.$$

According to the definitions of  $\beta_t = c_{\beta}\alpha_t$  and  $\gamma_t = c_{\gamma}\alpha_t$ , substituting  $c_3$  and  $c_4$  and choosing  $\gamma_t \leq \frac{1}{\mu_g + L_q}$ , we get 

$$\begin{aligned}
\mathbf{E} \begin{bmatrix} \|\mathbf{z} \left(\mathbf{x}_{t+1}^{0}, \mathbf{y}_{t+1}^{0}\right) - \mathbf{z}^{*} \left(\mathbf{x}_{t+1}^{0}, \mathbf{y}_{t+1}^{0}\right) \|^{2} - \|\mathbf{z} \left(\mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0}\right) - \mathbf{z}^{*} \left(\mathbf{x}_{t}^{0}, \mathbf{y}_{t}^{0}\right) \|^{2} \end{bmatrix} \\
& \leq \left( -\frac{c_{\gamma} L_{\mu_{q}}}{2} \alpha_{t} + \frac{8}{L_{\mu_{q}}} \sigma_{g_{yy}}^{2} c_{\gamma} \alpha_{t} \right) \mathbb{E} \left[ \|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\|^{2} \right] - \frac{c_{\gamma} \alpha_{t}}{\mu_{g} + L_{q}} \mathbb{E} \left[ \|\nabla_{\mathbf{z}}q \left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right) \|^{2} \right] \\
& + \frac{4L_{z}^{2} \alpha_{t} N}{c_{\gamma} L_{\mu_{q}}} \sum_{n=0}^{N-1} \mathbb{E} \left[ \left\| h_{t,n}^{f} \right\|^{2} \right] + \frac{8L_{z}^{2} c_{\beta}^{2} \alpha_{t} N}{c_{\gamma} L_{\mu_{q}}} \sum_{n=0}^{N-1} \mathbb{E} \left[ \|\nabla_{\mathbf{y}}g \left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}\right) \|^{2} \right] + \frac{8}{c_{\gamma} L_{\mu_{q}}} L_{z}^{2} c_{\beta}^{2} \alpha_{t} N^{2} \sigma_{g_{y}}^{2} \\
& + 2\sigma_{f_{y}}^{2} \frac{2}{L_{\mu_{q}}} c_{\gamma} \alpha_{t} + 4\sigma_{g_{yy}}^{2} \frac{B_{f_{y}}^{2}}{\mu_{g}^{2}} \frac{2}{L_{\mu_{q}}} c_{\gamma} \alpha_{t}.
\end{aligned}$$
(21)

Adding equations (19), (20) and (21), we get

$$\begin{split} & \mathbb{E}\left[W_{t+1} - W_{t}\right] \\ & \leq -\frac{\alpha_{t}}{2}\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\nabla l\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] + \bar{C}_{y}\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] + \bar{C}_{z}\mathbb{E}\left[\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2}\right] \\ & + \bar{C}_{g}\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}\right)\right\|^{2}\right] + \bar{C}_{h}\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] + \bar{C}_{q}\mathbb{E}\left[\left\|\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N}, \mathbf{z}_{t}\right)\right\|^{2}\right] \end{split}$$

$$+ \bar{C}_{g} \sum_{n=0}^{N-1} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} \right] + \bar{C}_{h} \sum_{n=0}^{N-1} \mathbb{E} \left[ \left\| h_{t,n}^{f} \right\|^{2} \right] + \bar{C}_{q} \mathbb{E} \left[ \left\| \nabla_{\mathbf{z}} q \left( \mathbf{x}_{t}^{N}, \mathbf{y}_{t}^{N} \right) + \frac{L_{l}^{2} c_{\beta}^{2}}{16 L^{2} L^{2} N} \sigma_{g_{y}}^{2} \alpha_{t} + 4 \sigma_{g_{xy}}^{2} \frac{B_{f_{y}}^{2}}{\mu^{2}} N \alpha_{t} + 2 \sigma_{f_{x}}^{2} N \alpha_{t} + K_{y} \frac{2 c_{\beta}}{L} \alpha_{t} \sigma_{g_{y}}^{2} N$$

$$+ \frac{-i c_{\beta}}{16L_{f}^{2}L_{z}^{2}N} \sigma_{g_{y}}^{2} \alpha_{t} + 4\sigma_{g_{xy}}^{2} \frac{J_{y}}{\mu_{g}^{2}} N \alpha_{t} + 2\sigma_{f_{x}}^{2} N \alpha_{t} + K_{y} \frac{2c_{\beta}}{L_{\mu_{g}}} \alpha_{t} \sigma_{g}^{2} \\ + K_{z} \left( \frac{8B_{f_{y}}^{2}}{\mu_{g}^{2}L_{\mu_{q}}} \sigma_{g_{yy}}^{2} c_{\gamma} \alpha_{t} + \frac{4c_{\gamma}}{L_{\mu_{q}}} \sigma_{f_{y}}^{2} \alpha_{t} + \frac{8}{c_{\gamma}L_{\mu_{q}}} L_{z}^{2} c_{\beta}^{2} \alpha_{t} N^{2} \sigma_{g_{y}}^{2} \right),$$

where

$$\bar{C}_y = 2\alpha_t L_f^2 - \frac{c_\beta L_{\mu_g}}{2} \alpha_t K_y$$

$$\begin{array}{ll} C_{z} = 4\alpha_{t}\sigma_{g,v}^{2}N + 4L_{f}^{2}N\alpha_{t} - \frac{c_{s}L_{\mu_{q}}}{2}\alpha_{t}K_{s} + \frac{8}{L_{\mu_{q}}}\sigma_{g,v}^{2}c_{\gamma}\alpha_{t}K_{z} \\ \tilde{C}_{g} = 16L_{f}^{2}L_{s}^{2}c_{g}^{2}N^{2}\alpha_{s}^{2} - \frac{c_{g}\alpha_{u}}{\mu_{g}}L_{y}^{2}\alpha_{t}K_{y} + \frac{8}{c_{s}L_{\mu_{q}}}L_{s}^{2}c_{g}^{2}\alpha_{t}NK_{z} \\ \tilde{C}_{h} = \alpha_{t}^{2}L_{t} - \frac{\alpha_{t}}{2} + \frac{2}{c_{g}L_{\mu_{g}}}L_{y}^{2}\alpha_{t}K_{y} + \frac{4}{c_{s}L_{\mu_{q}}}L_{s}^{2}\alpha_{t}NK_{z} \\ \tilde{C}_{q} = -\frac{c_{s}\alpha_{u}}{\mu_{g}+L_{q}}K_{z} \leq 0. \\ \tilde{C}_{q} = -\frac{c_{s}\alpha_{u}}{\mu_{g}+L_{q}}K_{z} \leq 0. \\ \tilde{C}_{z} = 4\alpha_{t}\sigma_{g,v}^{2}N + 4L_{f}^{2}N\alpha_{t} - \frac{c_{s}L_{\mu_{q}}}{2}\alpha_{t}K_{z} + \frac{8}{L_{\mu_{q}}}\sigma_{g,v}^{2}c_{r}\alpha_{t}K_{z} \\ = 4\alpha_{t}\sigma_{g,v}^{2}N + 4L_{f}^{2}N\alpha_{t} - \frac{c_{s}L_{\mu_{q}}}{2}\alpha_{t}K_{z} + \frac{8}{L_{\mu_{q}}}\sigma_{g,v}^{2}c_{r}\alpha_{t}K_{z} \\ = 4\alpha_{t}\sigma_{g,v}^{2}N + 4L_{f}^{2}N\alpha_{t} - \frac{c_{s}L_{\mu_{q}}}{2}\alpha_{t}K_{z} + \frac{8}{L_{\mu_{q}}}\sigma_{g,v}^{2}c_{r}\alpha_{t}K_{z} \\ = 4\alpha_{t}\sigma_{g,v}^{2}N + 4L_{f}^{2}N\alpha_{t} - \frac{c_{s}L_{\mu_{q}}}{2}\alpha_{t}K_{z} + \frac{8}{L_{\mu_{q}}}\sigma_{g,v}^{2}c_{r}\alpha_{t}K_{z} \\ = 4\alpha_{t}\sigma_{g,v}^{2}N + 4L_{f}^{2}N\alpha_{t} - \frac{c_{s}L_{\mu_{q}}}{2}\alpha_{t}K_{z} + \frac{8}{L_{\mu_{q}}}\sigma_{g,v}^{2}c_{r}\alpha_{t}K_{z} \\ = 4\alpha_{t}\sigma_{g,v}^{2}N + 4L_{f}^{2}N\alpha_{t} - \frac{c_{s}L_{\mu_{q}}}{2}\alpha_{t}K_{z} + \frac{8}{L_{\mu_{q}}}\sigma_{g,v}^{2}c_{r}\alpha_{t}K_{z} \\ = 4\alpha_{t}\sigma_{g,v}^{2}N + 4L_{f}^{2}N\alpha_{t} - \frac{c_{s}L_{\mu_{q}}}{2}\alpha_{t}K_{z} + \frac{8}{L_{\mu_{q}}}}\sigma_{g,v}^{2}C_{r}\alpha_{t}K_{z} \\ = 4\alpha_{t}\sigma_{g,v}^{2}N + 4L_{f}^{2}N\alpha_{t} - \frac{c_{s}L_{\mu_{q}}}{2}\alpha_{t}K_{x} + \frac{8}{L_{\mu_{q}}}}\lambda_{g,v}^{2}R_{t}K_{z} = 0, \\ \end{array}$$
where (a) utilizes  $K_{z} \ge max \left\{ \frac{24\sigma_{s}^{2}N^{2}N}{c_{\mu,\mu_{q}}}, \frac{24L_{f}^{2}N}{2}\alpha_{t}K_{y} + \frac{8}{c_{\mu_{\mu}}}L_{g,\sigma}^{2}L_{h}^{2}N_{h}K_{z} \\ \leq \alpha_{s}^{2}\frac{1}{2}\frac{c_{g}^{2}\Omega}N^{2}\alpha_{s}^{2} - \frac{c_{g}^{2}\Omega_{u}}{\mu_{g}+L_{g}}K_{g} + \frac{8}{L_{g}}\lambda_{g,\sigma}^{2}}N_{t}K_{x} \\ \leq \alpha_{s}^{2}\frac{1}{2}\frac{c_{g}^{2}\Omega}N^{2}\alpha_{s}^{2} - \frac{c_{g}^{2}\Omega_{u}}{\mu_{g}+L_{g}}K_{y} + \frac{8}{c_{\mu}}L_{\mu_{q}}}L_{g,\sigma}^{2}\Omega_{h}NK_{z} \\ \leq \alpha_{s}^{2}\frac{1}{2}\frac{c_{g}^{2}\Omega}N^{2}\alpha_{s}^{2} - \frac{c_{g}^{2}\Omega_{u}}{\mu_{g}+L_{g}}K_{y} + \frac{8}{2}\frac{c_{g}^{2}\Omega_{g}}N_{h}K_{z} \\ = 0, \\ \\ \text{where (a) results from \alpha_{t}$ 

Then, we get 

$$\begin{split} & \begin{array}{l} \mathbf{1514} \\ & \mathbf{1515} \\ & \mathbf{1515} \\ & \mathbf{1516} \\ \mathbf{1516} \\ & \mathbf{1516} \\ & \mathbf{1517} \\ & \mathbf{1518} \\ & \mathbf{1518} \\ & \mathbf{1518} \\ & \mathbf{1518} \\ \end{array} \\ & \begin{array}{l} \mathbf{\mathbb{E}} \left[ W_{t+1} - W_t \right] \leq -\frac{\alpha_t}{2} \sum_{n=0}^{N-1} \mathbb{E} \left[ \| \nabla l \left( \mathbf{x}_t^n \right) \|^2 \right] + \frac{L_l^2 c_\beta^2}{16 L_f^2 L_z^2 N} \sigma_{g_y}^2 \alpha_t + 4\sigma_{g_{xy}}^2 \frac{B_{f_y}^2}{\mu_g^2} N \alpha_t + 2\sigma_{f_x}^2 N \alpha_t \\ & + K_y \frac{2 c_\beta}{L_{\mu_g}} \alpha_t \sigma_{g_y}^2 N + K_z \left( \frac{8 B_{f_y}^2}{\mu_g^2 L_{\mu_q}} \sigma_{g_{yy}}^2 c_\gamma \alpha_t + \frac{4 c_\gamma}{L_{\mu_q}} \sigma_{f_y}^2 \alpha_t + \frac{8}{c_\gamma L_{\mu_q}} L_z^2 c_\beta^2 \alpha_t N^2 \sigma_{g_y}^2 \right) \end{split}$$

Therefore, the lemma is proved. 

#### C.2.5 PROOF OF THEOREM 5.3

**Theorem C.9** (Non-Convex  $\ell(\mathbf{x})$ ). Under Assumptions 3.1–3.4, choose step-sizes  $\alpha_t = \alpha$ ,  $\beta_t \triangleq c_{\beta}\alpha$ , and  $\gamma_t \triangleq c_{\gamma} \alpha$  for all  $t \in \{0, 1, \dots, T\}$  with 

$$c_{\beta} = \frac{12L_{y}^{2}K_{y}}{L_{\mu_{g}}}, \quad c_{\gamma} = \max\left\{\frac{24L_{z}^{2}NK_{z}}{L_{\mu_{q}}}, \frac{192\left(\mu_{g} + L_{g}\right)L_{z}^{2}NK_{z}L_{y}^{2}}{L_{\mu_{g}}L_{\mu_{q}}}\right\},\tag{22}$$

where 

$$K_{y} = \frac{L_{f}}{\sqrt{3}L_{y}}, \quad K_{z} = \max\left\{\frac{24\sigma_{g_{xy}}^{2}N}{c_{\gamma}L_{\mu_{q}}}, \frac{24L_{f}^{2}N}{c_{\gamma}L_{\mu_{q}}}\right\}, \quad L_{\mu_{g}} = \frac{\mu_{g}L_{g}}{\mu_{g} + L_{g}}, \quad L_{\mu_{q}} = \frac{\mu_{g}L_{q}}{\mu_{g} + L_{q}}.$$
 (23)

Moreover, choose  $\alpha$  such that

$$\begin{split} \alpha &\leq \min\left\{\frac{1}{6L_{l}}, \frac{1}{c_{\beta}\left(\mu_{g} + L_{g}\right)}, \frac{1}{c_{\gamma}\left(\mu_{g} + L_{q}\right)}, \frac{1}{2L_{\mu_{g}}c_{\beta}}, \frac{1}{2L_{\mu_{q}}c_{\gamma}}, \frac{L_{l}}{16L_{f}^{2}L_{z}^{2}N^{2}} \right. \\ &\left. \sqrt{\frac{K_{y}}{32\left(\mu_{g} + L_{g}\right)L_{f}^{2}L_{z}^{2}N^{2}c_{\beta}}}\right\}. \end{split}$$

Then, the iterates generated by LazyBLO satisfy: 

$$\frac{1}{TN} \sum_{t=0}^{T-1} \sum_{n=0}^{N-1} \mathbb{E}\left[ \|\nabla \ell\left(\mathbf{x}_{t}^{n}\right)\|^{2} \right] = \mathcal{O}\left(\frac{\Delta_{0}}{NT\alpha}\right) + \mathcal{O}\left(\sigma_{g_{y}}^{2} + \sigma_{g_{xy}}^{2} + \sigma_{f_{x}}^{2} + \sigma_{g_{yy}}^{2} + \sigma_{f_{y}}^{2}\right),$$
  
where  $\Delta_{0} = (\ell(\mathbf{x}_{0}^{0}) - \ell^{*}) + \|\mathbf{y}_{0}^{0} - \mathbf{y}^{*}(\mathbf{x}_{0}^{0})\|^{2} + \|\mathbf{z}_{0} - \mathbf{z}^{*}(\mathbf{x}_{0}^{0}, \mathbf{y}_{0}^{0})\|^{2}.$ 

> *Proof.* Choose  $\alpha_t$  as a constant stepsize  $\alpha_t = \alpha$ . Summing the result in Lemma C.4 from t = 0 to T-1, and then dividing by NT on both sides, we get

$$\frac{\mathbb{E}\left[W_{T} - W_{0}\right]}{NT} \leq -\frac{\alpha}{2TN} \sum_{t=0}^{T-1} \sum_{n=0}^{N-1} \mathbb{E}\left[\|\nabla \ell\left(\mathbf{x}_{t}^{n}\right)\|^{2}\right] + \frac{\alpha}{N} \left(\sigma_{f_{x}}^{2} C_{f_{x}} N + \sigma_{f_{y}}^{2} C_{f_{y}} N + \sigma_{g_{yy}}^{2} C_{g_{yy}} N + \sigma_{g_{yy}}^{$$

Rearranging the terms and multiplying by  $2/\alpha$  on both sides, we have

 $\frac{1}{TN}\sum_{i=1}^{T-1}\sum_{i=1}^{N-1}\mathbb{E}\left[\left\|\nabla \ell\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right]$  $\leq \frac{2\mathbb{E}\left[W_{0}-\ell^{*}\right]}{\alpha NT}+2\left(\sigma_{f_{x}}^{2}C_{f_{x}}+\sigma_{f_{y}}^{2}C_{f_{y}}+\sigma_{g_{y}}^{2}\left(C_{g_{1}}+C_{g_{2}}\frac{1}{N^{2}}\right)+\sigma_{g_{yy}}^{2}C_{g_{yy}}+\sigma_{g_{xy}}^{2}C_{g_{xy}}\right)$  $\leq \frac{2\left(W_{0}-\ell^{*}\right)}{\alpha NT}+2\left(\sigma_{f_{x}}^{2}C_{f_{x}}+\sigma_{f_{y}}^{2}C_{f_{y}}+\sigma_{g_{y}}^{2}\left(C_{g_{1}}+C_{g_{2}}\frac{1}{N^{2}}\right)+\sigma_{g_{yy}}^{2}C_{g_{yy}}+\sigma_{g_{xy}}^{2}C_{g_{xy}}\right),$ 

where  $W_0 = \ell(\mathbf{x}_0^0) + K_y \|\mathbf{y}_0^0 - \mathbf{y}^*(x_0^0)\|^2 + K_z \|\mathbf{z}_0 - \mathbf{z}^*(\mathbf{x}_0^0, \mathbf{y}_0^0)\|^2$ .

Therefore, 1567  $\frac{1}{TN} \sum_{i=0}^{T-1} \sum_{n=0}^{N-1} \mathbb{E}\left[ \|\nabla \ell \left( \mathbf{x}_{t}^{n} \right)\|^{2} \right]$ 1568 1569 1570  $\mathcal{L} = \mathcal{O}\left(rac{\ell\left(\mathbf{x}_{0}^{0}
ight) - \ell^{*}}{NTlpha}
ight) + \mathcal{O}\left(rac{\left\|\mathbf{y}_{0}^{0} - \mathbf{y}^{*}\left(\mathbf{x}_{0}^{0}
ight)
ight\|^{2}}{NTlpha}
ight) + \mathcal{O}\left(rac{\left\|\mathbf{z}_{0} - \mathbf{z}^{*}\left(\mathbf{x}_{0}^{0}, \mathbf{y}_{0}^{0}
ight)
ight\|^{2}}{NTlpha}
ight)$ 1571 1572  $+\mathcal{O}\left(\sigma_{q_{xy}}^2+\sigma_{f_x}^2+\sigma_{q_{yy}}^2+\sigma_{f_y}^2+\sigma_{q_y}^2\right)$ 1574 1575 The proof of the theorem is completed. 1576 1577 1578 1579 **PROOF OF THEOREM 5.6:** STRONGLY-CONVEX  $\ell(\mathbf{x})$ D 1580 1581 D.1 DESCENT IN THE UPPER-LEVEL OBJECTIVE FUNCTION 1582 **Lemma D.1.** Under Assumptions 3.1–3.4. For strongly-convex and smooth  $\ell(\mathbf{x})$ , the following inequality holds for successive iterations of Algorithm 1: 1585  $\mathbb{E}\left[\ell\left(\mathbf{x}_{t}^{n+1}\right)-\ell^{*}\right]$ 1586 1587  $\leq (1 - \mu_f \alpha_t) \mathbb{E}\left[\ell\left(\mathbf{x}_t^n\right) - \ell^*\right] - \left(\frac{\alpha_t}{2} - \frac{\alpha_t^2 L_l}{2}\right) \mathbb{E}\left[\left\|h_{t,n}^f\right\|^2\right] + 8L_f^2 L_z^2 \alpha_t^3 N \sum_{n=1}^{N-1} \mathbb{E}\left[\left\|h_{t,n}^f\right\|^2\right]$ 1589  $+2\alpha_{t}L_{f}^{2}\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right]+\left(4\sigma_{g_{xy}}^{2}\alpha_{t}+4L_{f}^{2}\alpha_{t}\right)\mathbb{E}\left[\left\|\mathbf{z}_{t}-\mathbf{z}_{t}^{*}\right\|^{2}\right]+16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N^{2}\sigma_{g_{y}}^{2}\alpha_{t}$ 1591 +  $16L_f^2 L_z^2 \beta_t^2 N \alpha_t \sum_{k=1}^{N-1} \mathbb{E}\left[ \left\| \nabla_{\mathbf{y}} g\left( \mathbf{x}_t^n, \mathbf{y}_t^n \right) \right\|^2 \right] + 2\sigma_{f_x}^2 \alpha_t + 4\sigma_{g_{xy}}^2 \frac{B_{f_y}^2}{\mu_z^2} \alpha_t,$ 1592 1593 1594 for all  $t \in \{0, 1, \dots, T-1\}$  and  $n \in \{0, 1, \dots, N-1\}$ , where the expectation is taken over the 1595 stochasticity of the algorithm. 1596 1597 *Proof.* From Lemma C.1, we have 1598 1599  $\mathbb{E}\left[\ell\left(\mathbf{x}_{t}^{n+1}\right)-\ell\left(\mathbf{x}_{t}^{n}\right)\right]$  $\leq -\frac{\alpha_t}{2} \mathbb{E}\left[\left\|\nabla \ell\left(\mathbf{x}_t^n\right)\right\|^2\right] - \left(\frac{\alpha_t}{2} - \frac{\alpha_t^2 L_l}{2}\right) \mathbb{E}\left[\left\|h_{t,n}^f\right\|^2\right] + 8L_f^2 L_z^2 \alpha_t^3 N \sum_{l=1}^{N-1} \mathbb{E}\left[\left\|h_{t,n}^f\right\|^2\right]$  $+2\alpha_{t}L_{f}^{2}\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right]+\left(4\sigma_{g_{xy}}^{2}\alpha_{t}+4L_{f}^{2}\alpha_{t}\right)\mathbb{E}\left[\left\|\mathbf{z}_{t}-\mathbf{z}_{t}^{*}\right\|^{2}\right]+2\sigma_{f_{x}}^{2}\alpha_{t}$ 1604  $+16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N\alpha_{t}\sum_{\alpha}^{N-1}\mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)\right\|^{2}\right]+16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}N^{2}\sigma_{g_{y}}^{2}\alpha_{t}+4\sigma_{g_{xy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{a}^{2}}\alpha_{t}.$ 1608 For a strongly convex function  $\ell(\mathbf{x})$ , we have the fact that for all  $\mathbf{x} \in \mathbb{R}^{u}$ , 1609 1610  $\|\nabla \ell(\mathbf{x})\|^2 > 2\mu_f \left(\ell(\mathbf{x}) - \ell^*\right).$ 1611 1612 Substitute (25) in (24) and subtract  $\ell^*$  from both sides. After rearranging the terms, the lemma is 1613 proved. 1614 1615 D.2 DESCENT IN THE ERROR OF  $y^*(x)$ 

(24)

(25)

1617 **Lemma D.2.** Under Assumptions 3.2–3.4, the approximation error of  $\mathbf{y}^*(\mathbf{x})$  of Algorithm 1 satisfies the following inequality: 1618

$$\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n+1}-\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n+1}\right)\right\|^{2}\right]$$

1616

1620  
1621 
$$\leq (1+c_1) (1-2\beta_t \mu_g) \mathbb{E} \left[ \left\| \mathbf{y}_t^n - \mathbf{y}^* \left( \mathbf{x}_t^n \right) \right\|^2 \right] + 2\beta_t^2 (1+c_1) \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_t^n, \mathbf{y}_t^n \right) \right\|^2 \right]$$

1625

1626 1627

 $+\left(1+\frac{1}{c_1}\right)L_y^2\alpha_t^2\mathbb{E}\left[\left\|h_{t,n}^f\right\|^2\right]+2\left(1+c_1\right)\beta_t^2\sigma_{g_y}^2,$ 

for all  $t \in \{0, 1, ..., T-1\}$  and  $n \in \{0, 1, ..., N-1\}$  with a constant  $c_1 > 0$ , where the expectation is taken over the stochasticity of the algorithm.

1628 Proof.

1629 1630 1631

1634

$$\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n+1} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n+1}\right)\right\|^{2}\right] \\
\stackrel{(a)}{\leq} \mathbb{E}\left[\left(1+c_{1}\right)\left\|\mathbf{y}_{t}^{n+1} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2} + \left(1+\frac{1}{c_{1}}\right)\left\|\mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right) - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n+1}\right)\right\|^{2}\right] \\
\stackrel{(b)}{\leq} \mathbb{E}\left[\left(1+c_{1}\right)\left\|\mathbf{y}_{t}^{n} - \beta_{t}h_{t,n}^{g} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2} + \left(1+\frac{1}{c_{1}}\right)L_{y}^{2}\left\|\mathbf{x}_{t}^{n+1} - \mathbf{x}_{t}^{n}\right\|^{2}\right] \\
\stackrel{(c)}{=} \mathbb{E}\left[\left(1+c_{1}\right)\left\|\mathbf{y}_{t}^{n} - \beta_{t}h_{t,n}^{g} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2} + \left(1+\frac{1}{c_{1}}\right)L_{y}^{2}\alpha_{t}^{2}\left\|h_{t,n}^{f}\right\|^{2}\right], \quad (26)$$

1635 1636 1637

where (a) results from Young's inequality. (b) is because of the update rule of Algorithm 1 and the Lipschitzness of  $y^*(\cdot)$  (see Lemma 5.1). (c) follows from the update rule of Algorithm 1.

1640 Next, we bound the first term of the above inequality.

$$\begin{aligned}
\mathbf{E} \begin{bmatrix} \left\| \mathbf{y}_{t}^{n} - \beta_{t} h_{t,n}^{g} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} \end{bmatrix} \\
= \mathbb{E} \begin{bmatrix} \left\| \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} \end{bmatrix} + \beta_{t}^{2} \mathbb{E} \begin{bmatrix} \left\| h_{t,n}^{g} \right\|^{2} \end{bmatrix} - 2\beta_{t} \mathbb{E} \left[ \left\langle h_{t,n}^{g}, \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\rangle \right] \\
= \mathbb{E} \begin{bmatrix} \left\| \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} \end{bmatrix} + 2\beta_{t}^{2} \mathbb{E} \begin{bmatrix} \left\| h_{t,n}^{g} - \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} \end{bmatrix} + 2\beta_{t}^{2} \mathbb{E} \begin{bmatrix} \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} \end{bmatrix} \\
= 2\beta_{t} \mathbb{E} \left[ \left\langle h_{t,n}^{g}, \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\rangle \end{bmatrix} \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} \right] + 2\beta_{t}^{2} \mathbb{E} \left[ \left\| h_{t,n}^{g} - \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} \right] + 2\beta_{t}^{2} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} \right] + 2\beta_{t}^{2} \mathbb{E} \left[ \left\| h_{t,n}^{g} - \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right\|^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right), \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right), \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right), \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right]^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right]^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right]^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right]^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right]^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right]^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right]^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right]^{2} \\
= 2\beta_{t} \mathbb{E} \left[ \left\| \nabla_{\mathbf{y}} g \left( \mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n} \right) \right]^{2} \right] \\
= 2\beta_{t} \mathbb{E} \left[ \left\|$$

1657 where (a) uses the fact that  $\mathbb{E}\left[h_{t,n}^{g}|\mathcal{F}_{t}^{n}\right] = \nabla_{\mathbf{y}}g(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n})$ , and  $\mathcal{F}_{t}^{n} \triangleq \sigma\left\{\mathbf{y}_{0}^{0},\mathbf{x}_{0}^{0},\cdots,\mathbf{y}_{t}^{n},\mathbf{x}_{t}^{n}\right\}$ 1658 is defined as the sigma algebra generated by the iteration sequence of Algorithm 1. (b) utilizes 1659 the fact that for  $\mu_{g}$ -strongly convex  $g(\mathbf{x},\mathbf{y})$ , we have  $\langle \nabla_{\mathbf{y}}g(\mathbf{x},\mathbf{y}_{1}) - \nabla_{\mathbf{y}}g(\mathbf{x},\mathbf{y}_{2}),\mathbf{y}_{1}-\mathbf{y}_{2}\rangle \geq$ 1660  $\mu_{g} \|\mathbf{y}_{1}-\mathbf{y}_{2}\|^{2}$ . (c) is because of the bounded variance in Assumption 3.4.

Substituting (27) in (26) yields the lemma.

 $\mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t+1}^*\right\|^2\right]$ 

# 1663 D.3 DESCENT IN THE ERROR OF $z^{*}(x, y)$

**Lemma D.3.** Under Assumptions 3.1–3.4, the following inequality of the approximation error of  $\mathbf{z}^*(\mathbf{x}, \mathbf{y})$  holds for Algorithm 1:

$$\leq (1+c_3) \left(1 - 2\gamma_t \mu_g + 8\sigma_{g_{yy}}^2 \gamma_t^2\right) \mathbb{E}\left[\|\mathbf{z}_t - \mathbf{z}_t^*\|^2\right] + \left(2 + \frac{2}{c_3}\right) L_z^2 \alpha_t^2 N \sum_{n=0}^{N-1} \mathbb{E}\left[\left\|h_{t,n}^f\right\|^2\right]$$

1672  
1673 
$$+ 2\gamma_t^2 (1+c_3) \mathbb{E}\left[ \left\| \nabla_{\mathbf{z}} q\left( \mathbf{x}_t^N, \mathbf{y}_t^N, \mathbf{z}_t \right) \right\|^2 \right] + 4\left( 1 + \frac{1}{c_3} \right) L_z^2 \beta_t^2 N \sum_{n=0}^{N-1} \mathbb{E}\left[ \left\| \nabla_{\mathbf{y}} g\left( \mathbf{x}_t^n, \mathbf{y}_t^n \right) \right\|^2 \right]$$

$$\begin{array}{l} {}^{1674}\\ {}^{1675}\\ {}^{1676}\end{array} + 4\left(1+\frac{1}{c_3}\right)L_z^2\beta_t^2N^2\sigma_{g_y}^2 + 4\sigma_{f_y}^2\left(1+c_3\right)\gamma_t^2 + 8\sigma_{g_{yy}}^2\frac{B_{f_y}^2}{\mu_g^2}\left(1+c_3\right)\gamma_t^2, \end{array}$$

*for all*  $t \in \{0, 1, ..., T-1\}$  *and*  $n \in \{0, 1, ..., N-1\}$  *with some constants*  $c_3, c_4 > 0$ *, where the* expectation is taken over the stochasticity of the algorithm.

*Proof.* With the results from the proof of Lemma C.7, we have

$$\mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t+1}^{*}\right\|^{2}\right] \leq (1+c_{3}) \mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t}^{*}\right\|^{2}\right] + 2\left(1+\frac{1}{c_{3}}\right) L_{z}^{2} \alpha_{t}^{2} N \sum_{n=0}^{N-1} \mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] + 4\left(1+\frac{1}{c_{3}}\right) L_{z}^{2} \beta_{t}^{2} N \sum_{n=0}^{N-1} \mathbb{E}\left[\left\|\nabla_{\mathbf{y}}g\left(\mathbf{x}_{t}^{n}, \mathbf{y}_{t}^{n}\right)\right\|^{2}\right] + 4\left(1+\frac{1}{c_{3}}\right) L_{z}^{2} \beta_{t}^{2} N^{2} \sigma_{g_{y}}^{2},$$
(28)

Then, we consider the first term on the right:

where (a) is due to the update rule of Algorithm 1. (b) follows from the fact that  $\mathbb{E}[h_t^q|\mathcal{F}_t] =$  $\nabla_{\mathbf{z}}q(\mathbf{x}_{t}^{N},\mathbf{y}_{t}^{N},\mathbf{z}_{t})$ , and (c) utilizes the fact that for  $\mu_{q}$ -strongly convex  $q(\mathbf{x},\mathbf{y},\mathbf{z})$ , we have  $\langle \nabla_{\mathbf{z}} q\left(\mathbf{x}, \mathbf{y}, \mathbf{z}_{1}\right) - \nabla_{\mathbf{z}} q\left(\mathbf{x}, \mathbf{y}, \mathbf{z}_{2}\right), \mathbf{z}_{1} - \mathbf{z}_{2} \rangle \geq \mu_{g} \|\mathbf{z}_{1} - \mathbf{z}_{2}\|^{2}.$ 

Next, we consider the second term  $\mathbb{E}\left[\left\|h_{t}^{q}-\nabla_{\mathbf{z}}q\left(\mathbf{x}_{t}^{N},\mathbf{y}_{t}^{N},\mathbf{z}_{t}\right)\right\|^{2}\right]$  of the above inequality: 

where (a) results from the definitions of  $h_t^q$  and  $\nabla_z q(\mathbf{x}_t^N, \mathbf{y}_t^N, \mathbf{z}_t)$ . (b) uses the bounded variance in Assumption 3.4. (c) is because of the bound of  $z^*(x, y)$  in Lemma 5.2. 

Substituting (30) into (29) and then substituting the obtained inequality into (28) proves the lemma. 

#### D.4 DESCENT IN THE POTENTIAL FUNCTION

We define a different potential function  $\hat{W}_t$  as follows:

$$\hat{W}_{t} = \sum_{n=0}^{N-1} \left( \ell \left( \mathbf{x}_{t}^{n} \right) - \ell^{*} \right) + \sum_{n=0}^{N-1} \left\| \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} + \left\| \mathbf{z}_{t} - \mathbf{z}_{t}^{*} \right\|^{2}.$$

**Lemma D.4.** Choose  $c_1 = \frac{\beta_t \mu_g}{2(1-\beta_t \mu_g)}$ , and  $c_3 = \frac{\gamma_t \mu_g}{2(1-\gamma_t \mu_g)}$ . Under the same conditions as described in Theorem D.5 and utilizing Lemmas B.D.1-B.D.3, the iterates generated by Algorithm 1 satisfies: 

$$\begin{split} & \sum_{\substack{1738\\1739}} & \mathbb{E}\left[\hat{W}_{t+1}\right] \leq (1-\mu_f\alpha_t) \, \mathbb{E}\left[\hat{W}_t\right] + 16L_f^2 L_z^2 \hat{c}_{\beta}^2 \alpha_t^3 N^3 \sigma_{g_y}^2 + 4\sigma_{g_{xy}}^2 \frac{B_{f_y}^2}{\mu_g^2} \alpha_t N + 2\sigma_{f_x}^2 \alpha_t N + 4\hat{c}_{\beta}^2 \sigma_{g_y}^2 N \alpha_t^2 \\ & + 8\sigma_{f_y}^2 \hat{c}_{\gamma}^2 \alpha_t^2 + 16\sigma_{g_{yy}}^2 \frac{B_{f_y}^2}{\mu_g^2} \hat{c}_{\gamma}^2 \alpha_t^2 + \frac{8}{\mu_g \hat{c}_{\gamma}} L_z^2 \hat{c}_{\beta}^2 N^2 \sigma_{g_y}^2 \alpha_t, \end{split}$$

for all  $t \in \{0, 1, \dots, T-1\}$ .

*Proof.* With the results from Lemma D.1, we have 

$$\begin{aligned} & \sum_{n=0}^{N-1} \mathbb{E}\left[\ell\left(\mathbf{x}_{t}^{n+1}\right) - \ell^{*}\right] \\ & \leq \mathbb{E}\left[\left(1 - \mu_{f}\alpha_{t}\right)\sum_{n=0}^{N-1}\left(\ell\left(\mathbf{x}_{t}^{n}\right) - \ell^{*}\right) + \left(\frac{\alpha_{t}^{2}L_{l}}{2} - \frac{\alpha_{t}}{2} + 8L_{f}^{2}L_{z}^{2}\alpha_{t}^{3}N^{2}\right)\sum_{n=0}^{N-1}\left\|h_{t,n}^{f}\right\|^{2} \\ & + \left(4\sigma_{g_{xy}}^{2}\alpha_{t}N + 4L_{f}^{2}\alpha_{t}N\right)\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2} + 16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}\alpha_{t}N^{2}\sum_{n=0}^{N-1}\left\|\nabla\mathbf{y}g\left(\mathbf{x}_{t}^{n},\mathbf{y}_{t}^{n}\right)\right\|^{2} \\ & + 2L_{f}^{2}\alpha_{t}\sum_{n=0}^{N-1}\left\|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] + 16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}\alpha_{t}N^{3}\sigma_{g_{y}}^{2} + 2\sigma_{f_{x}}^{2}\alpha_{t}N + 4\sigma_{g_{xy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\alpha_{t}N \\ & + 2L_{f}^{2}\alpha_{t}\sum_{n=0}^{N-1}\left\|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] + 16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}\alpha_{t}N^{3}\sigma_{g_{y}}^{2} + 2\sigma_{f_{x}}^{2}\alpha_{t}N + 4\sigma_{g_{xy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\alpha_{t}N \\ & \leq \mathbb{E}\left[\left(1 - \mu_{f}\alpha_{t}\right)\sum_{n=0}^{N-1}\left(\ell\left(\mathbf{x}_{t}^{n}\right) - \ell^{*}\right) + \left(\frac{\alpha_{t}^{2}L_{l}}{2} - \frac{\alpha_{t}}{2} + 8L_{f}^{2}L_{z}^{2}\alpha_{t}^{3}N^{2}\right)\sum_{n=0}^{N-1}\left\|h_{t,n}^{f}\right\|^{2} \\ & + \left(4\sigma_{g_{xy}}^{2}\alpha_{t}N + 4L_{f}^{2}\alpha_{t}N\right)\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2} + \left(2L_{f}^{2}\alpha_{t} + 16L_{f}^{2}L_{z}^{2}\beta_{t}^{2}\alpha_{t}N^{2}L_{g}^{2}\right)\sum_{n=0}^{N-1}\left\|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2} \right] \\ & + \left(4\sigma_{g_{xy}}^{2}\alpha_{t}N + 4L_{f}^{2}\alpha_{t}N\right)\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2} + \left(2L_{f}^{2}\alpha_{t} + 16L_{f}^{2}L_{z}^{2}\alpha_{t}^{3}N^{2}\right)\sum_{n=0}^{N-1}\left\|\mathbf{h}_{t,n}^{f}\right\|^{2} \\ & + \left(4\sigma_{g_{xy}}^{2}\alpha_{t}N + 4L_{f}^{2}\alpha_{t}N\right)\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2} + \left(2L_{f}^{2}\alpha_{t} + 16L_{f}^{2}L_{z}^{2}\alpha_{t}^{3}N^{2}L_{g}^{2}\right)\sum_{n=0}^{N-1}\left\|\mathbf{h}_{t,n}^{f}\right\|^{2} \\ & + \left(4\sigma_{g_{xy}}^{2}\alpha_{t}N + 4L_{f}^{2}\alpha_{t}N\right)\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2} + \left(2L_{f}^{2}\alpha_{t} + 16L_{f}^{2}L_{z}^{2}\alpha_{t}^{3}N^{2}L_{g}^{2}\right)\sum_{n=0}^{N-1}\left\|\mathbf{h}_{t,n}^{f}\right\|^{2} \\ & + \left(4\sigma_{g_{xy}}^{2}\alpha_{t}N + 4L_{f}^{2}\alpha_{t}N\right)\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2} + \left(2L_{f}^{2}\alpha_{t} + 16L_{f}^{2}L_{z}^{2}\alpha_{t}^{3}N^{2}L_{g}^{2}\right)\sum_{n=0}^{N-1}\left\|\mathbf{h}_{t,n}^{f}\right\|^{2} \\ & + \left(4\sigma_{g_{xy}}^{2}\alpha_{t}N + 4L_{f}^{2}\alpha_{t}N\right)\left\|\mathbf{z}_{t} - \mathbf{z}_{t$$

where (a) uses the fact that  $\nabla_{\mathbf{y}} g(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) = 0$  and utilizes the the Lipschitzness of  $\nabla_{\mathbf{y}} g(\mathbf{x}, \mathbf{y})$  (see Assumption 3.2). (b) follows from the definition of  $\beta_t = \hat{c}_{\beta} \alpha_t$ .

From Lemma D.2, we have 

$$\sum_{n=0}^{N-1} \mathbb{E}\left[\left\|\mathbf{y}_t^{n+1} - \mathbf{y}^*\left(\mathbf{x}_t^{n+1}\right)\right\|^2\right]$$

where (a) is because of the fact that  $\nabla_{\mathbf{y}} g(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) = 0$  and follows from the the Lipschitzness of  $\nabla_{\mathbf{y}} g(\mathbf{x}, \mathbf{y})$  (see Assumption 3.2). 

From the choice of  $c_1 = \frac{\beta_t \mu_g}{2(1-\beta_t \mu_g)}$ , we have  $1 + \frac{1}{c_1} \leq \frac{2}{\mu_g \beta_t}$ . Choosing  $\beta_t \leq \frac{\mu_g}{2L_g^2}$  and using the definition of  $\beta_t = \hat{c}_\beta \alpha_t$ , we get 

$$\sum_{n=0}^{N-1} \mathbb{E} \left[ \left\| \mathbf{y}_{t}^{n+1} - \mathbf{y}^{*} \left( x_{t}^{n+1} \right) \right\|^{2} \right] \\ \leq \mathbb{E} \left[ \left( 1 - \frac{\mu_{g} \hat{c}_{\beta} \alpha_{t}}{2} \right) \sum_{n=0}^{N-1} \left\| \mathbf{y}_{t}^{n} - \mathbf{y}^{*} \left( \mathbf{x}_{t}^{n} \right) \right\|^{2} + \frac{2}{\mu_{g} \hat{c}_{\beta}} L_{y}^{2} \alpha_{t} \sum_{n=0}^{N-1} \left\| h_{t,n}^{f} \right\|^{2} \right] + 4 \hat{c}_{\beta}^{2} \sigma_{g_{y}}^{2} N \alpha_{t}^{2}.$$
(32)

Following from Lemma D.3, we have

$$\begin{split} & \mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t+1}^{*}\right\|^{2}\right] \\ & \mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t+1}^{*}\right\|^{2}\right] \\ & \mathbb{E}\left[\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2}\right] + 2\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\alpha_{t}^{2}N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] \\ & \mathbb{E}\left[\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2}\right] + 2\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\alpha_{t}^{2}N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|h_{t,n}^{f}\right\|^{2}\right] \\ & \mathbb{E}\left[\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2}\right] + 4\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\alpha_{t}^{2}N\sum_{n=0}^{N-1}\mathbb{E}\left[\left\|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] + 4\left(1 + \frac{1}{c_{3}}\right)L_{z}^{2}\beta_{t}^{2}N^{2}\sigma_{g_{y}}^{2} + 4\sigma_{f_{y}}^{2}\left(1 + c_{3}\right)\gamma_{t}^{2} \\ & + 8\sigma_{g_{yy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\left(1 + c_{3}\right)\gamma_{t}^{2}, \end{split}$$

where (a) utilizes fact that  $\nabla_{\mathbf{y}} g(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) = 0$  and  $\nabla_{\mathbf{z}} q(\mathbf{x}, \mathbf{y}, \mathbf{z}^*) = 0$ . In addition, it uses the the Lipschitzness of  $\nabla_{\mathbf{y}} g(\mathbf{x}, \mathbf{y})$  in Assumption 3.2 and  $\nabla_{\mathbf{z}} q(\mathbf{x}, \mathbf{y}, \mathbf{z})$  proved as follows. 

$$\begin{aligned} &\|\nabla_{\mathbf{z}}q\left(\mathbf{x},\mathbf{y},\mathbf{z}_{1}\right) - \nabla_{\mathbf{z}}q\left(\mathbf{x},\mathbf{y},\mathbf{z}_{2}\right)\| \stackrel{(a)}{=} \left\|\nabla_{\mathbf{yy}}^{2}g(\mathbf{x},\mathbf{y})\mathbf{z}_{1} + \nabla_{\mathbf{y}}f(\mathbf{x},\mathbf{y}) - \nabla_{\mathbf{yy}}^{2}g(\mathbf{x},\mathbf{y})\mathbf{z}_{2} - \nabla_{\mathbf{y}}f(\mathbf{x},\mathbf{y})\right\| \\ &= \left\|\nabla_{\mathbf{yy}}^{2}g(\mathbf{x},\mathbf{y})\right\| \left\|\mathbf{z}_{1} - \mathbf{z}_{2}\right\| \stackrel{(b)}{\leq} B_{g_{yy}} \left\|\mathbf{z}_{1} - \mathbf{z}_{2}\right\| \stackrel{(c)}{=} L_{q} \left\|\mathbf{z}_{1} - \mathbf{z}_{2}\right\|, \end{aligned}$$

where (a) follows from the definition of  $\nabla_{\mathbf{z}} q(\mathbf{x}, \mathbf{y}, \mathbf{z})$ . (b) assumes  $\|\nabla_{yy}^2 g(x, y)\| \leq B_{g_{yy}}$ , and (c) defines  $L_q = B_{g_{yy}}$ . 

From the choice of  $c_3 = \frac{\gamma_t \mu_g}{2(1-\gamma_t \mu_g)}$ , we get  $1 + \frac{1}{c_3} \leq \frac{2}{\mu_g \gamma_t}$ . Selecting  $\gamma_t \leq \frac{\mu_g}{4L_q^2}$ ,  $\gamma_t \leq \frac{\mu_g}{16\sigma_{g_{u_u}}^2}$  and using the definition of  $\beta_t = \hat{c}_{\beta} \alpha_t$ ,  $\gamma_t = \hat{c}_{\gamma} \alpha_t$ , we have

$$\mathbb{E}\left[\left\|\mathbf{z}_{t+1} - \mathbf{z}_{t+1}^{*}\right\|^{2}\right] \leq \mathbb{E}\left[\left(1 - \frac{\mu_{g}\hat{c}_{\gamma}\alpha_{t}}{2}\right)\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2} + \frac{4}{\mu_{g}\hat{c}_{\gamma}}L_{z}^{2}\alpha_{t}N\sum_{n=0}^{N-1}\left\|h_{t,n}^{f}\right\|^{2} + \frac{8}{\mu_{g}\hat{c}_{\gamma}}L_{z}^{2}\hat{c}_{\beta}^{2}NL_{g}^{2}\alpha_{t}\sum_{n=0}^{N-1}\left\|\mathbf{y}_{t}^{n} - \mathbf{y}^{*}\left(\mathbf{x}_{t}^{n}\right)\right\|^{2}\right] + 16\sigma_{g_{yy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\hat{c}_{\gamma}^{2}\alpha_{t}^{2} + 8\sigma_{f_{y}}^{2}\hat{c}_{\gamma}^{2}\alpha_{t}^{2} + \frac{8}{\mu_{g}\hat{c}_{\gamma}}L_{z}^{2}\hat{c}_{\beta}^{2}N^{2}\sigma_{g_{y}}^{2}\alpha_{t}.$$
(33)

Combining equations (31), (32) and (33), we get

$$\mathbb{E}\left[\hat{W}_{t+1}\right] \le (1 - \mu_f \alpha_t) \mathbb{E}\left[\hat{W}_t\right] + \hat{C}_h \mathbb{E}\left[\sum_{n=0}^{N-1} \left\|h_{t,n}^f\right\|^2\right] + \hat{C}_y \mathbb{E}\left[\sum_{n=0}^{N-1} \left\|\mathbf{y}_t^n - \mathbf{y}^*\left(\mathbf{x}_t^n\right)\right\|^2\right]$$

$$+ \hat{C}_{z}\mathbb{E}\left[\left\|\mathbf{z}_{t} - \mathbf{z}_{t}^{*}\right\|^{2}\right] + 16L_{f}^{2}L_{z}^{2}\hat{c}_{\beta}^{2}\alpha_{t}^{3}N^{3}\sigma_{g_{y}}^{2} + 4\sigma_{g_{xy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\alpha_{t}N + 2\sigma_{f_{x}}^{2}\alpha_{t}N + 4\hat{c}_{\beta}^{2}\sigma_{g_{y}}^{2}N\alpha_{t}^{2}$$

Therefore, the lemma is proved. 

#### D.5 PROOF OF THEOREM 5.6

**Theorem D.5** (Strongly Convex  $\ell(\mathbf{x})$ ). Suppose the upper-level function  $\ell(\mathbf{x})$  is  $\mu_f$ -strongly-convex. Under Assumptions 3.1–3.4, choose the step-sizes  $\alpha_t = \alpha$ ,  $\beta_t \triangleq \hat{c}_{\beta}\alpha$  and  $\gamma_t \triangleq \hat{c}_{\gamma}\alpha$  for all 

1891  
1892  
1893  
1894  

$$\hat{c}_{\beta} = \max\left\{\frac{16L_{y}^{2}}{\mu_{g}}, \frac{6\left(2L_{f}^{2} + \mu_{f}\right)}{\mu_{g}}\right\}, \ \hat{c}_{\gamma} = \max\left\{\frac{32L_{z}^{2}}{\mu_{g}}, \frac{\mu_{g}^{2}}{48L_{z}^{2}L_{g}^{2}N\hat{c}_{\beta}}, \frac{8\sigma_{g_{xy}}^{2}N + 8L_{f}^{2}N + 2\mu_{f}}{\mu_{g}}\right\}$$
(34)

Moreover, choose  $\alpha$  such that

 $t \in \{0, 1, \dots, T-1\}$ , where

$$\alpha \le \min\left\{\frac{1}{4L_l}, \frac{1}{8L_f L_z N}, \sqrt{\frac{\mu_g}{96L_f^2 L_z^2 L_g^2 N^2 \hat{c}_\beta}}, \frac{\mu_g}{2L_g^2 \hat{c}_\beta}, \frac{2}{3\mu_g \hat{c}_\beta}, \frac{\mu_g}{16\sigma_{g_{yy}}^2 \hat{c}_\gamma}, \frac{\mu_g}{4L_q^2 \hat{c}_\gamma}, \frac{2}{3\mu_g \hat{c}_\gamma}\right\}.$$

Then, the iterates generated by LazyBLO satisfy:

$$\begin{split} \sum_{n=0}^{N-1} \mathbb{E} \bigg[ \ell \left( \mathbf{x}_{t}^{n} \right) - \ell^{*} \bigg] &\leq \left( 1 - \mu_{f} \alpha \right)^{t} \hat{\Delta}_{0} + \frac{1}{\mu_{f}} \bigg( 4 \sigma_{g_{xy}}^{2} \frac{B_{f_{y}}^{2}}{\mu_{g}^{2}} N + 2 \sigma_{f_{x}}^{2} N + \frac{8}{\mu_{g} \hat{c}_{\gamma}} L_{z}^{2} \hat{c}_{\beta}^{2} N^{2} \sigma_{g_{y}}^{2} \bigg) \\ &+ \frac{\alpha}{\mu_{f}} \bigg( 16 \sigma_{g_{yy}}^{2} \frac{B_{f_{y}}^{2}}{\mu_{g}^{2}} \hat{c}_{\gamma}^{2} + 8 \sigma_{f_{y}}^{2} \hat{c}_{\gamma}^{2} + 4 \hat{c}_{\beta}^{2} \sigma_{g_{y}}^{2} N \bigg) + \frac{16 \alpha^{2}}{\mu_{f}} L_{f}^{2} L_{z}^{2} \hat{c}_{\beta}^{2} N^{3} \sigma_{g_{y}}^{2}, \end{split}$$

for any  $t \ge 1$ , where  $\hat{\Delta}_0 = \sum_{n=0}^{N-1} \left( \ell\left(\mathbf{x}_0^n\right) - \ell^* \right) + \sum_{n=0}^{N-1} \|\mathbf{y}_0^n - \mathbf{y}^*\left(\mathbf{x}_0^n\right)\|^2 + \|\mathbf{z}_0 - \mathbf{z}_0^*\|^2.$ 

*Proof.* Selecting a constant step-size  $\alpha_t = \alpha$  for all  $t \in \{0, 1, \dots, T-1\}$  and from Lemma D.4, we have 

$$\begin{split} & \begin{array}{l} & \begin{array}{l} \mathbf{1913} \\ & \mathbf{1914} \\ & \mathbf{1914} \\ & \begin{array}{l} \mathbf{1914} \\ & \mathbf{1914} \\ \end{array} \\ & \begin{array}{l} \mathbf{1915} \\ & \mathbf{1915} \\ & \begin{array}{l} \mathbf{1915} \\ & \mathbf{1916} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1917} \\ \end{array} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1917} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1917} \\ \end{array} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1917} \\ \end{array} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1917} \\ \end{array} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1917} \\ \end{array} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1917} \\ \end{array} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1917} \\ \end{array} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1917} \\ \end{array} \\ & \begin{array}{l} \mathbf{1916} \\ & \mathbf{1918} \end{array} \\ & \begin{array}{l} \mathbf{1917} \\ \mathbf{1918} \end{array} \\ \\ & \begin{array}{l} \mathbf{1918} \\ \mathbf{1918} \end{array} \\ \\ & \begin{array}{l} \mathbf{1917} \\ \mathbf{1918} \end{array} \\$$

Applying the above inequality recursively yields

$$\mathbb{E}\left[\hat{W}_{t}\right] \leq (1-\mu_{f}\alpha)^{t} \mathbb{E}\left[\hat{W}_{0}\right] + \sum_{k=0}^{t-1} (1-\mu_{f}\alpha)^{k} \left(+16L_{f}^{2}L_{z}^{2}\hat{c}_{\beta}^{2}\alpha^{3}N^{3}\sigma_{g_{y}}^{2} + 2\sigma_{f_{x}}^{2}\alpha N + 4\hat{c}_{\beta}^{2}\sigma_{g_{y}}^{2}N\alpha^{2} + 8\sigma_{f_{y}}^{2}\hat{c}_{\gamma}^{2}\alpha^{2} + 16\sigma_{g_{yy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\hat{c}_{\gamma}^{2}\alpha^{2} + \frac{8}{\mu_{g}\hat{c}_{\gamma}}L_{z}^{2}\hat{c}_{\beta}^{2}N^{2}\sigma_{g_{y}}^{2}\alpha + 4\sigma_{g_{xy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\alpha N\right)$$

$$\stackrel{(a)}{\leq} (1-\mu_{f}\alpha)^{t} \mathbb{E}\left[\hat{W}_{0}\right] + \frac{1}{\mu_{f}}\left(16L_{f}^{2}L_{z}^{2}\hat{c}_{\beta}^{2}\alpha^{2}N^{3}\sigma_{g_{y}}^{2} + 4\sigma_{g_{xy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}N + 2\sigma_{f_{x}}^{2}N + 4\hat{c}_{\beta}^{2}\sigma_{g_{y}}^{2}N\alpha^{2}\right)$$

$$+16\sigma_{g_{yy}}^{2}\frac{B_{f_{y}}^{2}}{\mu_{g}^{2}}\hat{c}_{\gamma}^{2}\alpha+8\sigma_{f_{y}}^{2}\hat{c}_{\gamma}^{2}\alpha+\frac{8}{\mu_{g}\hat{c}_{\gamma}}L_{z}^{2}\hat{c}_{\beta}^{2}N^{2}\sigma_{g_{y}}^{2}\right),$$

where (a) follows from the summation of a geometric progression. 

Utilizing the definition of the potential function  $\hat{W}_t$  and Jenson's inequality finishes the proof of the theorem.