

Stochastic Regret Guarantees for Online Zeroth- and First-Order Bilevel Optimization

Anonymous Authors¹

Abstract

Online bilevel optimization (OBO) is a powerful framework for machine learning problems where both outer and inner objectives evolve over time, requiring dynamic updates. Current OBO approaches rely on deterministic *window-smoothed* regret minimization, which may not accurately reflect system performance when functions change rapidly. In this work, we introduce a novel search direction and show that both first- and zeroth-order (ZO) stochastic OBO algorithms leveraging this direction achieve sublinear stochastic bilevel regret without window smoothing. Beyond these guarantees, our framework enhances efficiency by: (i) reducing oracle dependence in hypergradient estimation, (ii) updating inner and outer variables alongside the linear system solution, and (iii) employing ZO-based estimation of Hessians, Jacobians, and gradients. Experiments on online parametric loss tuning and black-box adversarial attacks validate our approach.

1. Introduction

Bilevel optimization (BO) minimizes an outer objective dependent on an inner problem’s solution. Originating in game theory (Stackelberg, 1952) and formalized in mathematical optimization (Bracken & McGill, 1973), BO finds applications in operations research, engineering, economics (Dempe, 2002), and image processing (Crockett et al., 2022). Recently, BO has gained traction in machine learning, including hyperparameter optimization (Franceschi et al., 2018), meta-learning (Finn et al., 2017), reinforcement learning (Stadie et al., 2020), and neural architecture search (Liu et al., 2018a).

In the *offline setting*, BO solves the following problem:

$$\begin{aligned} \mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^{d_1}} f(\mathbf{x}, \mathbf{y}^*(\mathbf{x})) \\ \text{subj. to } \mathbf{y}^*(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^{d_2}} g(\mathbf{x}, \mathbf{y}), \end{aligned} \quad (\text{BO})$$

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

where f and g are the outer and inner objectives, and \mathbf{x} and \mathbf{y} are their respective optimization variables.

OBO (Tarzanagh et al., 2024) addresses dynamic scenarios where objectives evolve over time, requiring the agent to update the outer decision in response to the optimal inner decision. Similar to online single-level optimization (OSO) (Zinkevich, 2003), OBO involves iterative decision-making without prior knowledge of outcomes (Tarzanagh et al., 2024; Lin et al., 2024; Bohne et al., 2024). Let T be the total number of rounds. Define $\mathbf{x}_t \in \mathcal{X} \subset \mathbb{R}^{d_1}$ as the decision variable and $f_t : \mathcal{X} \times \mathbb{R}^{d_2} \rightarrow \mathbb{R}$ as the outer function. Similarly, define $\mathbf{y}_t \in \mathbb{R}^{d_2}$ and $g_t : \mathcal{X} \times \mathbb{R}^{d_2} \rightarrow \mathbb{R}$ for the inner problem, where $\mathbf{y}_t^*(\mathbf{x}) = \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^{d_2}} g_t(\mathbf{x}, \mathbf{y})$. OBO can be seen as a *single-player* problem, where the player selects \mathbf{x}_t without knowing $\mathbf{y}_t^*(\mathbf{x})$, using \mathbf{y}_t as an estimate based on g_t . Alternatively, it can be framed as a *two-player* game (Stackelberg, 1952), where the leader (\mathbf{x}_t) competes with the follower (\mathbf{y}_t), who selects $\mathbf{y}_t^*(\mathbf{x})$ based on limited knowledge of g_t ; see Section 2. This framework includes online and adversarial variants of (BO), such as online actor-critic algorithms (Zhou et al., 2020), online meta-learning (Finn et al., 2019), and online hyperparameter optimization (Lin et al., 2024). The inner and outer functions may be time-varying, adversarial, unavailable *a priori*, and require *nonstationary* optimization.

1.1. Our Contributions

This paper addresses stochastic OBO, introducing novel first- and zeroth-order methods to minimize stochastic bilevel regret. Key contributions are summarized below.

- **Stochastic regret minimization without window-smoothing.** Existing OBO methods (Tarzanagh et al., 2024; Lin et al., 2024; Huang et al., 2023; Bohne et al., 2024) rely on deterministic *window-smoothed* regret minimization, which may not accurately reflect system performance when functions change rapidly. We address these limitations by introducing a novel search direction (Section 3) and proving that both first-order and ZO methods achieve sublinear *stochastic bilevel regret without window-smoothing* ($w = 1$); see Theorems 3.6 and 4.2 and Table 1.

• **OBO with function value oracle feedback.** In large-scale

¹First-order refers to the setting where only partial gradients of the leader objective f_t are accessible, while second-order information is still required for the follower objective g_t ; refer to Section 3.

| OBO Method | Window Size in Regret (w) | System Iterations | Stochastic Regret | Const. Regret Min. | Only Func. Feedback | Local Regret Bound |
|------------|-------------------------------|---------------------------------------|-------------------|--------------------|---------------------|--|
| OAGD | $o(T)$ | N.A. (Exact) | ✗ | ✗ | ✗ | $\frac{T}{w} + H_{1,T} + H_{2,T}$ |
| SOBOW | $o(T)$ | $\mathcal{O}(\kappa_g \log \kappa_g)$ | ✗ | ✗ | ✗ | $\frac{T}{w} + V_T + H_{2,T}$ |
| SOBBO | $o(T)$ | $\mathcal{O}(\kappa_g \log \kappa_g)$ | ✓ | ✓ | ✗ | $\frac{T}{w} \sigma^2 + V_T + H_{2,T}$ |
| SOGD | 1 | 1 | ✓ | ✓ | ✗ | $T^{\frac{1}{3}}(\sigma^2 + \Delta_T) + T^{\frac{2}{3}}\Psi_T$ |
| ZO-SOGD | 1 | 1 | ✓ | ✓ | ✓ | $(d_1 + d_2)^{\frac{3}{4}}T^{\frac{1}{4}}(\hat{\sigma}^2 + \hat{\Delta}_T) + (d_1 + d_2)^{\frac{3}{2}}T^{\frac{3}{4}}\hat{\Psi}_T$ |

Table 1. Comparison of OBO algorithms based on regret window size (w), system solver iterations, stochastic regret, constrained regret minimization, function feedback settings, and local regret bounds. Here, κ_g denotes the condition number of g_t , while V_T , $H_{p,T}$, Δ_T , Ψ_T , $\hat{\Delta}_T$, and $\hat{\Psi}_T$ are defined in (10), (13), and (25), respectively. The compared algorithms include OAGD (Tarzanagh et al., 2024), SOBOW (Lin et al., 2024), and SOBBO (Bohne et al., 2024).

and black-box settings (Chen et al., 2017; Nesterov, 2005), first- and second-order information is often unavailable or costly. Constructing accurate (hyper)-gradient estimators using only function value oracles is particularly challenging due to BO’s nested structure. Existing methods rely on gradient, Hessian, and Jacobian oracles, limiting scalability (Franceschi et al., 2017; Ghadimi & Wang, 2018). We propose Algorithm 2, which estimates Hessians, Jacobians, and gradients using function value oracles, achieving sublinear local regret (Theorem 4.2).

• **OBO with one subproblem solver iteration.** A major challenge in BO is solving implicit systems to approximate the hypergradient (Ji et al., 2021; Chen et al., 2021). While efficient offline BO methods exist (Ji et al., 2021; Dagréou et al., 2022), extending them to OBO is difficult due to time-varying objectives. SOBOW (Lin et al., 2024) partially addresses this using a conjugate gradient (CG) algorithm with increasing iterations (Table 1). We improve upon SOBOW by introducing Algorithms 1 and 2, which require only a *single* subproblem solver iteration.

2. Preliminaries

Notation. \mathbb{R}^d denotes the d -dimensional real space, with \mathbb{R}_+^d and \mathbb{R}_{++}^d as its positive and negative orthants. Vectors are bold lower-case letters (e.g., \mathbf{x}, \mathbf{y}), with $\langle \mathbf{x}, \mathbf{y} \rangle$ for inner product and $\|\cdot\|$ for Euclidean norm. A gradient is $\nabla_{\mathbf{x}}$, with $\nabla_{\mathbf{x}\mathbf{y}}^2 = \nabla_{\mathbf{x}}\nabla_{\mathbf{y}}$. A function is L -smooth if its gradient is L -Lipschitz. The Euclidean projection onto a convex set \mathcal{X} is $\Pi_{\mathcal{X}}(\mathbf{z}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} (1/2)\|\mathbf{x} - \mathbf{z}\|^2$. The set $\{1, \dots, T\}$ is denoted by $[T]$, and $\mathbb{E}[\cdot]$ represents expectation. Lastly, $\mathcal{O}(\cdot)$ hides problem-independent constants.

Stochastic OBO Setting. Let T be the total rounds (Tarzanagh et al., 2024). Define $\mathbf{x}_t \in \mathcal{X} \subset \mathbb{R}^{d_1}$ as the decision variable and $f_t : \mathcal{X} \times \mathbb{R}^{d_2}$ as the outer objective. The inner decision variable and objective are $\mathbf{y}_t \in \mathbb{R}^{d_2}$ and $g_t : \mathcal{X} \times \mathbb{R}^{d_2}$, where the optimal inner decision is:

$$\mathbf{y}_t^*(\mathbf{x}) \in \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^{d_2}} \left\{ g_t(\mathbf{x}, \mathbf{y}) := \mathbb{E}_{\zeta_t \sim \mathcal{D}_g} [g_t(\mathbf{x}, \mathbf{y}; \zeta_t)] \right\}. \quad (1)$$

Further, we have

$$f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) := \mathbb{E}_{\xi_t \sim \mathcal{D}_f} [f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}); \xi_t)].$$

Here, $(\mathcal{D}_f, \mathcal{D}_g)$ are data distributions. Note that our setting is stochastic, and only noisy evaluations of the function, gradient, and Hessian are accessible.

Unlike OSO, where true losses are revealed immediately, in OBO, the outer function $f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))$ is unavailable for updating \mathbf{x}_t . Moreover, $f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))$ is typically non-convex in \mathbf{x} , making standard regret definitions from online convex optimization (Hazan, 2016b) inapplicable.

Given a sequence $\{\alpha_t \in \mathbb{R}_{++}\}_{t=1}^T$, we define the following notion of *bilevel local regret*:

$$\text{BL-Reg}_T := \sum_{t=1}^T \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} (\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \right\|^2 \right], \quad (2a)$$

with

$$\begin{aligned} & \mathcal{P}_{\mathcal{X}, \alpha_t} (\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \\ &= \frac{1}{\alpha_t} \left(\mathbf{x}_t - \Pi_{\mathcal{X}} \left[\mathbf{x}_t - \alpha_t \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \right] \right). \end{aligned} \quad (2b)$$

The local regret (2) compares the leader’s decision \mathbf{x}_t to the stationary points \mathbf{x}_t^* satisfying $\mathcal{P}_{\mathcal{X}, \alpha_t} (\mathbf{x}_t^*; \nabla f_t(\mathbf{x}_t^*, \mathbf{y}_t^*(\mathbf{x}_t^*))) = 0$. This can also be viewed as dynamic local regret, as the baseline corresponds to a stationary point of the leader’s objective f_t .

Previous work on (nonconvex) OBO examined unconstrained local regret using window-smoothed objectives: $F_{t,w}(\mathbf{x}, \mathbf{y}) = (1/w) \sum_{i=0}^{w-1} f_{t-i}(\mathbf{x}, \mathbf{y})$. For $w = 1$ and $\mathcal{X} = \mathbb{R}^{d_1}$, this reduces to (2). Tarzanagh et al. (2024); Lin et al. (2024) showed that $w = o(T)$ ensures sublinear regret under slow variations in $\{F_{t,w}\}_{t=1}^T$, while rapid changes can lead to deviations. However, smoothing may misrepresent regret (Figure 1). This paper introduces a new projection-based local regret notion (2) without smoothing, and establishes sublinear regret for constrained OBO.

Online Gradient Descent (OGD). One of the most widely used algorithms for online (single-level) optimization is

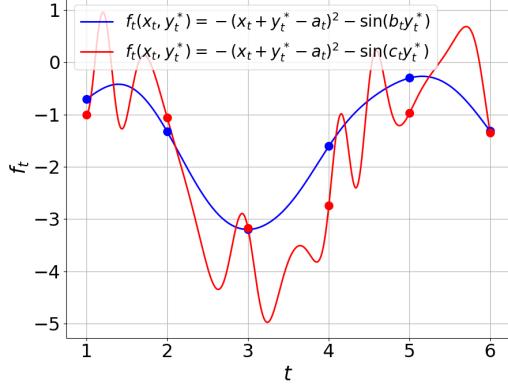


Figure 1. Smoothly and rapidly changing f_t in OBO with $g_t(x_t, y_t) = (y_t - \cos(x_t))^2$, $a_t = 1 + 0.5 \sin(t)$, $b_t = 1 + \sin(0.5t)$, and $c_t = 10b_t$.

OGD (Zinkevich, 2003). The procedure for OGD is as follows: For each $t \in [T]$, the algorithm selects $\mathbf{x}_t \in \mathcal{X}$, observes the function $f_t : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}$, and updates according to

$$\mathbf{x}_{t+1} = \Pi_{\mathcal{X}}(\mathbf{x}_t - \alpha_t \nabla f_t(\mathbf{x}_t)), \quad \alpha_t > 0. \quad (\text{OGD})$$

In the following, we adapt OGD to OBO and introduce a novel framework that requires limited feedback and can utilize ZO updates within a single-loop structure.

3. Stochastic OBO with Access to First and Second Order Oracles

To adapt OGD to OBO, Tarzanagh et al. (2024); Lin et al. (2024); Bohne et al. (2024) developed a variant alternating between inner and outer OGD, achieving sublinear bilevel regret bounds. We introduce a new search direction that enables sublinear bilevel regret without window smoothing.

To compute the hypergradient $\nabla f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))$ where $\mathbf{y}_t^*(\mathbf{x})$ is defined in (1), since $\nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) = 0$, using the implicit function theorem, yields

$$\begin{aligned} \nabla f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) &= \nabla_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) \\ &\quad + \nabla_{\mathbf{y}_t^*(\mathbf{x})} \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})), \end{aligned} \quad (3)$$

where $\nabla_{\mathbf{y}_t^*(\mathbf{x})} \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) + \nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) = 0$.

As the exact $\mathbf{y}_t^*(\mathbf{x})$ is not available, we estimate the hypergradient of f_t at (\mathbf{x}, \mathbf{y}) by

$$\tilde{\nabla} f_t(\mathbf{x}, \mathbf{y}) := \nabla_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{y}) + \nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y}) \mathbf{v}_t^*(\mathbf{x}), \quad (4a)$$

where

$$\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y}) \mathbf{v}_t^*(\mathbf{x}) + \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y}) = 0. \quad (4b)$$

An accurate solution of (4b) is crucial for tight regret bounds. Tarzanagh et al. (2024) assumes an exact solution, which is restrictive in large-scale settings. To address this, Lin et al. (2024) proposed an efficient OBO algorithm with window averaging, using CG methods to solve (4b), which

Algorithm 1 SOGD

Require: $(\mathbf{x}_1, \mathbf{y}_1, \mathbf{v}_1) \in \mathcal{X} \times \mathbb{R}^{d_2} \times \mathbb{R}^{d_2}; T \in \mathbb{N}; p \in \mathbb{R}_{++};$ stepsizes $\{(\alpha_t, \beta_t, \delta_t) \in \mathbb{R}_{++}^3\}_{t=1}^T;$ parameters $\{(\gamma_t, \lambda_t, \eta_t)\}_{t=1}^T \in (0, 1); \mathbf{z}_t := (\mathbf{x}_t, \mathbf{y}_t).$

For $t = 1$ **to** T **do:**

S1. Draw samples \mathcal{B}_t and $\bar{\mathcal{B}}_t$ with batch sizes b and \bar{b} . Get search directions $\mathbf{d}_t^{\mathbf{y}}, \mathbf{d}_t^{\mathbf{v}}$, and $\mathbf{d}_t^{\mathbf{x}}$:

$$\mathbf{d}_t^{\mathbf{yy}}(\mathbf{z}_t; \bar{\mathcal{B}}_t) = \nabla_{\mathbf{y}} g_t(\mathbf{z}_t; \bar{\mathcal{B}}_t), \quad (7a)$$

$$\mathbf{d}_t^{\mathbf{yy}}(\mathbf{z}_t; \mathcal{B}_t) = \nabla_{\mathbf{y}} f_t(\mathbf{z}_t; \mathcal{B}_t) + \nabla_{\mathbf{y}}^2 g_t(\mathbf{z}_t; \bar{\mathcal{B}}_t) \mathbf{v}_t, \quad (7b)$$

$$\mathbf{d}_t^{\mathbf{yy}}(\mathbf{z}_t; \mathcal{B}_t) = \nabla_{\mathbf{x}} f_t(\mathbf{z}_t; \mathcal{B}_t) + \nabla_{\mathbf{xy}}^2 g_t(\mathbf{z}_t; \bar{\mathcal{B}}_t) \mathbf{v}_t, \quad (7c)$$

$$\mathbf{d}_t^{\mathbf{xx}}(\mathbf{z}_t; \mathcal{B}_t) = \nabla_{\mathbf{x}} f_t(\mathbf{z}_t; \mathcal{B}_t) + \nabla_{\mathbf{xy}}^2 g_t(\mathbf{z}_t; \bar{\mathcal{B}}_t) \mathbf{v}_t, \quad (7c)$$

$$\mathbf{d}_t^{\mathbf{x}} = \mathbf{d}_t^{\mathbf{xx}}(\mathbf{z}_t; \mathcal{B}_t) + (1 - \eta_t)(\mathbf{d}_{t-1}^{\mathbf{x}} - \mathbf{d}_t^{\mathbf{xx}}(\mathbf{z}_{t-1}; \mathcal{B}_t)).$$

S2. Update inner, system, and outer solutions:

$$\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{y}_t - \beta_t \mathbf{d}_t^{\mathbf{y}}, \quad \mathbf{v}_{t+1} = \Pi_{\mathcal{Z}_p}[\mathbf{v}_t - \delta_t \mathbf{d}_t^{\mathbf{y}}], \\ \mathbf{x}_{t+1} &= \Pi_{\mathcal{X}}[\mathbf{x}_t - \alpha_t \mathbf{d}_t^{\mathbf{x}}]. \end{aligned}$$

is equivalent to:

$$\min_{\mathbf{v}_t \in \mathbb{R}^{d_2}} (1/2) \|\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y}) \mathbf{v}_t + \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\|^2. \quad (5)$$

New Search Direction for OBO. Next, we introduce a novel search direction that enables both first- and ZO stochastic OBO algorithms to achieve sublinear bilevel regret without smoothing. We first state the following lemma:

Lemma 3.1. Let $w = t$ and $W = 1/\eta$ in the window-smoothed gradient $\hat{\nabla} F_{t,\nu}(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) = (1/W) \sum_{i=0}^{w-1} \nu^i \hat{\nabla} f_{t-i}(\mathbf{x}_{t-i}, \mathbf{y}_{t-i}; \mathcal{B}_{t-i})$, where $\mathcal{B}_t := \{\xi_1, \dots, \xi_b\}$ is drawn i.i.d. from \mathcal{D}_f . Then,

$$\hat{\nabla} F_{t,\nu}(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) = \sum_{j=1}^t \eta(1 - \eta)^{t-j} \hat{\nabla} f_j(\mathbf{x}_j, \mathbf{y}_j; \mathcal{B}_j).$$

Furthermore, we have $\hat{\nabla} F_{t,\nu}(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) = \hat{\mathbf{d}}_t^{\mathbf{x}}$ with $\hat{\mathbf{d}}_t^{\mathbf{x}} = \eta \hat{\nabla} f_t(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) + (1 - \eta) \hat{\mathbf{d}}_{t-1}^{\mathbf{x}}$, and $\hat{\mathbf{d}}_1 = (1/W) \hat{\nabla} f_1(\mathbf{x}_1, \mathbf{y}_1; \mathcal{B}_1)$ for all $t \geq 2$.

As shown in Lemma 3.1, for a specific choice of w and W , the time-smoothed gradient forms a recursive momentum-type search direction. However, achieving sublinear regret in stochastic OBO requires large-window smoothing ($w = o(T)$). To address this, we propose the following search direction:

$$\mathbf{d}_t^{\mathbf{x}} = \eta \nabla f_t(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) + (1 - \eta) \mathbf{d}_{t-1}^{\mathbf{x}} \quad (6a)$$

$$+ (1 - \eta)(\nabla f_t(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) - \nabla f_t(\mathbf{x}_{t-1}, \mathbf{y}_{t-1}; \mathcal{B}_t)). \quad (6b)$$

This direction is used for updating \mathbf{x} , with similar updates for \mathbf{y} and \mathbf{v} , as discussed below and detailed in Algorithm 1.

The quadratic optimization formulation of (4b) in (5) leads to single-loop frameworks such as Dagréou et al. (2022). Inspired by this, we present Simultaneous Online Gradient Descent (SOGD) for constrained OBO, outlined in Algorithm 1. SOGD evolves the follower's decision (inner) variable, the linear system solution, and the leader's decision (outer) variable simultaneously at each step for given batches $\mathcal{B} := \{\xi_1, \dots, \xi_b\}$ and $\bar{\mathcal{B}} := \{\zeta_1, \dots, \zeta_{\bar{b}}\}$, which are drawn i.i.d. from unknown distributions \mathcal{D}_f and \mathcal{D}_g with batch sizes b and \bar{b} . Computing directions in S1. of Algorithm 1 does not require $\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)$ and $\nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)$, only their product with a vector, at the same cost as computing a gradient. Technically, it utilizes an auxiliary variable \mathbf{v}_t and concurrently updates \mathbf{y}_t , \mathbf{v}_t , and \mathbf{x}_t at each local iteration t . Moreover, S2. of Algorithm 1 introduces an auxiliary projection $\Pi_{\mathcal{Z}_p}$ on the ball \mathcal{Z}_p defined as follows:

$$\Pi_{\mathcal{Z}_p}(\mathbf{v}) := \min \left\{ 1, \frac{p}{\|\mathbf{v}\|} \right\} \mathbf{v}, \quad (8)$$

where $\mathcal{Z}_p := \{\mathbf{v} \in \mathbb{R}^{d_2} \mid \|\mathbf{v}\| \leq p\}$.

Unlike OAGD (Tarzanagh et al., 2024), which updates \mathbf{x} and \mathbf{y} in separate loops, SOGD updates both simultaneously. Compared to SOBOW (Lin et al., 2024), which uses multiple CG updates, our method employs a single OGD to update the inner solution, linear system, and outer variable.

Assumption 3.2. $g_t(\mathbf{x}, \mathbf{y})$ is twice continuously differentiable and μ_g -strongly convex in \mathbf{y} for all $\mathbf{x} \in \mathcal{X}, t \in [T]$.

Assumption 3.3. Let $\mathbf{z} = [\mathbf{x}; \mathbf{y}]$ and $\mathbf{z}' = [\mathbf{x}'; \mathbf{y}']$, where $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ and $\mathbf{y}, \mathbf{y}' \in \mathbb{R}^{d_2}$. For any \mathbf{z}, \mathbf{z}' , and $t \in [T]$:

- B1. $\exists \ell_{f,0} \in \mathbb{R}_+$ s.t. $\|f_t(\mathbf{z}; \xi) - f_t(\mathbf{z}'; \xi)\| \leq \ell_{f,0} \|\mathbf{z} - \mathbf{z}'\|$;
- B2. $\exists \ell_{f,1} \in \mathbb{R}_+$ s.t. $\|\nabla f_t(\mathbf{z}; \xi) - \nabla f_t(\mathbf{z}'; \xi)\| \leq \ell_{f,1} \|\mathbf{z} - \mathbf{z}'\|$;
- B3. $\exists \ell_{g,1} \in \mathbb{R}_+$ s.t. $\|\nabla g_t(\mathbf{z}; \zeta) - \nabla g_t(\mathbf{z}'; \zeta)\| \leq \ell_{g,1} \|\mathbf{z} - \mathbf{z}'\|$;
- B4. $\exists \ell_{g,2} \in \mathbb{R}_+$ s.t. $\|\nabla^2 g_t(\mathbf{z}; \zeta) - \nabla^2 g_t(\mathbf{z}'; \zeta)\| \leq \ell_{g,2} \|\mathbf{z} - \mathbf{z}'\|$.

Assumption 3.4. For any $t \in [T]$, $|f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))| \leq M$ for some finite constant $M \in \mathbb{R}_{++}$ and any $\mathbf{x} \in \mathcal{X}$.

Assumption 3.5. There exist constants $\sigma_{g_y}, \sigma_{g_{yy}}, \sigma_{g_{xy}}, \sigma_{f_y}, \sigma_{f_x}$ such that, for all $\mathbf{z} = [\mathbf{x}, \mathbf{y}]$:

- C1. $\mathbb{E} \|\nabla_{\mathbf{y}} g_t(\mathbf{z}; \zeta) - \nabla_{\mathbf{y}} g_t(\mathbf{z})\|^2 \leq \sigma_{g_y}^2$,
- C2. $\mathbb{E} \|\nabla_{\mathbf{y}}^2 g_t(\mathbf{z}; \zeta) - \nabla_{\mathbf{y}}^2 g_t(\mathbf{z})\|^2 \leq \sigma_{g_{yy}}^2$,
- C3. $\mathbb{E} \|\nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{z}; \zeta) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{z})\|^2 \leq \sigma_{g_{xy}}^2$,
- C4. $\mathbb{E} \|\nabla_{\mathbf{y}} f_t(\mathbf{z}; \xi) - \nabla_{\mathbf{y}} f_t(\mathbf{z})\|^2 \leq \sigma_{f_y}^2$,
- C5. $\mathbb{E} \|\nabla_{\mathbf{x}} f_t(\mathbf{z}; \xi) - \nabla_{\mathbf{x}} f_t(\mathbf{z})\|^2 \leq \sigma_{f_x}^2$.

Throughout this paper, we define

$$\sigma^2 := \sigma_{g_y}^2 + \sigma_{g_{yy}}^2 + \sigma_{g_{xy}}^2 + \sigma_{f_y}^2 + \sigma_{f_x}^2. \quad (9)$$

Assumptions 3.2 and 3.3 are widely used in both BO (Chen et al., 2021; Ji et al., 2021) and OBO (Tarzanagh et al., 2024), and many bilevel machine learning problems satisfy it (Franceschi et al., 2018). Further, Assumption 3.4 is widely used in the study of non-convex online optimization (Hazan et al., 2017; Lin et al., 2024). Assumption 3.5

assumes that we have access to an unbiased stochastic gradient, Hessian and Jacobian with bounded variance, which is standard in the literature (Chen et al., 2021).

Achieving sublinear dynamic regret is generally impossible due to arbitrary fluctuations in time-varying functions (Besbes et al., 2015). Existing analyses (Tarzanagh et al., 2024; Lin et al., 2024) bound regret by imposing regularity constraints on the comparator sequence. To achieve sublinear regret, we introduce the following regularities:

- **Path-length (of order p) and function variation:** Tarzanagh et al. (2024) defines the following metrics for bilevel sequences:

$$\begin{aligned} H_{p,T} &:= \sum_{t=2}^T \sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{y}_{t-1}^*(\mathbf{x}) - \mathbf{y}_t^*(\mathbf{x})\|^p, \\ V_T &:= \sum_{t=2}^T \sup_{\mathbf{x} \in \mathcal{X}} |f_{t-1}(\mathbf{x}, \mathbf{y}_{t-1}^*(\mathbf{x})) - f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))|. \end{aligned} \quad (10)$$

Path-length $H_{p,T}$ measures changes in the follower's costs, while V_T captures the smoothness of the leader's objective. We use path-length for the follower and function variation for the leader, as the follower's objective is strongly convex (see Assumption 3.2), while the leader's is nonconvex.

- **Inner and Outer Gradient Variations:** Another regularity is the sequential difference between the individual gradients of the upper-level loss function:

$$\begin{aligned} D_{\mathbf{x},T} &:= \sum_{t=2}^T \sup_{\mathbf{x}, \mathbf{y}} \|\nabla_{\mathbf{x}} f_{t-1}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{y})\|^2, \\ D_{\mathbf{y},T} &:= \sum_{t=2}^T \sup_{\mathbf{x}, \mathbf{y}} \|\nabla_{\mathbf{y}} f_{t-1}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\|^2. \end{aligned} \quad (11)$$

As in Huang et al.; Hallak et al. (2021), $D_{\mathbf{x},T}$ and $D_{\mathbf{y},T}$ measure the gradient drift of f_t relative to f_{t-1} for \mathbf{x} and \mathbf{y} , respectively. We further define deviations in the gradient, Hessian, and Jacobian of the lower-level objective as:

$$\begin{aligned} G_{\mathbf{y},T} &:= \sum_{t=2}^T \|\nabla_{\mathbf{y}} g_{t-1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2, \\ G_{\mathbf{y}\mathbf{y},T} &:= \sum_{t=2}^T \|\nabla_{\mathbf{y}}^2 g_{t-1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2, \\ G_{\mathbf{x}\mathbf{y},T} &:= \sum_{t=2}^T \|\nabla_{\mathbf{x}\mathbf{y}}^2 g_{t-1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2. \end{aligned} \quad (12)$$

We introduce the following notations for simplicity:

$$\Delta_T := E_1 + V_T, \quad \Psi_T := H_{2,T} + G_T + D_T, \quad (13)$$

where $(V_T, H_{p,T})$ are defined in (10), and

$$\begin{aligned} E_1 &:= \|\mathbf{y}_1 - \mathbf{y}_1^*(\mathbf{x}_1)\|^2 + \|\mathbf{v}_1 - \mathbf{v}_1^*(\mathbf{x}_1)\|^2, \\ G_T &:= G_{\mathbf{y},T} + G_{\mathbf{y}\mathbf{y},T} + G_{\mathbf{x}\mathbf{y},T}, \\ D_T &:= D_{\mathbf{x},T} + D_{\mathbf{y},T}. \end{aligned} \quad (14)$$

By accounting for both D_T and G_T , we can represent the variations in the environments of OBO.

Theorem 3.6. Let $\{(f_t, g_t)\}_{t=1}^T$ be the sequence of functions presented to Algorithm 1, satisfying Assumptions 3.2–3.5. For all $t \in [T]$, let

$$\begin{aligned}\alpha_t &= \frac{1}{(c+t)^{1/3}}, & \beta_t &= c_\beta \alpha_t, & \delta_t &= c_\delta \alpha_t, & b = \bar{b} &= 1, \\ \gamma_{t+1} &= c_\gamma \alpha_t^2, & \eta_{t+1} &= c_\eta \alpha_t^2, & \lambda_{t+1} &= c_\lambda \alpha_t^2.\end{aligned}\quad (15)$$

Here, $c, c_\beta, c_\delta, c_\gamma, c_\eta$, and c_λ are specified in (104). Algorithm 1 guarantees:

$$\text{BL-Reg}_T \leq \mathcal{O}\left(T^{1/3}(\sigma^2 + \Delta_T) + T^{2/3}\Psi_T\right), \quad (16)$$

where σ and (Δ_T, Ψ_T) are defined in (9) and (13).

Theorem 3.6 bounds the regret of Algorithm 1 without window-smoothing, based on the regularities in (14). We note that the average dynamic regret $\text{BL-Reg}_T/T \leq \mathcal{O}(T^{-2/3}(\sigma^2 + \Delta_T) + T^{-1/3}\Psi_T)$ remains sublinear under suitable conditions on Δ_T and Ψ_T .

Remark 3.7 (Stochastic Regret Guarantee for OBO and OSO with $w = 1$). The additional terms in (6b) improve the average regret dependence on variance, achieving a $T^{-2/3}\sigma^2$ bound, better than the $T^{-1/2}\sigma^2$ bound for stochastic OBO (Bohne et al., 2024). This also provides the first regret bound without window-smoothing, unlike (Bohne et al., 2024; Tarzanagh et al., 2024; Lin et al., 2024; Huang et al., 2023). For OSO, our approach improves the $T^{-1/2}\sigma^2$ dependence from (Hallak et al., 2021).

4. OBO with Zeroth Order Oracles

Black-box optimization arises in machine learning when explicit gradients are unavailable (Chen et al., 2017). We study ZO-type OBO algorithms with limited access to the leader’s and follower’s objective values. Let $\mathbf{s} \in \mathbb{R}^{d_1}$ and $\mathbf{r} \in \mathbb{R}^{d_2}$ be vectors uniformly generated from the unit balls B_1 and B_2 , respectively. Given positive smoothing parameters $\rho = (\rho_s, \rho_r)$, we use the Gaussian smoothing function (Nesterov & Spokoiny, 2017) to define the OBO objectives:

$$f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) = \mathbb{E}_{(\mathbf{s}, \mathbf{r})} [f_t(\mathbf{x} + \rho_s \mathbf{s}, \hat{\mathbf{y}}_t^*(\mathbf{x}) + \rho_r \mathbf{r}; \xi)], \quad (17)$$

where

$$\begin{aligned}\hat{\mathbf{y}}_t^*(\mathbf{x}) &\in \operatorname{argmin}_{\mathbf{y} \in \mathbb{R}^{d_2}} \{g_{t,\rho}(\mathbf{x}, \mathbf{y}) \\ &:= \mathbb{E}_{(\mathbf{s}, \mathbf{r})} [g_t(\mathbf{x} + \rho_s \mathbf{s}, \mathbf{y} + \rho_r \mathbf{r}; \xi)]\}.\end{aligned}\quad (18)$$

Using (17), we provide methodology to approximate each term in (7) using ZO oracles. Specifically, following Shamir (2017), we estimate the gradient of a function $h : \mathbb{R}^d \rightarrow \mathbb{R}$, querying at $\mathbf{x} - \lambda \mathbf{s}$ and $\mathbf{x} + \lambda \mathbf{s}$, yielding an estimator $(d/2\lambda)(h(\mathbf{x} + \lambda \mathbf{s}) - h(\mathbf{x} - \lambda \mathbf{s})) \mathbf{s}$. Using this strategy, the finite-difference estimation of $\nabla g_{t,\rho}(\mathbf{x}, \mathbf{y})$, denoted as $\hat{\nabla} g_t(\mathbf{x}, \mathbf{y})$, is constructed for given smoothing

parameters $\rho = (\rho_s, \rho_r)$, and batches $\mathcal{B} := \{\xi_1, \dots, \xi_b\}$ and $\bar{\mathcal{B}} := \{\zeta_1, \dots, \zeta_{\bar{b}}\}$, drawn i.i.d. from \mathcal{D}_f and \mathcal{D}_g , as:

$$\begin{aligned}\hat{\nabla}_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}}) &:= \frac{d_2}{2\bar{b}\rho_r} \sum_{i=1}^{\bar{b}} (g_t(\mathbf{x}, \mathbf{y} + \rho_r \mathbf{r}_i; \zeta_i) \\ &\quad - g_t(\mathbf{x}, \mathbf{y} - \rho_r \mathbf{r}_i; \zeta_i)) \mathbf{r}_i,\end{aligned}\quad (19a)$$

$$\begin{aligned}\hat{\nabla}_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}}) &:= \frac{d_1}{2\bar{b}\rho_s} \sum_{i=1}^{\bar{b}} (g_t(\mathbf{x} + \rho_s \mathbf{s}_i, \mathbf{y}; \zeta_i) \\ &\quad - g_t(\mathbf{x} - \rho_s \mathbf{s}_i, \mathbf{y}; \zeta_i)) \mathbf{s}_i.\end{aligned}\quad (19b)$$

Similarly, we estimate $\nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}; \mathcal{B})$ and $\nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \mathbf{y}; \mathcal{B})$, respectively, by

$$\begin{aligned}\hat{\nabla}_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y}; \mathcal{B}) &:= \frac{d_2}{2b\rho_r} \sum_{i=1}^b (f_t(\mathbf{x}, \mathbf{y} + \rho_r \mathbf{r}_i; \xi_i) \\ &\quad - f_t(\mathbf{x}, \mathbf{y} - \rho_r \mathbf{r}_i; \xi_i)) \mathbf{r}_i,\end{aligned}\quad (20a)$$

$$\begin{aligned}\hat{\nabla}_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{y}; \mathcal{B}) &:= \frac{d_1}{2b\rho_s} \sum_{i=1}^b (f_t(\mathbf{x} + \rho_s \mathbf{s}_i, \mathbf{y}; \xi_i) \\ &\quad - f_t(\mathbf{x} - \rho_s \mathbf{s}_i, \mathbf{y}; \xi_i)) \mathbf{s}_i.\end{aligned}\quad (20b)$$

Further, given a smoothing parameter $\rho_v > 0$, we can approximate the Hessian-vector product $\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) \mathbf{v}$ and the Jacobian-vector product $\nabla_{\mathbf{xy}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) \mathbf{v}$ as the finite difference between two gradients, respectively, as

$$\begin{aligned}\hat{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}}) &:= \frac{1}{2\bar{b}\rho_v} \sum_{i=1}^{\bar{b}} (\hat{\nabla}_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y} + \rho_v \mathbf{v}; \zeta_i) \\ &\quad - \hat{\nabla}_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y} - \rho_v \mathbf{v}; \zeta_i)),\end{aligned}\quad (21a)$$

$$\begin{aligned}\hat{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}}) &:= \frac{1}{2\bar{b}\rho_v} \sum_{i=1}^{\bar{b}} (\hat{\nabla}_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y} + \rho_v \mathbf{v}; \zeta_i) \\ &\quad - \hat{\nabla}_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y} - \rho_v \mathbf{v}; \zeta_i)).\end{aligned}\quad (21b)$$

Using (19)–(21), the first-order terms in (7) are approximated as $\hat{\mathbf{d}}_t^y$, $\hat{\mathbf{d}}_t^v$, and $\hat{\mathbf{d}}_t^x$ in (22). The approximations in (21a) and (21b) introduce errors in the hypergradient, which must be controlled. (21) depends on the dimension of \mathbf{y} , as in ZO optimization (Nesterov & Spokoiny, 2017; Shamir, 2017). The projection $\Pi_{\mathcal{Z}_p}$ in (8) bounds \mathbf{v} , controlling variance in \mathbf{v} and \mathbf{x} updates for convergence.

Assumption 4.1. There exist constants $\hat{\sigma}_{gy}, \hat{\sigma}_{gx}, \hat{\sigma}_{fy}, \hat{\sigma}_{fx}$ such that, for all $\mathbf{z} = [\mathbf{x}, \mathbf{y}]$, the following holds:

- D1. $\mathbb{E}\|\hat{\nabla}_{\mathbf{y}} g_t(\mathbf{z}; \zeta) - \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{z})\|^2 \leq \hat{\sigma}_{gy}^2$,
- D2. $\mathbb{E}\|\hat{\nabla}_{\mathbf{x}} g_t(\mathbf{z}; \zeta) - \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{z})\|^2 \leq \hat{\sigma}_{gx}^2$,
- D3. $\mathbb{E}\|\hat{\nabla}_{\mathbf{y}} f_t(\mathbf{z}; \xi) - \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{z})\|^2 \leq \hat{\sigma}_{fy}^2$,
- D4. $\mathbb{E}\|\hat{\nabla}_{\mathbf{x}} f_t(\mathbf{z}; \xi) - \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{z})\|^2 \leq \hat{\sigma}_{fx}^2$.

Assumption 4.1 is analogous to the upper bound on the variance of stochastic partial gradients discussed in Luo et al. (2020); Wang et al. (2020). We simplify the notation by introducing the following shorthand.

$$\hat{\sigma}^2 := \hat{\sigma}_{gy}^2 + \hat{\sigma}_{gx}^2 + \hat{\sigma}_{fy}^2 + \hat{\sigma}_{fx}^2. \quad (23)$$

Next, we establish a regret bound for ZO-SOGD. Similar to

Algorithm 2 ZO-SOGD

Require: In addition to parameters in SOGD, choose

$$\rho_v, \rho_r, \rho_s, \in \mathbb{R}_{++}.$$

For $t = 1$ **to** T **do:**

- S1.** Draw samples \mathcal{B}_t and $\bar{\mathcal{B}}_t$ with batch sizes b and \bar{b} . Using (19)–(21), get ZO search directions $\hat{\mathbf{d}}_t^y, \hat{\mathbf{d}}_t^v, \hat{\mathbf{d}}_t^x$:

$$\mathbf{d}_t^y(\mathbf{z}_t; \bar{\mathcal{B}}_t) = \hat{\nabla}_y g_t(\mathbf{z}_t; \bar{\mathcal{B}}_t), \quad (22a)$$

$$\hat{\mathbf{d}}_t^y = \mathbf{d}_t^y(\mathbf{z}_t; \bar{\mathcal{B}}_t) + (1 - \gamma_t)(\hat{\mathbf{d}}_{t-1}^y - \mathbf{d}_t^y(\mathbf{z}_{t-1}; \bar{\mathcal{B}}_t)),$$

$$\mathbf{d}_t^{vv}(\mathbf{z}_t; \mathcal{B}_t) = \hat{\nabla}_y f_t(\mathbf{z}_t; \mathcal{B}_t) + \hat{\nabla}_y^2 g_t(\mathbf{z}_t; \bar{\mathcal{B}}_t), \quad (22b)$$

$$\hat{\mathbf{d}}_t^v = \mathbf{d}_t^{vv}(\mathbf{z}_t; \mathcal{B}_t) + (1 - \lambda_t)(\hat{\mathbf{d}}_{t-1}^v - \mathbf{d}_t^{vv}(\mathbf{z}_{t-1}; \mathcal{B}_t)),$$

$$\mathbf{d}_t^{xy}(\mathbf{z}_t; \mathcal{B}_t) = \hat{\nabla}_x f_t(\mathbf{z}_t; \mathcal{B}_t) + \hat{\nabla}_x^2 g_t(\mathbf{z}_t; \bar{\mathcal{B}}_t), \quad (22c)$$

$$\hat{\mathbf{d}}_t^x = \mathbf{d}_t^{xy}(\mathbf{z}_t; \mathcal{B}_t) + (1 - \eta_t)(\hat{\mathbf{d}}_{t-1}^x - \mathbf{d}_t^{xy}(\mathbf{z}_{t-1}; \mathcal{B}_t)),$$

- S2.** Update inner, system, and outer solutions:

$$\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{y}_t - \beta_t \hat{\mathbf{d}}_t^y, & \mathbf{v}_{t+1} &= \Pi_{\mathcal{Z}_p}[\mathbf{v}_t - \delta_t \hat{\mathbf{d}}_t^v], \\ \mathbf{x}_{t+1} &= \Pi_{\mathcal{X}}[\mathbf{x}_t - \alpha_t \hat{\mathbf{d}}_t^x]. \end{aligned}$$

previous results, we introduce regularity conditions for the smoothed functions in (17) and (18).

Inner and Outer Perturbed Gradient Variations: We define the gradient variations at the perturbed point as follows:

$$G_{v,T} := \sum_{t=2}^T (\chi_{1t} + \chi_{2t}), \quad G_{x,T} := \sum_{t=2}^T (\chi_{3t} + \chi_{4t}). \quad (24)$$

where $\mathbf{z}_t^+ := (\mathbf{x}_{t-1}, \mathbf{y}_{t-1} + \rho_v \mathbf{v}_{t-1})$, $\mathbf{z}_t^- := (\mathbf{x}_{t-1}, \mathbf{y}_{t-1} - \rho_v \mathbf{v}_{t-1})$, and

$$\begin{aligned} \chi_{1t} &:= \|\nabla_y g_t(\mathbf{z}_t^+) - \nabla_y g_{t-1}(\mathbf{z}_t^+)\|^2, \\ \chi_{2t} &:= \|\nabla_y g_t(\mathbf{z}_t^-) - \nabla_y g_{t-1}(\mathbf{z}_t^-)\|^2, \\ \chi_{3t} &:= \|\nabla_x g_t(\mathbf{z}_t^+) - \nabla_x g_{t-1}(\mathbf{z}_t^+)\|^2, \\ \chi_{4t} &:= \|\nabla_x g_t(\mathbf{z}_t^-) - \nabla_x g_{t-1}(\mathbf{z}_t^-)\|^2. \end{aligned}$$

Further, for simplicity of notation, we define

$$\begin{aligned} \hat{\Delta}_T &:= E_1 + V_T + D_T + G_{y,T}, \\ \hat{\Psi}_T &:= H_{2,T} + G_{v,T} + G_{x,T}, \end{aligned} \quad (25)$$

where $(V_T, H_{p,T})$ and (E_1, D_T) are defined in (10), and (14), respectively. Moreover, $G_{y,T}$ and $(G_{v,T}, G_{x,T})$, are defined in (12) and (24), respectively.

Theorem 4.2. Let $\{(f_t, g_t)\}_{t=1}^T$ be the sequence of functions presented to Algorithm 2, satisfying Assumptions 3.2–3.4 and 4.1. For all $t \in [T]$, let

$$\begin{aligned} \alpha_t &= \frac{1}{(d_1 + d_2)^{3/4}(c + t)^{1/3}}, & \beta_t &= c_\beta \alpha_t, & \delta_t &= c_\delta \alpha_t, \\ \gamma_{t+1} &= c_\gamma \alpha_t, & \eta_{t+1} &= c_\eta \alpha_t, & \lambda_{t+1} &= c_\lambda \alpha_t, \\ \rho_v^2 &= c_v \alpha_t, & \rho_r^2 &= \frac{1}{d_2^2 T}, & \rho_s^2 &= \frac{1}{d_1^2 T}, \\ b &= \frac{T^{1/3}}{(d_1 + d_2)^{3/2}}, & \bar{b} &= \frac{T^{2/3}}{(d_1 + d_2)^{3/4}}, \end{aligned} \quad (26)$$

where $c, c_\beta, c_\delta, c_\gamma, c_\eta, c_v$, and c_λ are specified in (232). Let $p = \ell_{f,0}/\mu_g$ for the set \mathcal{Z}_p defined in (8). Then, Algorithm 2 guarantees:

$$\begin{aligned} \text{BL-Reg}_T &\leq \mathcal{O}\left((d_1 + d_2)^{3/4} T^{1/3} \left(\hat{\sigma}^2 + \hat{\Delta}_T\right)\right. \\ &\quad \left.+ (d_1 + d_2)^{3/2} T^{2/3} \hat{\Psi}_T\right). \end{aligned}$$

where $\hat{\sigma}^2$ and $(\hat{\Delta}_T, \hat{\Psi}_T)$ are defined in (23) and (25).

Theorem 4.2 bounds the regret of Algorithm 2 without window-smoothing, based on the regularities in (25). We note that the average dynamic regret $\text{BL-Reg}_T/T \leq \mathcal{O}((d_1+d_2)^{3/4}T^{-2/3}(\hat{\sigma}^2 + \hat{\Delta}_T) + (d_1+d_2)^{3/2}T^{-1/3}\hat{\Psi}_T)$ remains sublinear under suitable conditions on $\hat{\Delta}_T$ and $\hat{\Psi}_T$.

Remark 4.3 (Regret Guarantee for Zeroth Order OBO). Theorem 4.2 provides the first regret guarantee for OBO with access only to noisy function evaluations of the leader and follower. The dimensional dependence $\mathcal{O}(d_1 + d_2)$ in Theorem 4.2 aligns with optimal results for simpler offline min-max problems (Huang et al., 2022). The bound also depends on the sample sizes b, \bar{b} and smoothing parameters ρ_v, ρ_r, ρ_s at each iteration.

Remark 4.4 (Improved Regret for OSO). Our dynamic regret for single-level non-stationary optimization is $\mathcal{O}((d_1 + d_2)^{3/4}T^{-2/3}(\hat{\sigma}^2 + E_1 + V_T + D_T))$, improving the result in Roy et al. (2022), which is $\mathcal{O}(T^{-1/2}\sigma^2\sqrt{d})$. Roy et al. (2022) proposed a zeroth-order stochastic gradient descent algorithm for unconstrained, non-convex, time-varying objective functions, achieving a regret bound of $\mathcal{O}(T^{-1/2}\sigma^2\sqrt{dW_T})$ using a two-point gradient estimator, where W_T bounds the nonstationarity. Additionally, Guan et al. (2023a) showed that the local regret for standard online stochastic gradient descent with the standard two-point gradient estimator (Agarwal et al., 2010) is $\mathcal{O}(T^{-1/2}d\sqrt{V_T})$.

5. Experimental Results

In this section, we provide experimental results on bilevel optimization-based black-box attacks on deep neural networks and parametric loss tuning for imbalanced data.

5.1. Bilevel Optimization-Based Black-Box Attacks

Deep neural network classifiers are vulnerable to adversarial examples—images subtly modified to mislead the classifier. These examples can deceive classifiers even without knowledge of the model, as seen in black-box adversarial attacks (BBAAs) (Chen et al., 2017; Liu et al., 2018b; Chen et al.,

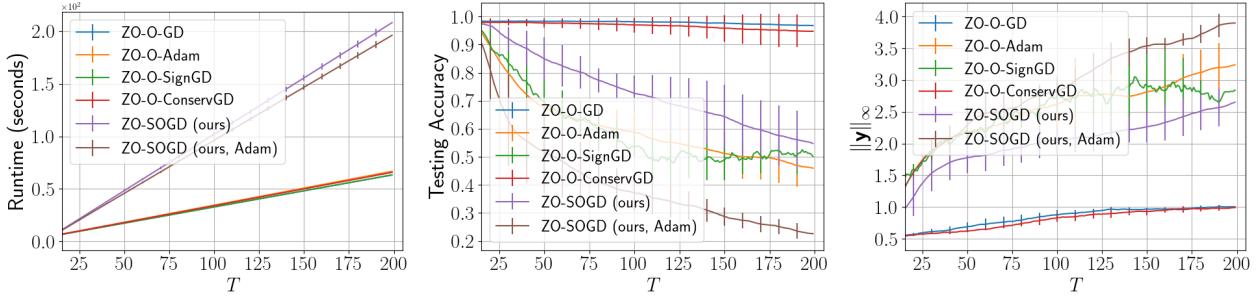


Figure 2. Performance comparison (mean \pm std) of optimizers including ZO-O-GD, ZO-O-Adam, ZO-O-SignSGD, ZO-O-ConservSGD, ZO-SOGD, and ZO-SOGD (Adam) on **online adversarial attack** for MNIST data across five runs.

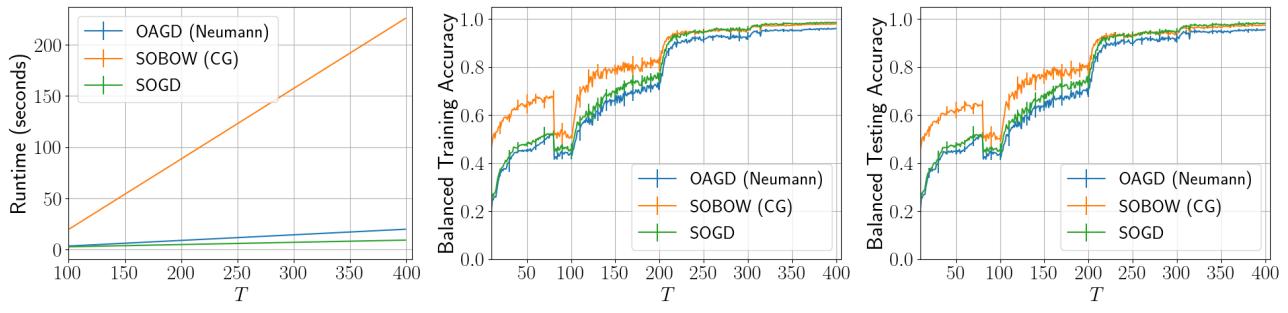


Figure 3. Performance comparison (mean \pm std) on **imbalanced loss tuning with distribution shift** for MNIST data across five runs between OGD (Zinkevich, 2003), OAGD (Tarzanagh et al., 2024), SOBOW (Lin et al., 2024), and our SOGD.

2019).

We first review the ZO single-level optimization for BBAA (Chen et al., 2017). Let (\mathbf{a}, b) denote a legitimate image $\mathbf{a} \in \mathbb{R}^d$ with true label $b \in \{1, 2, \dots, J\}$, where J is the total number of classes. Define $\mathbf{a}' = \mathbf{a} + \mathbf{y}$ as an adversarial example, with \mathbf{y} as the adversarial perturbation. Let $\mathcal{Y} := [-5, 5]^d \subset \mathbb{R}^d$, and $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$ denote the black-box attack loss. The goal of BBAA (Chen et al., 2017) is to design \mathbf{y} for images $\{\mathbf{a}_i\}_{i=1}^m$ by solving:

$$\min_{\mathbf{y} \in \mathcal{Y}} \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{a}_i + \mathbf{y}) + \lambda \|\mathbf{y}\|^2. \quad (27)$$

Here, $\lambda > 0$ is a hyperparameter balancing attack loss minimization and ℓ_2 regularization.

To adapt (27) to our OBO, consider OBO for supervised learning: at each timestep t , new samples $(\mathbf{a}_t, b_t) \in \mathcal{D}_t := \{\mathcal{D}_t^{\text{val}}, \mathcal{D}_t^{\text{tr}}\}$ are received, where $\mathbf{a}_t \in \mathbb{R}^{d_2}$ is the feature vector (image) and $b_t \in \mathbb{R}$ is the corresponding target. Note that the correct decision can change abruptly. We consider an S -stage scenario where $(\mathbf{x}_s^*, \mathbf{y}_s^*(\mathbf{x}_s^*))$ represents the best decisions for the s -th stage, for all $s \in [S]$.

$$\begin{aligned} \mathbf{x}_s^* &\in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \sum_{t=1}^{T_s} f(\mathbf{y}_s^*(\mathbf{x}); \mathcal{D}_t^{\text{val}}) \\ \text{s.t. } \mathbf{y}_s^*(\mathbf{x}) &\in \operatorname{argmin}_{\mathbf{y} \in \mathcal{Y}} \sum_{t=1}^{T_s} g(\mathbf{x}, \mathbf{y}; \mathcal{D}_t^{\text{tr}}), \end{aligned} \quad (28)$$

where

$$\begin{aligned} g(\mathbf{x}_t, \mathbf{y}_t; \mathcal{D}_t^{\text{tr}}) &= \frac{1}{|\mathcal{D}_t^{\text{tr}}|} \sum_{i \in \mathcal{D}_t^{\text{tr}}} \ell(\mathbf{a}_t^{(i)} + \mathbf{y}_t) \\ &+ \frac{1}{2} \sum_{\iota=1}^p e^{[\mathbf{x}_t]_\iota} [\mathbf{y}_t]_\iota^2, \end{aligned} \quad (29a)$$

and

$$f(\mathbf{y}_t(\mathbf{x}_t); \mathcal{D}_t^{\text{val}}) = \frac{1}{|\mathcal{D}_t^{\text{val}}|} \sum_{i \in \mathcal{D}_t^{\text{val}}} \ell(\mathbf{a}_t^{(i)} + \mathbf{y}_t). \quad (29b)$$

Here, $\{\mathbf{a}_t^{(i)}\}_{i \in \mathcal{D}_t^{\text{tr}}}$ and $\{\mathbf{a}_t^{(i)}\}_{i \in \mathcal{D}_t^{\text{val}}}$ are batches of training and validation samples at timestep t ; $\mathbf{a}_t^{(i)}$ is the i th sample in that batch; and $[\mathbf{x}_t]_\iota$ and $[\mathbf{y}_t]_\iota$ denote the ι th component of \mathbf{x}_t and \mathbf{y}_t , respectively.

We normalize the pixel values to \mathcal{Y} . For an untargeted attack, the loss in (29) is $\ell(\mathbf{a}_t') = \max\{Z(\mathbf{a}_t')_{b_t} - \max_{j \neq b_t} Z(\mathbf{a}_t')_j, -\kappa\}$, where $Z(\mathbf{a}_t')_j$ is the prediction score for class j given input $\mathbf{a}_t' = \mathbf{a}_t + \mathbf{y}_t$, and $\kappa > 0$ controls the confidence gap. In our experiments, we set $\kappa = 0$.

Eq. (28) introduces the first OBO formulation of BBAA. Using a vector $\mathbf{x} \in \mathbb{R}_+^d$ for hyperparameters instead of $\lambda \in \mathbb{R}_{++}$ in (27) enables finer control over model components, enhancing performance for complex models and heterogeneous data (Lorraine et al., 2020). For a fair comparison with single-level BBAA, we replace λ with a fixed

vector multiplied by each component of \mathbf{y} in (27).
 We compare our ZO-SOGD and ZO-SOGD (Adam) with the following competing methods in the online setting:

ZO-O-GD: A single-level method that updates \mathbf{y}_t with a fixed \mathbf{x} at each timestep using ZO gradient descent (Nesterov & Spokoiny, 2017).

ZO-O-Adam: A single-level method that updates \mathbf{y}_t with a fixed \mathbf{x} at each timestep using ZO Adam (Kingma & Ba, 2014; Chen et al., 2019).

ZO-O-SignSGD: A single-level method that updates \mathbf{y}_t with a fixed \mathbf{x} at each timestep using ZO SignSGD (Bernstein et al., 2018).

ZO-O-ConservSGD: A single-level method that updates \mathbf{y}_t with a fixed \mathbf{x} at each timestep using ZO Conservative SGD (Cutkosky & Boahen, 2019).

Note that ZO-SOGD (Adam) is a variant of our algorithm with an adaptive stepsize, similar to that of (Kingma & Ba, 2014).

We evaluated the proposed algorithms based on runtime, test accuracy on perturbed samples, and the infinity norm of \mathbf{y}_t . Figure 2 compares the methods. The left panel shows that ZO-SOGD has similar runtime to single-level baselines, despite outer-level optimization on \mathbf{x} . The middle panel shows that all methods' accuracy decreases as the adversarial attack \mathbf{y} strengthens, with ZO-SOGD outperforming ZO-O-GD and ZO-O-ConservGD, and ZO-SOGD (Adam) outperforming ZO-O-Adam and all baselines. The right panel shows that the increasing infinity norm of \mathbf{y}_t over time for all methods, which reduces accuracy. However, the perturbations remain unnoticeable with a max \mathbf{y}_t no larger than 4, demonstrating that ZO-SOGD achieves effective attacks with better performance than other methods.

5.2. Parametric Loss Tuning for Imbalanced Data

Imbalanced datasets are common in modern machine learning, causing challenges in generalization and fairness due to underrepresented classes and sensitive attributes. Deep NNs often overfit, seeming accurate and fair during training but performing poorly during testing. A common solution is designing a parametric training loss that balances accuracy and fairness while preventing overfitting (Li et al., 2021).

We consider an optimization problem similar to (28). For a new sample (\mathbf{a}_t, b_t) , the follower and leader incur a parametric and balanced cross-entropy loss, respectively:

$$\begin{aligned} g(\mathbf{x}_t, \mathbf{y}_t; \mathcal{D}_t^{\text{tr}}) &= -\log \frac{e^{\gamma_{b_t} [\mathbf{y}_t(\mathbf{a}_t)]_{b_t} + \Delta_{b_t}}}{\sum_{j=1}^J e^{\gamma_j [\mathbf{y}_t(\mathbf{a}_t)]_j + \Delta_j}}, \quad \text{and} \\ f(\mathbf{y}_t(\mathbf{x}_t); \mathcal{D}_t^{\text{val}}) &= -u_{b_t} \log \frac{e^{[\mathbf{y}_t(\mathbf{a}_t)]_{b_t}}}{\sum_{j=1}^J e^{[\mathbf{y}_t(\mathbf{a}_t)]_j}}. \end{aligned} \quad (30)$$

Here, $\mathbf{x}_t := (\Delta_j, \gamma_j)_{j=1}^J$ represents the logits adjustments, with j indexing the J classes, and u_j is the reciprocal of the proportion of samples from the j -th class to the total

number of samples (Li et al., 2021).

To clarify the notation in (30): $\mathbf{y}_t(\mathbf{x}_t)$ denotes the follower \mathbf{y}_t conditioned on the leader \mathbf{x}_t , while $[\mathbf{y}_t(\mathbf{a}_t)]_{b_t}$ represents the predicted logit for class b_t on sample \mathbf{a}_t . The backbone model for \mathbf{y}_t is a 4-layer CNN, leading to a nonconvex bilevel objective.

We compare SOGD with the following methods:

OAGD (Tarzanagh et al., 2024): A state-of-the-art static online bilevel gradient descent method using the Neumann series for hypergradient approximation.

SOBOW (Lin et al., 2024): A dynamic online bilevel gradient descent method using conjugate gradients (CG) for hypergradient approximation.

We conducted experiments on the MNIST (LeCun et al., 2010). We used a batch size of 64 per timestep. We evaluated cumulative runtime, balanced accuracy, and test accuracy, where balanced accuracy is the class-specific average accuracy:

$$\frac{1}{J} \sum_{j=1}^J \mathbb{P}_{\mathbf{a}_t \sim \mathcal{D}_j} [\text{argmax}_i ([\mathbf{y}_t(\mathbf{a}_t)]_i) = j],$$

with \mathcal{D}_j denoting the distribution over samples of class j (Li et al., 2021). Learning rates were tuned as $\beta_t = \delta_t = \beta \in \{0.001, 0.005, 0.01, 0.05, 0.1\}$ and $\alpha_t = \alpha \in \{0.0001, 0.0005, 0.001, 0.005, 0.01\}$ for all $t \in [T]$. The parameters $\gamma_t, \lambda_t, \eta_t$ were tuned as $\gamma_t = \lambda_t = \eta_t = \gamma \in \{0.9, 0.99, 0.999\}$. The Neumann series iterations in OAGD and CG iterations in SOBOW were set to 5.

We evaluated performance over 400 timesteps in four 100-timestep phases, transitioning from an imbalanced (0.4^i) to a balanced (0.8^i) distribution for each class ($i = 0, 1, \dots, 9$). Figure 3 (left) shows SOBOW's longer runtime due to CG complexity, while SOGD is the fastest with simultaneous updates. Figures 3 (middle, right) show accuracy gains as balance increases, with SOGD achieving competitive accuracy.

6. Conclusion

We introduced a novel online bilevel optimization (OBO) framework that overcomes the limitations of existing algorithms, which often rely on extensive oracle information and incur high computational costs. Our approach uses limited feedback and zeroth-order updates for efficient hypergradient estimation and simultaneous updates of decision variables, achieving sublinear bilevel regret without window smoothing. Experiments on online parametric loss tuning and black-box adversarial attacks confirm its effectiveness.

Impact Statements

This paper develops methods to advance online learning. While our work has societal implications, none require specific emphasis here.

References

- 440 Agarwal, A., Dekel, O., and Xiao, L. Optimal algorithms
441 for online convex optimization with multi-point bandit
442 feedback. In *Colt*, pp. 28–40. Citeseer, 2010.
- 443
- 444 Agarwal, N., Gonen, A., and Hazan, E. Learning in non-
445 convex games with an optimization oracle. In *Conference
446 on Learning Theory*, pp. 18–29. PMLR, 2019.
- 447
- 448 Aghasi, A. and Ghadimi, S. Fully zeroth-order bilevel
449 programming via gaussian smoothing. *arXiv preprint
450 arXiv:2404.00158*, 2024.
- 451
- 452 Allen-Zhu, Z. and Li, Y. Neon2: Finding local minima
453 via first-order oracles. *Advances in Neural Information
454 Processing Systems*, 31, 2018.
- 455
- 456 Bach, F. and Perchet, V. Highly-smooth zero-th order online
457 optimization. In *Conference on Learning Theory*, pp.
458 257–283. PMLR, 2016.
- 459
- 460 Bernstein, J., Wang, Y.-X., Azizzadenesheli, K., and Anand-
461 kumar, A. Signsdg: Compressed optimisation for non-
462 convex problems. In *International Conference on Ma-
463 chine Learning*, pp. 560–569. PMLR, 2018.
- 464
- Besbes, O., Gur, Y., and Zeevi, A. Non-stationary stochastic
465 optimization. *Operations research*, 63(5):1227–1244,
466 2015.
- 467
- Bohne, J., Rosenberg, D., Kazantsev, G., and Polak, P. On-
468 line nonconvex bilevel optimization with bregman diver-
469 gences. *arXiv preprint arXiv:2409.10470*, 2024.
- 470
- Bracken, J. and McGill, J. T. Mathematical programs with
471 optimization problems in the constraints. *Operations
472 Research*, 21(1):37–44, 1973.
- 473
- Bubeck, S., Stoltz, G., Szepesvári, C., and Munos, R. Online
474 optimization in x-armed bandits. *Advances in Neural
475 Information Processing Systems*, 21, 2008.
- 476
- Chen, P.-Y., Zhang, H., Sharma, Y., Yi, J., and Hsieh, C.-
477 J. Zoo: Zeroth order optimization based black-box at-
478 tacks to deep neural networks without training substitute
479 models. In *Proceedings of the 10th ACM workshop on
480 artificial intelligence and security*, pp. 15–26, 2017.
- 481
- Chen, T., Sun, Y., and Yin, W. Closing the gap: Tighter anal-
482 ysis of alternating stochastic gradient methods for bilevel
483 problems. *Advances in Neural Information Processing
484 Systems*, 34, 2021.
- 485
- Chen, X., Liu, S., Xu, K., Li, X., Lin, X., Hong, M., and
486 Cox, D. Zo-adammm: Zeroth-order adaptive momentum
487 method for black-box optimization. *Advances in neural
488 information processing systems*, 32, 2019.
- 489
- Crockett, C., Fessler, J. A., et al. Bilevel methods for im-
490 age reconstruction. *Foundations and Trends® in Signal
491 Processing*, 15(2-3):121–289, 2022.
- 492
- 493
- 494 Cutkosky, A. and Boahen, K. Anytime online-to-batch
495 conversions and the conservative algorithm. *Advances in
496 Neural Information Processing Systems*, 32, 2019.
- 497
- Dagréou, M., Ablin, P., Vaïter, S., and Moreau, T. A frame-
498 work for bilevel optimization that enables stochastic and
499 global variance reduction algorithms. *arXiv preprint
500 arXiv:2201.13409*, 2022.
- 501
- Dempe, S. *Foundations of bilevel programming*. Springer
502 Science & Business Media, 2002.
- 503
- Duchi, J. C., Jordan, M. I., Wainwright, M. J., and Wibisono,
504 A. Optimal rates for zero-order convex optimization: The
505 power of two function evaluations. *IEEE Transactions
506 on Information Theory*, 61(5):2788–2806, 2015.
- 507
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-
508 learning for fast adaptation of deep networks. In *Inter-
509 national Conference on Machine Learning*, pp. 1126–1135.
510 PMLR, 2017.
- 511
- Finn, C., Rajeswaran, A., Kakade, S., and Levine, S. Online
512 meta-learning. In *International Conference on Machine
513 Learning*, pp. 1920–1930. PMLR, 2019.
- 514
- Flaxman, A. D., Kalai, A. T., and McMahan, H. B. On-
515 line convex optimization in the bandit setting: gradient
516 descent without a gradient. *arXiv preprint cs/0408007*,
517 2004.
- 518
- Franceschi, L., Donini, M., Frasconi, P., and Pontil, M.
519 Forward and reverse gradient-based hyperparameter opti-
520 mization. In *International Conference on Machine Learn-
521 ing*, pp. 1165–1173. PMLR, 2017.
- 522
- Franceschi, L., Frasconi, P., Salzo, S., Grazzi, R., and Pontil,
523 M. Bilevel programming for hyperparameter optimiza-
524 tion and meta-learning. In *International Conference on
525 Machine Learning*, pp. 1568–1577. PMLR, 2018.
- 526
- Gao, X., Li, X., and Zhang, S. Online learning with non-
527 convex losses and non-stationary regret. In *International
528 Conference on Artificial Intelligence and Statistics*, pp.
529 235–243. PMLR, 2018.
- 530
- Ghadimi, S. and Lan, G. Stochastic first-and zeroth-order
531 methods for nonconvex stochastic programming. *SIAM
532 journal on optimization*, 23(4):2341–2368, 2013.
- 533
- Ghadimi, S. and Wang, M. Approximation methods for
534 bilevel programming. *arXiv preprint arXiv:1802.02246*,
535 2018.
- 536
- Ghadimi, S., Lan, G., and Zhang, H. Mini-batch stochastic
537 approximation methods for nonconvex stochastic com-
538 posite optimization. *Mathematical Programming*, 155
539 (1-2):267–305, 2016.
- 540
- Goel, G., Lin, Y., Sun, H., and Wierman, A. Beyond online
541 balanced descent: An optimal algorithm for smoothed
542 online optimization. *Advances in Neural Information
543 Processing Systems*, 32, 2019.
- 544

- 495 Guan, Z., Zhou, Y., and Liang, Y. On the hardness of online
 496 nonconvex optimization with single oracle feedback. In
 497 *The Twelfth International Conference on Learning Representations*, 2023a.
 498
- 499 Guan, Z., Zhou, Y., and Liang, Y. Online nonconvex optimiza-
 500 tion with limited instantaneous oracle feedback. In
 501 *The Thirty Sixth Annual Conference on Learning Theory*,
 502 pp. 3328–3355. PMLR, 2023b.
- 503 Hallak, N., Mertikopoulos, P., and Cevher, V. Regret
 504 minimization in stochastic non-convex learning via a
 505 proximal-gradient approach. In *International Conference
 506 on Machine Learning*, pp. 4008–4017. PMLR, 2021.
- 507 Hansen, P., Jaumard, B., and Savard, G. New branch-and-
 508 bound rules for linear bilevel programming. *SIAM Journal
 509 on scientific and Statistical Computing*, 13(5):1194–
 510 1217, 1992.
- 511 Hazan, E. Introduction to online convex optimization. *Founda-
 512 tions and Trends® in Optimization*, 2(3-4):157–325,
 513 2016a. URL <http://ocobook.cs.princeton.edu/OCOb book.pdf>.
- 514 Hazan, E. Introduction to online convex optimization. *Founda-
 515 tions and Trends in Optimization*, 2(3-4):157–325,
 516 2016b.
- 517 Hazan, E., Agarwal, A., and Kale, S. Logarithmic regret al-
 518 gorithms for online convex optimization. *Machine Learn-
 519 ing*, 69(2):169–192, 2007.
- 520 Hazan, E., Singh, K., and Zhang, C. Efficient regret min-
 521 imization in non-convex games. In *International Confer-
 522 ence on Machine Learning*, pp. 1433–1441. PMLR,
 523 2017.
- 524 Héliou, A., Martin, M., Mertikopoulos, P., and Rahier, T.
 525 Online non-convex optimization with imperfect feedback.
 526 *Advances in Neural Information Processing Systems*, 33:
 527 17224–17235, 2020.
- 528 Héliou, A., Martin, M., Mertikopoulos, P., and Rahier,
 529 T. Zeroth-order non-convex learning via hierarchical
 530 dual averaging. In *International Conference on Machine
 531 Learning*, pp. 4192–4202. PMLR, 2021.
- 532 Huang, F., Gao, S., Pei, J., and Huang, H. Accelerated
 533 zeroth-order and first-order momentum methods from
 534 mini to minimax optimization. *Journal of Machine Learn-
 535 ing Research*, 23(36):1–70, 2022.
- 536 Huang, Y., Cheng, Y., Liang, Y., and Huang, L. Online min-
 537 max problems with non-convexity and non-stationarity.
 538 *Transactions on Machine Learning Research*.
- 539 Huang, Y., Cheng, Y., Liang, Y., and Huang, L. Online min-
 540 max problems with non-convexity and non-stationarity.
 541 *Transactions on Machine Learning Research*, 2023.
- 542 Ji, K., Wang, Z., Zhou, Y., and Liang, Y. Improved zeroth-
 543 order variance reduced algorithms and analysis for non-
 544 convex optimization. In *International conference on ma-
 545 chine learning*, pp. 3100–3109. PMLR, 2019.
- 546 Ji, K., Yang, J., and Liang, Y. Bilevel optimization: Con-
 547 vergence analysis and enhanced design. In *International
 548 Conference on Machine Learning*, pp. 4882–4892.
 549 PMLR, 2021.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic
 optimization. In *International Conference on Learning
 Representations*, 2014.
- Kleinberg, R., Slivkins, A., and Upfal, E. Multi-armed
 bandits in metric spaces. In *Proceedings of the fortieth
 annual ACM symposium on Theory of computing*, pp.
 681–690, 2008.
- Krichene, W., Balandat, M., Tomlin, C., and Bayen, A.
 The hedge algorithm on a continuum. In *International
 Conference on Machine Learning*, pp. 824–832. PMLR,
 2015.
- LeCun, Y., Cortes, C., and Burges, C. Mnist hand-
 written digit database. *ATT Labs [Online]. Available:
<http://yann.lecun.com/exdb/mnist>*, 2, 2010.
- Li, M., Zhang, X., Thrampoulidis, C., Chen, J., and Oymak,
 S. Autobalance: Optimized loss functions for imbal-
 anced data. *Advances in Neural Information Processing
 Systems*, 34:3163–3177, 2021.
- Lin, S., Sow, D., Ji, K., Liang, Y., and Shroff, N. Non-
 convex bilevel optimization with time-varying objective
 functions. *Advances in Neural Information Processing
 Systems*, 36, 2024.
- Liu, H., Simonyan, K., and Yang, Y. Darts: Differentiable
 architecture search. *arXiv preprint arXiv:1806.09055*,
 2018a.
- Liu, S., Chen, J., Chen, P.-Y., and Hero, A. Zeroth-order on-
 line alternating direction method of multipliers: Conver-
 gence analysis and applications. In *International Confer-
 ence on Artificial Intelligence and Statistics*, pp. 288–297.
 PMLR, 2018b.
- Lorraine, J., Vicol, P., and Duvenaud, D. Optimizing mil-
 lions of hyperparameters by implicit differentiation. In
*International conference on artificial intelligence and
 statistics*, pp. 1540–1552. PMLR, 2020.
- Luo, L., Ye, H., Huang, Z., and Zhang, T. Stochastic re-
 cursive gradient descent ascent for stochastic nonconvex-
 strongly-concave minimax problems. *Advances in Neural
 Information Processing Systems*, 33:20566–20577, 2020.
- Lv, Y., Hu, T., Wang, G., and Wan, Z. A penalty func-
 tion method based on kuhn–tucker condition for solving
 linear bilevel programming. *Applied Mathematics and
 Computation*, 188(1):808–813, 2007.
- Nesterov, Y. Smooth minimization of non-smooth functions.
Mathematical programming, 103:127–152, 2005.
- Nesterov, Y. and Spokoiny, V. Random gradient-free mini-
 mization of convex functions. *Foundations of Computa-
 tional Mathematics*, 17(2):527–566, 2017.

- 550 Roy, A., Balasubramanian, K., Ghadimi, S., and Mohapatra,
 551 P. Stochastic zeroth-order optimization under nonstationar-
 552 ity and nonconvexity. *Journal of Machine Learning
 553 Research*, 23(64):1–47, 2022.
- 554 Shalev-Shwartz, S. et al. Online learning and online con-
 555 vex optimization. *Foundations and trends in Machine
 556 Learning*, 4(2):107–194, 2011.
- 557 Shamir, O. An optimal algorithm for bandit and zero-order
 558 convex optimization with two-point feedback. *Journal of
 559 Machine Learning Research*, 18(52):1–11, 2017.
- 560 Sow, D., Ji, K., and Liang, Y. On the convergence theory
 561 for hessian-free bilevel algorithms. *Advances in Neural
 562 Information Processing Systems*, 35:4136–4149, 2022.
- 563 Stackelberg, H. v. Theory of the market economy. *Oxford
 564 University Press*, 1952.
- 565 Stadie, B., Zhang, L., and Ba, J. Learning intrinsic rewards
 566 as a bi-level optimization problem. In *Conference on Un-
 567 certainty in Artificial Intelligence*, pp. 111–120. PMLR,
 568 2020.
- 569 Suggala, A. S. and Netrapalli, P. Online non-convex learn-
 570 ing: Following the perturbed leader is optimal. In *Algo-
 571 rithmic Learning Theory*, pp. 845–861. PMLR, 2020.
- 572 Tarzanagh, D. A., Nazari, P., Hou, B., Shen, L., and Balzano,
 573 L. Online bilevel optimization: Regret analysis of online
 574 alternating gradient methods. In *International Conference
 575 on Artificial Intelligence and Statistics*, pp. 2854–2862.
 576 PMLR, 2024.
- 577 Wang, Z., Balasubramanian, K., Ma, S., and Razaviyayn,
 578 M. Zeroth-order algorithms for nonconvex minimax
 579 problems with improved complexities. *arXiv preprint
 580 arXiv:2001.07819*, 2020.
- 581 Zhang, Y., Zhou, Y., Ji, K., and Zavlanos, M. M. Boosting
 582 one-point derivative-free online optimization via residual
 583 feedback. *arXiv preprint arXiv:2010.07378*, 2020.
- 584 Zhou, W., Li, Y., Yang, Y., Wang, H., and Hospedales,
 585 T. Online meta-critic learning for off-policy actor-critic
 586 methods. *Advances in Neural Information Processing
 587 Systems*, 33:17662–17673, 2020.
- 588 Zinkevich, M. Online convex programming and generalized
 589 infinitesimal gradient ascent. In *Proceedings of the 20th
 590 international conference on machine learning (icml-03)*,
 591 pp. 928–936, 2003.
- 592
- 593
- 594
- 595
- 596
- 597
- 598
- 599
- 600
- 601
- 602
- 603
- 604

A. Related Work

BO was introduced in game theory by (Stackelberg, 1952) and modeled mathematically in (Bracken & McGill, 1973). Initial works (Hansen et al., 1992; Lv et al., 2007) reduced it to single-level optimization. Recently, gradient-based approaches have gained popularity for their simplicity and efficacy (Franceschi et al., 2017; Ghadimi & Wang, 2018; Ji et al., 2021; Chen et al., 2021), though they assume offline objectives.

OBO was initiated by Tarzanagh et al. (2024), proposing the OAGD method with regret bounds. (Huang et al., 2023) developed algorithms for online minimax optimization, special cases of OBO with local regret guarantees. (Lin et al., 2024) introduced SOBOW, a single-loop optimizer using window-smoothed functions and multiple CGs for nonconvex-strongly-convex cases. Unlike these works, we propose using *projected gradient* as a more general performance measure for constrained objectives, focusing on the original functions and their regret; See Table 1 for a comparison.

Single-Level Regret Minimization. Single-level online optimization predominantly focuses on convex problems, either with static or dynamic convex regret minimization (Zinkevich, 2003; Hazan, 2016a; Shalev-Shwartz et al., 2011). Non-convex online optimization (Hazan et al., 2017; Guan et al., 2023b;a) poses greater challenges than its convex counterparts (Shalev-Shwartz et al., 2011; Zinkevich, 2003; Hazan et al., 2007; Besbes et al., 2015). Notable contributions in this field include adversarial multi-armed bandit algorithms (Bubeck et al., 2008; Héliou et al., 2020; 2021; Krichene et al., 2015) and the Follow-the-Perturbed-Leader approach (Agarwal et al., 2019; Kleinberg et al., 2008; Suggala & Netrapalli, 2020). Hazan et al. (Hazan et al., 2017) introduced window-smoothed local regret for gradient averaging in non-convex models, which Hallak et al. (Hallak et al., 2021) extended to non-smooth, non-convex problems. Inspired by their work, we employ local regret for Online Bandit Optimization (OBO) without window-smoothing.

Zereth-Order Optimization. Single-Level ZO Optimization has been widely studied in both offline (Ghadimi & Lan, 2013; Duchi et al., 2015; Agarwal et al., 2010; Nesterov & Spokoiny, 2017) and online settings (Liu et al., 2018b; Guan et al., 2023a;b; Zhang et al., 2020; Bach & Perchet, 2016). We next review closely related work. Liu et al. (Liu et al., 2018b) proposed ZOO-ADMM, a gradient-free online optimization algorithm utilizing ADMM. Guan et al. (Guan et al., 2023b) studied online non-convex optimization with limited oracle feedback. Research on online non-convex optimization with bandit feedback includes work by Heliou et al. (Héliou et al., 2020), which established bounds on global static and dynamic regret using dual averaging, further refined in (Héliou et al., 2021). Gao et al. (Gao et al., 2018) extended these ideas to ZO algorithms. Flaxman et al. (Flaxman et al., 2004) provided algorithms for bandit online optimization of convex functions using ZO gradient approximation. Our work closely relates to (Sow et al., 2022), which proposes a Hessian-free method approximating the Jacobian matrix using a ZO method based on finite differences of gradients. In contrast, our method uses function oracles to approximate both the Hessian and gradients and is derivative-free. We also point out the recent work (Aghasi & Ghadimi, 2024) on ZO stochastic algorithms for solving bilevel problems when neither the upper/lower objective values nor their unbiased gradient estimates are available. Their approach, limited to the *offline* setting, does not include numerical results, thus leaving its practical efficiency unclear.

B. Additional Preliminaries and Notations

B.1. Preliminary Lemmas

We first provide several useful lemmas for the main proofs.

Definition B.1 (Projected gradient (Ghadimi et al., 2016)). Let $\mathcal{X} \subset \mathbb{R}^{d_1}$ be a closed convex set. Then, the projected gradient for any $\alpha_t > 0$ and $\mathbf{p} \in \mathbb{R}^{d_1}$ is defined as

$$\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p}) := \frac{1}{\alpha_t} (\mathbf{x} - \mathbf{x}^+), \quad (31a)$$

where

$$\mathbf{x}^+ = \Pi_{\mathcal{X}}(\mathbf{x} - \alpha_t \mathbf{p}), \quad (31b)$$

and $\Pi_{\mathcal{X}}[\cdot]$ denotes the orthogonal projection operator onto set \mathcal{X} .

Lemma B.2. Goel et al. (2019, Lemma 13) If $f : \mathcal{X} \rightarrow \mathbb{R}$ is a μ_f -strongly convex function with respect to some norm $\|\cdot\|$, and \mathbf{x}^* is the minimizer of f (i.e. $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$), then we have $\forall \mathbf{x} \in \mathcal{X}$,

$$\frac{\mu_f}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 \leq f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{1}{2\mu_f} \|\nabla f(\mathbf{x})\|^2.$$

Lemma B.3. Suppose $f(\mathbf{x})$ is L -smooth, and $\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$. Then, we can upper bound the magnitude of the

gradient at any given point $\mathbf{x} \in \mathbb{R}^d$ in terms of the objective sub optimality at \mathbf{x} , as follows:

$$\frac{1}{2L} \|\nabla f(\mathbf{x})\|^2 \leq f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{L}{2} \|\mathbf{x} - \mathbf{x}^*\|^2. \quad (32)$$

Lemma B.4. For any set of vectors $\{\mathbf{x}_i\}_{i=1}^m$ with $\mathbf{x}_i \in \mathbb{R}^d$, we have

$$\left\| \sum_{i=1}^m \mathbf{x}_i \right\|^2 \leq m \sum_{i=1}^m \|\mathbf{x}_i\|^2.$$

Lemma B.5. For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, the following holds for any $c > 0$:

$$\|\mathbf{x} + \mathbf{y}\|^2 \leq (1 + c) \|\mathbf{x}\|^2 + \left(1 + \frac{1}{c}\right) \|\mathbf{y}\|^2, \text{ and} \quad (33)$$

$$\|\mathbf{x} - \mathbf{y}\|^2 \geq (1 - c) \|\mathbf{x} - \mathbf{z}\|^2 + \left(1 - \frac{1}{c}\right) \|\mathbf{z} - \mathbf{y}\|^2. \quad (34)$$

We provide a set of auxiliary lemmas that will be used in establishing the proofs for the main theorems.

Lemma B.6. Ghadimi et al. (2016, Proposition 1) Let $\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p})$ be defined in Definition B.1. Then, for any \mathbf{p}_1 and \mathbf{p}_2 in \mathbb{R}^d , we

$$\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p}_1) - \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p}_2)\| \leq \|\mathbf{p}_1 - \mathbf{p}_2\|.$$

Lemma B.7. Hazan et al. (2017, Proposition 2.4) Let $\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p})$ be the projected gradient as per Definition B.1. For any $\mathbf{x}, \mathbf{p}_1, \mathbf{p}_2 \in \mathbb{R}^d$ and $\alpha_t > 0$ it holds that

$$\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p}_1 + \mathbf{p}_2)\| \leq \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p}_1)\| + \|\mathbf{p}_2\|.$$

Lemma B.8. Let $\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p})$ be as given in Definition B.1. Then, for any $\mathbf{p} \in \mathbb{R}^d$ and $\alpha_t > 0$, we have

$$\langle \mathbf{p}, \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p}) \rangle \geq \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p})\|^2.$$

Proof. By the definition of \mathbf{x}^+ , the optimality condition of (31b) is

$$\left\langle \mathbf{p} + \frac{1}{\alpha_t} (\mathbf{x}^+ - \mathbf{x}), \mathbf{z} - \mathbf{x}^+ \right\rangle \geq 0, \quad \forall \mathbf{z} \in \mathcal{X}.$$

Letting $\mathbf{z} = \mathbf{x}$, we obtain

$$\langle \mathbf{p}, \mathbf{x} - \mathbf{x}^+ \rangle \geq \frac{1}{\alpha_t} \langle \mathbf{x} - \mathbf{x}^+, \mathbf{x} - \mathbf{x}^+ \rangle,$$

which can be rearranged to

$$\begin{aligned} \langle \mathbf{p}, \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p}) \rangle &= \frac{1}{\alpha_t} \langle \mathbf{p}, \mathbf{x} - \mathbf{x}^+ \rangle \geq \frac{1}{\alpha_t^2} \langle \mathbf{x} - \mathbf{x}^+, \mathbf{x} - \mathbf{x}^+ \rangle \\ &= \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}; \mathbf{p})\|^2. \end{aligned}$$

□

B.2. Examples

Theorem 3.6 achieves sublinear bilevel regret when the variations V_T and $H_{2,T}$ are both $o(T)$. Below, we provide some examples of online optimization in both single-level and bilevel settings to illustrate when this occurs.

Example B.9. Consider function $f_t(\mathbf{x}) = \|\mathbf{A}_t \mathbf{x} - \mathbf{b}_t\|^2$, where $\mathbf{A}_t = [1, 0; 0, 1 + \frac{1}{t}]$, $\mathbf{x} = \mathbf{b}_t = (1, 1)$. Then, $V_T :=$

715 $\sum_{t=2}^T \max_{\mathbf{x}} |f_t(\mathbf{x}) - f_{t-1}(\mathbf{x})| = \sum_{t=2}^T \left| \left(\frac{1}{t} \right)^2 - \left(\frac{1}{t-1} \right)^2 \right|$. By $a^2 - b^2 = (a-b)(a+b)$, we have
 716

$$\begin{aligned} V_T &= \sum_{t=2}^T \left| \left(\frac{1}{t} - \frac{1}{t-1} \right) - \left(\frac{1}{t} + \frac{1}{t-1} \right) \right| \\ &= \sum_{t=2}^T \left| \left(\frac{t-1-t}{t(t-1)} \right) - \left(\frac{1}{t} + \frac{1}{t-1} \right) \right| \\ &= \sum_{t=2}^T \left| -\frac{1}{t(t-1)} - \left(\frac{1}{t} + \frac{1}{t-1} \right) \right| \\ &= \sum_{t=2}^T \left| \frac{1}{t(t-1)} \right| \left| \frac{t-1+t}{t(t-1)} \right| \\ &= \sum_{t=2}^T \left| \frac{2}{t(t-1)^2} \right|. \end{aligned}$$

732 Then, $V_T \leq \sum_{t=2}^T \frac{2}{t^3} \approx \int_2^T \frac{2}{t^3} dt = \frac{1}{4} - \frac{1}{T^2}$. As $T \rightarrow \infty$, V_T becomes bounded and approaches a constant value, indicating
 733 that V_T grows slower than T itself.
 734

735
 736
 737 *Example B.10.* Let $f_t(\mathbf{x}) = (-\frac{1}{T}, 0, 0, 0)$ if t is even, and $f_t(\mathbf{x}) = (0, -\frac{1}{T}, 0, 0)$ if t is odd. Then, $V_T =$
 738 $\sum_{t=2}^T \max_{\mathbf{x}} |f_t(\mathbf{x}) - f_{t-1}(\mathbf{x})| = \mathcal{O}(1)$.

739
 740
 741 *Example B.11.* Let $x \in \mathcal{X} = [-1, 1] \subset \mathbb{R}$, $y \in \mathbb{R}$, and consider a sequence of quadratic cost functions
 742

$$\begin{aligned} f_t(x, y) &= \frac{1}{2} \left(x + 2a_t^{(1)} \right)^2 + \frac{1}{2} \left(y - a_t^{(2)} \right)^2, \\ g_t(x, y) &= \frac{1}{2} y^2 - \left(x - a_t^{(2)} \right) y, \end{aligned}$$

747 where $a_t^{(1)} = 1/t$ and $a_t^{(2)} = 1/\sqrt{t}$ for all $t \in [T]$.
 748

749 We have

$$y_t^*(x) = x - a_t^{(2)}.$$

750
 751 We have
 752

$$\begin{aligned} &f_t(x, y_t^*(x)) - f_{t-1}(x, y_{t-1}^*(x)) \\ &= \frac{1}{2} \left[\left(x + 2a_t^{(1)} \right)^2 - \left(x + 2a_{t-1}^{(1)} \right)^2 \right] + \frac{1}{2} \left[\left(y_t^*(x) - a_t^{(2)} \right)^2 - \left(y_{t-1}^*(x) - a_{t-1}^{(2)} \right)^2 \right] \\ &= \frac{1}{2} \left[\left(x^2 + 4xa_t^{(1)} + 4(a_t^{(1)})^2 \right) - \left(x^2 + 4xa_{t-1}^{(1)} + 4(a_{t-1}^{(1)})^2 \right) \right] \\ &\quad + \frac{1}{2} \left[\left((x - a_t^{(2)})^2 - 2(x - a_t^{(2)})a_t^{(2)} + (a_t^{(2)})^2 \right) - \left((x - a_{t-1}^{(2)})^2 - 2(x - a_{t-1}^{(2)})a_{t-1}^{(2)} + (a_{t-1}^{(2)})^2 \right) \right] \\ &= 2x \left(a_t^{(1)} - a_{t-1}^{(1)} - a_t^{(2)} + a_{t-1}^{(2)} \right) + 2 \left((a_t^{(1)})^2 - (a_{t-1}^{(1)})^2 + (a_t^{(2)})^2 - (a_{t-1}^{(2)})^2 \right). \end{aligned}$$

763 Taking the maximum over x and using $x \in [-1, 1]$:

$$\begin{aligned} \sup_x |f_t(x, y_t^*(x)) - f_{t-1}(x, y_{t-1}^*(x))| &= 2 \left| a_t^{(1)} - a_{t-1}^{(1)} \right| + 2 \left| -a_t^{(2)} + a_{t-1}^{(2)} \right| \\ &\quad + 2 \left| (a_t^{(1)})^2 - (a_{t-1}^{(1)})^2 \right| + 2 \left| (a_t^{(2)})^2 - (a_{t-1}^{(2)})^2 \right|. \end{aligned}$$

770 Since $a_t^{(1)} = 1/t$ and $a_t^{(2)} = 1/\sqrt{t}$ for all $t \in [T]$, then we have

$$\begin{aligned} 771 \quad |a_t^{(1)} - a_{t-1}^{(1)}| &\approx \frac{1}{t^2}, \quad |a_t^{(2)} - a_{t-1}^{(2)}| \approx \frac{1}{2t^{3/2}}, \\ 772 \quad |(a_t^{(1)})^2 - (a_{t-1}^{(1)})^2| &\approx \frac{1}{t^3}, \quad |(a_t^{(2)})^2 - (a_{t-1}^{(2)})^2| \approx \frac{1}{t^2}. \end{aligned}$$

776 Then, we get

$$778 \quad V_T := \sum_{t=2}^T \sup_x |f_t(x, y_t^*(x)) - f_{t-1}(x, y_{t-1}^*(x))| = \sum_{t=2}^T \left(\frac{2}{t^2} + \frac{1}{2t^{3/2}} + \frac{1}{t^3} \right).$$

781 The series $\sum_{t=2}^T (\frac{2}{t^2} + \frac{1}{2t^{3/2}} + \frac{1}{t^3})$ converges, implying $V_T = \mathcal{O}(1)$. Moreover, we have

$$\begin{aligned} 783 \quad H_{2,T} &= \sum_{t=2}^T \sup_x \|y_t^*(x) - y_{t-1}^*(x)\|^2 = \sum_{t=2}^T \sup_x \|x - a_t^{(2)} - x + a_{t-1}^{(2)}\|^2 \\ 784 \quad &= \sum_{t=2}^T | -a_t^{(2)} + a_{t-1}^{(2)} |^2 = \sum_{t=2}^T |a_t^{(2)} - a_{t-1}^{(2)} |^2 \approx \sum_{t=2}^T \frac{1}{4t^3}, \end{aligned}$$

789 which implies $H_{2,T} = \mathcal{O}(1)$.

790 To achieve $V_T = o(T)$, the changes in the cost functions $f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))$ and $\mathbf{y}_t^*(\mathbf{x})$ should decay to zero faster than $\mathcal{O}(1/t)$.
 791 For example, if the coefficients in the functions change as $\mathcal{O}(1/t^a)$ with $a > 1$, then the cumulative sum over T will be
 792 $o(T)$. When $f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))$ and $\mathbf{y}_t^*(\mathbf{x})$ decay as $\mathcal{O}(1/\sqrt{t})$, then the total variation grows at most as $\mathcal{O}(\sqrt{T})$.

796 C. Proof of Regret Bounds for Simultaneous Online Gradient Descent (SOGD)

797 **Proof Roadmap.** We introduce Lemma C.2, which quantifies the error between the approximated direction of the
 798 momentum-based gradient estimator, \mathbf{d}_t^y , and the true direction, $\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)$, at each iteration. To bound the error of
 799 the lower-level variable, we provide Lemma C.4, which captures the gap $\|\mathbf{y}_{t+1} - \mathbf{y}_t^*(\mathbf{x}_t)\|^2$ and incorporates the error
 800 introduced in Lemma C.2. Moreover, we provide Lemma C.5, which quantifies the error between the approximated direction
 801 of the momentum-based gradient estimator, \mathbf{d}_t^y , and the true direction, $\nabla_y^2 g_t(\mathbf{z}_t) \mathbf{v}_t + \nabla_y f_t(\mathbf{z}_t)$, at each iteration. To bound
 802 the error of the system solution, we provide Lemma C.8, which captures the gap $\|\mathbf{v}_{t+1} - \mathbf{v}_t^*(\mathbf{x}_t)\|^2$ and incorporates the
 803 error introduced in Lemma C.5. Moreover, we provide Lemma C.9, which quantifies the error between the approximated
 804 direction of the momentum-based hypergradient estimator, \mathbf{d}_t^x , and the true direction, $\nabla_x f_t(\mathbf{z}_t) + \nabla_{xy}^2 g_t(\mathbf{z}_t) \mathbf{v}_t$, at each
 805 iteration. We also present Lemma C.11, which provides an upper bound for the projection mapping and relates to the three
 806 errors discussed in Lemmas C.4, C.8, and C.9. Finally, by combining these lemmas and appropriately setting the parameters,
 807 we achieve the desired result.

808 C.1. Proof of Lemma 3.1

809 *Proof.* SOBOW (Lin et al., 2024) has estimated the hypergradient as the weighted average of previous ones over a sliding
 810 window of size w for a given $\mathcal{B}_t := \{\xi_1, \dots, \xi_b\}$ drawn i.i.d. from the distribution \mathcal{D}_f , as follows:

$$813 \quad \hat{\nabla} F_{t,\nu}(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) = \frac{1}{W} \sum_{i=0}^{w-1} \nu^i \hat{\nabla} f_{t-i}(\mathbf{x}_{t-i}, \mathbf{y}_{t-i}; \mathcal{B}_{t-i}),$$

816 with $W = \sum_{i=0}^{w-1} \nu^i$, $\nu \in (0, 1)$. Let $\nu = 1 - \eta$ for $\eta \in (0, 1)$.

817 Then, the above equality is equivalent to

$$818 \quad \hat{\nabla} F_{t,\nu}(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) = \frac{1}{W} \sum_{j=t-w+1}^t (1 - \eta)^{t-j} \hat{\nabla} f_j(\mathbf{x}_j, \mathbf{y}_j; \mathcal{B}_j), \tag{35}$$

823 with $W = \sum_{j=t-w+1}^t (1 - \eta)^{t-j}$.

825 Let $\hat{\mathbf{d}}_t^{\mathbf{x}} := \hat{\nabla}F_{t,\nu}(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t)$. Then (35) is equivalent to

$$827 \quad \hat{\mathbf{d}}_t^{\mathbf{x}} = \frac{1}{W} \hat{\nabla}f_t(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) + (1 - \eta)\hat{\mathbf{d}}_{t-1}^{\mathbf{x}} - \frac{(1 - \eta)^w}{W} \hat{\nabla}f_{t-w}(\mathbf{x}_{t-w}, \mathbf{y}_{t-w}; \mathcal{B}_{t-w}), \quad (36)$$

829 with $f_i(\cdot) = 0$ for all $i \leq 0$.

830 If $w = t$ and $W = \frac{1}{\eta}$, then, we have

$$832 \quad \hat{\mathbf{d}}_t^{\mathbf{x}} = \eta \hat{\nabla}f_t(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_t) + (1 - \eta)\hat{\mathbf{d}}_{t-1}^{\mathbf{x}}.$$

833 \square

836 C.2. Bounds on the Inner Decision Variable

837 We first provide a lemma that characterizes the Lipschitz continuity of approximate gradients, inner, and system solutions.

838 **Lemma C.1.** Under Assumptions 3.2 and 3.3, for all $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, and the search directions $\{\mathbf{d}_t^{\mathbf{x}}\}_{t=1}^T$ and $\{\mathbf{d}_t^{\mathbf{v}}\}_{t=1}^T$ generated by Algorithm 1, we have

$$842 \quad \|\mathbf{d}_t^{\mathbf{x}} - \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\|^2 \leq M_f^2 \left(\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \right), \quad (37a)$$

$$844 \quad \|\mathbf{d}_t^{\mathbf{v}}\|^2 \leq M_{\mathbf{v}}^2 \left(\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \right), \quad (37b)$$

$$846 \quad \|\nabla f_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) - \nabla f_t(\mathbf{x}', \mathbf{y}_t^*(\mathbf{x}'))\| \leq L_f \|\mathbf{x} - \mathbf{x}'\|, \quad (37c)$$

$$847 \quad \|\mathbf{y}_t^*(\mathbf{x}) - \mathbf{y}_t^*(\mathbf{x}')\| \leq L_y \|\mathbf{x} - \mathbf{x}'\|, \quad (37d)$$

$$848 \quad \|\mathbf{v}_t^*(\mathbf{x}) - \mathbf{v}_t^*(\mathbf{x}')\| \leq L_v \|\mathbf{x} - \mathbf{x}'\|, \quad (37e)$$

849 where M_f , $M_{\mathbf{v}}$, and (L_y, L_v, L_f) are defined in (40), (41), and (42), respectively.

855 *Proof.* We first show (37a).

856 Using Assumptions 3.2 and 3.3, we have $\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \succeq \mu_g$, and

$$857 \quad \|\mathbf{v}_t^*(\mathbf{x}_t)\| = \|(\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))^{-1} \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\| \leq \frac{\ell_{f,0}}{\mu_g}. \quad (38)$$

861 Observe that

$$863 \quad \begin{aligned} \|\mathbf{d}_t^{\mathbf{x}} - \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\| &\leq \|\nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\| \\ &\quad + \|\mathbf{v}_t \nabla_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{v}_t^*(\mathbf{x}_t) \nabla_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\| \\ &\leq \|\nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\| \\ &\quad + \|\nabla_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)\| \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\| \\ &\quad + \|\mathbf{v}_t^*(\mathbf{x}_t)\| \|\nabla_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\| \\ &\leq \left(\ell_{f,1} + \frac{\ell_{g,2}\ell_{f,0}}{\mu_g} \right) \|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\| + \ell_{g,1} \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\| \\ &\leq M_f^2 (\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\| + \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|), \end{aligned} \quad (39)$$

874 where

$$875 \quad M_f := \sqrt{2} \max \left\{ \ell_{f,1} + \frac{\ell_{g,2}\ell_{f,0}}{\mu_g}, \ell_{g,1} \right\}, \quad (40)$$

876 the third inequality is by Assumption 3.3, and the last inequality follows from (38).

877 Next, we establish (37b).

880 Since $\mathbf{d}_t^{\mathbf{v}^*} := \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) + \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \mathbf{v}_t^*(\mathbf{x}_t) = 0$, we have

$$\begin{aligned} \| \mathbf{d}_t^{\mathbf{v}} \| &= \| \mathbf{d}_t^{\mathbf{v}} - \mathbf{d}_t^{\mathbf{v}^*} \| \\ &= \| \mathbf{v}_t \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) + \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t) \\ &\quad - (\mathbf{v}_t^*(\mathbf{x}_t) \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) + \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \| \\ &\leq \| (\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \mathbf{v}_t^*(\mathbf{x}_t) \| \\ &\quad + \| \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) (\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)) \| \\ &\quad + \| \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \| . \end{aligned}$$

890 Then, from Assumption 3.3 and (38), we have

$$\begin{aligned} \| \mathbf{d}_t^{\mathbf{v}} \| &\leq \ell_{g,2} \| \mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t) \| \| \mathbf{v}_t^*(\mathbf{x}_t) \| + \ell_{g,1} \| \mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t) \| + \ell_{f,1} \| \mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t) \| \\ &\leq \left(\frac{\ell_{g,2} \ell_{f,0}}{\mu_g} + \ell_{f,1} \right) \| \mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t) \| + \ell_{g,1} \| \mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t) \| \\ &\leq M_{\mathbf{v}} (\| \mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t) \| + \| \mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t) \|) , \end{aligned}$$

897 where

$$M_{\mathbf{v}} := \sqrt{2} \max \left\{ \frac{\ell_{g,2} \ell_{f,0}}{\mu_g} + \ell_{f,1}, \ell_{g,1} \right\}. \quad (41)$$

900 The proofs of Eqs. (37c)-(37e) follow from Tarzanagh et al. (2024, Lemma 17) by setting

$$\begin{aligned} L_{\mathbf{y}} &:= \frac{\ell_{g,1}}{\mu_g}, \\ L_{\mathbf{v}} &:= \ell_{f,1} + \frac{\ell_{g,1} \ell_{f,1}}{\mu_g} + \frac{\ell_{f,0}}{\mu_g} \left(\ell_{g,2} + \frac{\ell_{g,1} \ell_{g,2}}{\mu_g} \right), \\ L_f &:= \ell_{f,1} + \frac{\ell_{g,1} (\ell_{f,1} + M_f)}{\mu_g} + \frac{\ell_{f,0}}{\mu_g} \left(\ell_{g,2} + \frac{\ell_{g,1} \ell_{g,2}}{\mu_g} \right), \end{aligned} \quad (42)$$

909 where the other constants are defined in Assumption 3.3. \square

911 **Lemma C.2.** Suppose Assumptions 3.5, B3, and C1. hold. Let $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ be generated according to Algorithm 1.
912 For e_t^g defined as

$$e_t^g := \mathbf{d}_t^{\mathbf{y}} - \nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t), \quad (43)$$

916 we have:

$$\begin{aligned} \mathbb{E} \| e_{t+1}^g \|^2 &\leq (1 - \gamma_{t+1})^2 (1 + 48\ell_{g,1}^2 \beta_t^2) \mathbb{E} \| e_t^g \|^2 + 2\gamma_{t+1}^2 \frac{\sigma_{g_y}^2}{b} + 24(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \mathbb{E} \| \mathbf{x}_{t+1} - \mathbf{x}_t \|^2 \\ &\quad + 6(1 - \gamma_{t+1})^2 \mathbb{E} \| \nabla_{\mathbf{y}} g_t(\mathbf{z}_{t+1}) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_{t+1}) \|^2 \\ &\quad + 48(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \beta_t^2 \mathbb{E} \| \nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t) \|^2. \end{aligned} \quad (44)$$

924 *Proof.* From Algorithm 1, we have

$$\mathbf{d}_{t+1}^{\mathbf{y}} = \nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) + (1 - \gamma_{t+1})(\mathbf{d}_t^{\mathbf{y}} - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})).$$

927 Then, we have

$$\begin{aligned} \mathbb{E} \| e_{t+1}^g \|^2 &= \mathbb{E} \| \mathbf{d}_{t+1}^{\mathbf{y}} - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_{t+1}) \|^2 \\ &= \mathbb{E} \| \nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) + (1 - \gamma_{t+1})(\mathbf{d}_t^{\mathbf{y}} - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_{t+1}) \|^2 \\ &= \mathbb{E} \| (1 - \gamma_{t+1}) e_t^g + (\nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_{t+1})) \\ &\quad - (1 - \gamma_{t+1}) (\nabla_{\mathbf{y}} g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})) - \nabla_{\mathbf{y}} g_t(\mathbf{z}_t) \|^2, \end{aligned}$$

935 which implies that

$$\begin{aligned}
 936 \quad \mathbb{E}\|e_{t+1}^g\|^2 &= (1 - \gamma_{t+1})^2 \mathbb{E}\|e_t^g\|^2 + \mathbb{E}\|(\nabla_y g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_{t+1})) \\
 937 \quad &\quad - (1 - \gamma_{t+1})(\nabla_y g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})) - \nabla_y g_t(\mathbf{z}_t)\|^2 \\
 938 \quad &\leq (1 - \gamma_{t+1})^2 \mathbb{E}\|e_t^g\|^2 + 2\gamma_{t+1}^2 \mathbb{E}\|\nabla_y g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_{t+1})\|^2 \\
 939 \quad &\quad + 2(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) \\
 940 \quad &\quad - \nabla_y g_{t+1}(\mathbf{z}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1}) + \nabla_y g_t(\mathbf{z}_t)\|^2 \\
 941 \quad &\leq (1 - \gamma_{t+1})^2 \mathbb{E}\|e_t^g\|^2 + 2\gamma_{t+1}^2 \frac{\sigma_{g_y}^2}{b} \\
 942 \quad &\quad + 2(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) \\
 943 \quad &\quad - \nabla_y g_{t+1}(\mathbf{z}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1}) + \nabla_y g_t(\mathbf{z}_t)\|^2,
 \end{aligned}$$

944 where the second inequality follows from Cauchy–Schwartz inequality and Assumption 3.5.
 945 Moreover, from Cauchy–Schwartz inequality, we have

$$\begin{aligned}
 946 \quad \mathbb{E}\|e_{t+1}^g\|^2 &\leq (1 - \gamma_{t+1})^2 \mathbb{E}\|e_t^g\|^2 + 2\gamma_{t+1}^2 \frac{\sigma_{g_y}^2}{b} \\
 947 \quad &\quad + 6(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_t(\mathbf{z}_t) - \nabla_y g_t(\mathbf{z}_{t+1})\|^2 \\
 948 \quad &\quad + 6(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_t(\mathbf{z}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_{t+1})\|^2 \\
 949 \quad &\quad + 6(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2.
 \end{aligned}$$

950 From Assumption B3., we have

$$\begin{aligned}
 951 \quad \mathbb{E}\|\nabla_y g_t(\mathbf{z}_{t+1}) - \nabla_y g_t(\mathbf{z}_t)\|^2 \\
 952 \quad &\leq 2\mathbb{E}\|\nabla_y g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_y g_t(\mathbf{x}_{t+1}, \mathbf{y}_t)\|^2 + 2\mathbb{E}\|\nabla_y g_t(\mathbf{x}_{t+1}, \mathbf{y}_t) - \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 953 \quad &\leq 2\ell_{g,1}^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\ell_{g,1}^2 \mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 954 \quad &= 2\ell_{g,1}^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\ell_{g,1}^2 \beta_t^2 \mathbb{E}\|\mathbf{d}_t^y\|^2,
 \end{aligned}$$

955 and

$$\begin{aligned}
 956 \quad \mathbb{E}\|\nabla_y g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 957 \quad &\leq 2\mathbb{E}\|\nabla_y g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_y g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 958 \quad &\quad + 2\mathbb{E}\|\nabla_y g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1}) - \nabla_y g_{t+1}(\mathbf{x}_t, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 959 \quad &\leq 2\ell_{g,1}^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\ell_{g,1}^2 \mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 960 \quad &= 2\ell_{g,1}^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\ell_{g,1}^2 \beta_t^2 \mathbb{E}\|\mathbf{d}_t^y\|^2.
 \end{aligned}$$

961 From the two inequalities above, we have

$$\begin{aligned}
 962 \quad \mathbb{E}\|e_{t+1}^g\|^2 &\leq (1 - \gamma_{t+1})^2 \mathbb{E}\|e_t^g\|^2 + 2\gamma_{t+1}^2 \frac{\sigma_{g_y}^2}{b} \\
 963 \quad &\quad + 6(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_t(\mathbf{z}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_{t+1})\|^2 \\
 964 \quad &\quad + 24(1 - \gamma_{t+1})^2 \ell_{g,1}^2 (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \beta_t^2 \mathbb{E}\|\mathbf{d}_t^y\|^2).
 \end{aligned}$$

990 Since $e_t^g := \mathbf{d}_t^y - \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)$, we have

$$\begin{aligned}
 992 \quad & \mathbb{E}\|e_{t+1}^g\|^2 \leq (1 - \gamma_{t+1})^2 \mathbb{E}\|e_t^g\|^2 + 2\gamma_{t+1}^2 \frac{\sigma_{gy}^2}{\bar{b}} + 24(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\
 993 \quad & + 6(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_t(\mathbf{z}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_{t+1})\|^2 \\
 994 \quad & + 48(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \beta_t^2 \mathbb{E}\|e_t^g\|^2 + 48(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \beta_t^2 \mathbb{E}\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 995 \quad & \leq (1 - \gamma_{t+1})^2 (1 + 48\ell_{g,1}^2 \beta_t^2) \mathbb{E}\|e_t^g\|^2 + 2\gamma_{t+1}^2 \frac{\sigma_{gy}^2}{\bar{b}} + 24(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\
 996 \quad & + 6(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_t(\mathbf{z}_{t+1}) - \nabla_y g_{t+1}(\mathbf{z}_{t+1})\|^2 \\
 997 \quad & + 48(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \beta_t^2 \mathbb{E}\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2.
 \end{aligned}$$

1002 \square

1003 **Lemma C.3.** Suppose Assumptions 3.2, and B3. hold. Then, for the sequence $\{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^T$ generated by Algorithm 1, we
1004 have

$$\begin{aligned}
 1006 \quad & \mathbb{E}[\|\mathbf{y}_{t+1} - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \leq (1 + a) \left(1 - 2\beta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \right) \mathbb{E}[\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \\
 1007 \quad & + \left(-(1 + a) \left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2 \right) \right) \mathbb{E}[\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\
 1008 \quad & + (1 + \frac{1}{a}) \beta_t^2 \mathbb{E}[\|e_t^g\|^2],
 \end{aligned}$$

1013 where e_t^g defined in (43) and $a > 0$ is a constant.

1016 **Proof.** From Lemma B.5, we have

$$\begin{aligned}
 1017 \quad & \mathbb{E}[\|\mathbf{y}_{t+1} - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] = \mathbb{E}[\|\mathbf{y}_t - \beta_t \mathbf{d}_t^y - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \\
 1018 \quad & \leq (1 + a) \mathbb{E}[\|\mathbf{y}_t - \beta_t \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \\
 1019 \quad & + (1 + \frac{1}{a}) \beta_t^2 \mathbb{E}[\|\mathbf{d}_t^y - \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2]. \tag{45}
 \end{aligned}$$

1022 Next, we will bound the first term on the RHS of (45).

1023 We have

$$\begin{aligned}
 1025 \quad & \mathbb{E}[\|\mathbf{y}_t - \beta_t \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] = \mathbb{E}[\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] + \beta_t^2 \mathbb{E}[\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\
 1026 \quad & - 2\beta_t \mathbb{E}[\langle \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t), \mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t) \rangle] \\
 1027 \quad & \leq \left(1 - 2\beta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \right) \mathbb{E}[\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \\
 1028 \quad & - \left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2 \right) \mathbb{E}[\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2], \tag{46}
 \end{aligned}$$

1033 where the inequality results from the strong convexity of g_t by Assumption 3.2, which implies

$$\langle \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t), \mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t) \rangle \geq \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \frac{1}{\mu_g + \ell_{g,1}} \|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2.$$

1037 Substituting (46) into (45), gives the desired result.

1039 \square

1041 To simplify the notation in the analysis, we introduce the definitions

$$\theta_t^y := \|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2, \quad \text{and} \quad \theta_t^v := \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2. \tag{47}$$

1042

1045 **Lemma C.4.** Suppose Assumptions 3.2, and B2., B3. hold. Let θ_t^y be defined as in (47). Then, for the sequence
 1046 $\{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^T$ generated by Algorithm 1, the following bound is guaranteed:

$$\begin{aligned} & \sum_{t=1}^T (\mathbb{E}[\theta_{t+1}^y] - \mathbb{E}[\theta_t^y]) \\ & \leq -\frac{L_{\mu_g}}{2} \sum_{t=1}^T \beta_t \mathbb{E}[\theta_t^y] + \frac{2}{L_{\mu_g}} \sum_{t=1}^T \beta_t \mathbb{E}[\|e_t^g\|^2] + \frac{4L_y^2}{L_{\mu_g}} \sum_{t=1}^T \frac{1}{\beta_t} \mathbb{E}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ & + \frac{4}{L_{\mu_g}} \sum_{t=2}^T \frac{1}{\beta_t} \sup_{\mathbf{x} \in \mathcal{X}} \mathbb{E}\|\mathbf{y}_{t-1}^*(\mathbf{x}) - \mathbf{y}_t^*(\mathbf{x})\|^2 + \sum_{t=1}^T \left(-\frac{2\beta_t}{\mu_g + \ell_{g,1}} + \beta_t^2 \right) \mathbb{E}[\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2], \end{aligned} \quad (48)$$

1057 where $L_{\mu_g} = \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}$, $L_y = \frac{\ell_{g,1}}{\mu_g}$ is defined as in (42); $H_{2,T}$ is defined in (10). Moreover, e_t^g is defined in (43).

1061 *Proof.* From Lemma B.5, we have for any $\hat{c} > 0$

$$\begin{aligned} \mathbb{E}[\|\mathbf{y}_{t+1} - \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})\|^2] &= \mathbb{E}[\|\mathbf{y}_{t+1} - \mathbf{y}_t^*(\mathbf{x}_t) + \mathbf{y}_t^*(\mathbf{x}_t) - \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})\|^2] \\ &\leq (1 + \hat{c}) \mathbb{E}[\|\mathbf{y}_{t+1} - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \\ &+ \left(1 + \frac{1}{\hat{c}}\right) \mathbb{E}[\|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2]. \end{aligned} \quad (49)$$

1068 From Lemma C.3, we have for any $a > 0$

$$\begin{aligned} \mathbb{E}[\|\mathbf{y}_{t+1} - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] &\leq (1 + a) \left(1 - 2\beta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}\right) \mathbb{E}[\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \\ &+ \left(-(1 + a) \left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2\right)\right) \mathbb{E}[\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\ &+ \left(1 + \frac{1}{a}\right) \beta_t^2 \mathbb{E}[\|e_t^g\|^2]. \end{aligned} \quad (50)$$

1077 Substituting (50) into (49), we get

$$\begin{aligned} & \mathbb{E}[\|\mathbf{y}_{t+1} - \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})\|^2] \\ & \leq (1 + \hat{c})(1 + a) \left(1 - 2\beta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}\right) \mathbb{E}[\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \\ & + \left(-(1 + \hat{c})(1 + a) \left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2\right)\right) \mathbb{E}[\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\ & + (1 + \hat{c})(1 + \frac{1}{a}) \beta_t^2 \mathbb{E}[\|e_t^g\|^2] \\ & + \left(1 + \frac{1}{\hat{c}}\right) \mathbb{E}[\|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2]. \end{aligned} \quad (51)$$

1090 Choose $\hat{c} = \frac{\beta_t L_{\mu_g}/2}{1 - \beta_t L_{\mu_g}}$ and $a = \frac{\beta_t L_{\mu_g}}{1 - 2\beta_t L_{\mu_g}}$. Then, the following equations and inequalities are satisfied.

$$\begin{aligned} (1 + \hat{c})(1 + a)(1 - 2\beta_t L_{\mu_g}) &= 1 - \frac{\beta_t L_{\mu_g}}{2}, \\ (1 + a)(1 - 2\beta_t L_{\mu_g}) &= 1 - \beta_t L_{\mu_g}, \\ (1 + \hat{c})(1 - \beta_t L_{\mu_g}) &= 1 - \frac{\beta_t L_{\mu_g}}{2}, \\ 1 + \frac{1}{a} &\leq \frac{1}{\beta_t L_{\mu_g}}, \quad 1 + \frac{1}{\hat{c}} \leq \frac{2}{\beta_t L_{\mu_g}}, \end{aligned} \quad (52)$$

1100 where $L_{\mu_g} = \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}$. Based on (51) and (52), we get
 1101

$$\begin{aligned} & \mathbb{E} [\|\mathbf{y}_{t+1} - \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})\|^2] - \mathbb{E} [\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \\ & \leq -\frac{\beta_t L_{\mu_g}}{2} \mathbb{E} [\|\mathbf{y}_t - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] + \left(-\left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2 \right) \right) \mathbb{E} [\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\ & + \frac{2}{\beta_t L_{\mu_g}} \beta_t^2 \mathbb{E} [\|e_t^g\|^2] + \frac{2}{\beta_t L_{\mu_g}} \mathbb{E} [\|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2]. \end{aligned} \quad (53)$$

1108 Next, we upper-bound the last term of the above inequality.
 1109

$$\begin{aligned} & \mathbb{E} [\|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] \\ & \leq 2 (\mathbb{E} [\|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_{t+1}^*(\mathbf{x}_t)\|^2] + \mathbb{E} [\|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2]) \\ & \leq 2 (L_y^2 \mathbb{E} [\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2]), \end{aligned} \quad (54)$$

1114 where the second inequality is by Lemma D.2.
 1115

1116 Substituting (54) into (53) and summing over $t \in [T]$, give the desired result.
 1117

□

C.3. Bounds on the Linear System Solution

1121 **Lemma C.5.** Suppose Assumptions B2., B3., B4., C2. and C4. hold. Let $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ be generated according to
 1122 Algorithm 1. For e_{t+1}^v defined as

$$1123 e_t^v := \mathbf{d}_t^v - \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t), \quad \text{where } \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t) := \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t + \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t). \quad (55)$$

1125 we have:

$$\begin{aligned} 1126 \mathbb{E} \|e_{t+1}^v\|^2 & \leq (1 - \lambda_{t+1})^2 (1 + 72\ell_{g,1}^2 \delta_t^2) \mathbb{E} \|e_t^v\|^2 + 4\lambda_{t+1}^2 \left(\frac{\sigma_{g_{yy}}^2}{b} p^2 + \frac{\sigma_{f_y}^2}{b} \right) \\ 1127 & + 12p^2 (1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\ 1128 & + 12(1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\ 1129 & + 72(1 - \lambda_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E} \|e_t^g\|^2 + 2\beta_t^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\ 1130 & + 72(1 - \lambda_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \theta_t^v, \end{aligned} \quad (56)$$

1134 for all $t \in [T]$ and θ_t^v is defined in (47).

1137 *Proof.* Note that

$$1139 e_{t+1}^v := \mathbf{d}_{t+1}^v - \nabla P_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}, \mathbf{v}_{t+1}),$$

1141 where

$$1142 \nabla P_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}, \mathbf{v}_{t+1}) := \nabla_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) \mathbf{v}_{t+1} + \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}).$$

1144 From Algorithm 1, we have

$$1146 \mathbf{d}_{t+1}^v = \mathbf{d}_{t+1}^{vv}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1}) + (1 - \lambda_{t+1})(\mathbf{d}_t^v - \mathbf{d}_{t+1}^{vv}(\mathbf{x}_t, \mathbf{y}_t; \mathcal{B}_{t+1})).$$

1148 Let $\mathbf{u} = [\mathbf{x}; \mathbf{y}; \mathbf{v}]$. Then, we have

$$\begin{aligned} 1149 \mathbb{E} \|e_{t+1}^v\|^2 & = \mathbb{E} \|\mathbf{d}_{t+1}^v - \nabla P_{t+1}(\mathbf{u}_{t+1})\|^2 \\ 1150 & = \mathbb{E} \|\nabla P_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) + (1 - \lambda_{t+1})(\mathbf{d}_t^v - \nabla P_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1})) - \nabla P_{t+1}(\mathbf{u}_{t+1})\|^2 \\ 1151 & = \mathbb{E} \|(1 - \lambda_{t+1})e_t^v + \nabla P_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_{t+1}) \\ 1152 & \quad - (1 - \lambda_{t+1})(\nabla P_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1}) - \nabla P_t(\mathbf{u}_t))\|^2, \end{aligned}$$

which implies that

$$\begin{aligned}
 & \mathbb{E}\|e_{t+1}^{\mathbf{v}}\|^2 \\
 &= (1 - \lambda_{t+1})^2 \mathbb{E}\|e_t^{\mathbf{v}}\|^2 + \mathbb{E}\|\lambda_{t+1}(\nabla P_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_{t+1})) \\
 &\quad - (1 - \lambda_{t+1})(\nabla P_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) + \nabla P_{t+1}(\mathbf{u}_{t+1}) - \nabla P_t(\mathbf{u}_t))\|^2 \\
 &\leq (1 - \lambda_{t+1})^2 \mathbb{E}\|e_t^{\mathbf{v}}\|^2 + 2\lambda_{t+1}^2 \mathbb{E}\|\nabla P_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_{t+1})\|^2 \\
 &\quad + 2(1 - \lambda_{t+1})^2 \mathbb{E}\|\nabla P_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1}) + \nabla P_t(\mathbf{u}_t)\|^2,
 \end{aligned}$$

where the inequality follows from Cauchy–Schwartz inequality.

For the first term, from Assumptions C2. and C4., we have

$$\begin{aligned}
 & \mathbb{E}\|\nabla P_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_{t+1})\|^2 \\
 &= \mathbb{E}\|\left(\nabla_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\right) \mathbf{v}_{t+1} \\
 &\quad + \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1}) - \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 &\leq 2\mathbb{E}\|\left(\nabla_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\right) \mathbf{v}_{t+1}\|^2 \\
 &\quad + 2\mathbb{E}\|\nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1}) - \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 &\leq 2\left(\frac{\sigma_{g_{\mathbf{y}\mathbf{y}}}^2}{b} p^2 + \frac{\sigma_{f_{\mathbf{y}}}^2}{b}\right),
 \end{aligned}$$

where the last inequality follows from (8).

Then, from the above inequality and $\|a + b + c\|^2 \leq 3(\|a\|^2 + \|b\|^2 + \|c\|^2)$, we have

$$\begin{aligned}
 \mathbb{E}\|e_{t+1}^{\mathbf{v}}\|^2 &\leq (1 - \lambda_{t+1})^2 \mathbb{E}\|e_t^{\mathbf{v}}\|^2 + 4\lambda_{t+1}^2 \left(\frac{\sigma_{g_{\mathbf{y}\mathbf{y}}}^2}{b} p^2 + \frac{\sigma_{f_{\mathbf{y}}}^2}{b}\right) \\
 &\quad + 6(1 - \lambda_{t+1})^2 \mathbb{E}\|\nabla P_t(\mathbf{u}_t) - \nabla P_t(\mathbf{u}_{t+1})\|^2 \\
 &\quad + 6(1 - \lambda_{t+1})^2 \mathbb{E}\|\nabla P_t(\mathbf{u}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_{t+1})\|^2 \\
 &\quad + 6(1 - \lambda_{t+1})^2 \mathbb{E}\|\nabla P_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1})\|^2. \tag{57}
 \end{aligned}$$

Moreover, from $\|a + b + c\|^2 \leq 3(\|a\|^2 + \|b\|^2 + \|c\|^2)$, we have

$$\begin{aligned}
 & \mathbb{E}\|\nabla P_t(\mathbf{u}_{t+1}) - \nabla P_t(\mathbf{u}_t)\|^2 \\
 &\leq 3\mathbb{E}\|\nabla P_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}, \mathbf{v}_{t+1}) - \nabla P_t(\mathbf{x}_t, \mathbf{y}_{t+1}, \mathbf{v}_{t+1})\|^2 \\
 &\quad + 3\mathbb{E}\|\nabla P_t(\mathbf{x}_t, \mathbf{y}_{t+1}, \mathbf{v}_{t+1}) - \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_{t+1})\|^2 \\
 &\quad + 3\mathbb{E}\|\nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_{t+1}) - \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2 \\
 &\leq 3\mathbb{E}\|(\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_{t+1})) \mathbf{v}_{t+1} + \nabla_{\mathbf{y}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_{t+1})\|^2 \\
 &\quad + 3\mathbb{E}\|(\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)) \mathbf{v}_{t+1} + \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 &\quad + 3\mathbb{E}\|\nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_{t+1}) - \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2 \\
 &\leq 6(\ell_{g,2}^2 \mathbb{E}\|\mathbf{v}_{t+1}\|^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2) + 3\ell_{g,1}^2 \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 \\
 &\leq 6(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \beta_t^2 \mathbb{E}\|\mathbf{d}_t^{\mathbf{y}}\|^2) + 3\ell_{g,1}^2 \delta_t^2 \mathbb{E}\|\mathbf{d}_t^{\mathbf{v}}\|^2 \\
 &\leq 6(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\
 &\quad + 6\ell_{g,1}^2 \delta_t^2 (\mathbb{E}\|e_t^{\mathbf{v}}\|^2 + \mathbb{E}\|\nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2) \\
 &\leq 6(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\
 &\quad + 6\ell_{g,1}^2 \delta_t^2 (\mathbb{E}\|e_t^{\mathbf{v}}\|^2 + \ell_{g,1}^2 \mathbb{E}\|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2), \tag{58}
 \end{aligned}$$

where the third inequality follows from Assumptions B2., B3. and B4.; the last inequality follows from (62).

1210 Similarly, we have

$$\begin{aligned} & \mathbb{E}\|\nabla P_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1})\|^2 \\ & \leq 6(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\ & \quad + 6\ell_{g,1}^2 \delta_t^2 (\mathbb{E}\|e_t^v\|^2 + \ell_{g,1}^2 \mathbb{E}\|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2). \end{aligned} \quad (59)$$

1216 Substituting (59) and (58) into (57), we have

$$\begin{aligned} & \mathbb{E}\|e_{t+1}^v\|^2 \leq (1 - \lambda_{t+1})^2 (1 + 72\ell_{g,1}^2 \delta_t^2) \mathbb{E}\|e_t^v\|^2 + 4\lambda_{t+1}^2 \left(\frac{\sigma_{gyx}^2}{b} p^2 + \frac{\sigma_{fy}^2}{b} \right) \\ & \quad + 6(1 - \lambda_{t+1})^2 \mathbb{E}\|\nabla P_t(\mathbf{u}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_{t+1})\|^2 \\ & \quad + 72(1 - \lambda_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\ & \quad + 72(1 - \lambda_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \mathbb{E}\|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2. \end{aligned}$$

1224 From $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$ and (8), we have

$$\begin{aligned} & \mathbb{E}\|\nabla P_t(\mathbf{u}_{t+1}) - \nabla P_{t+1}(\mathbf{u}_{t+1})\|^2 = \mathbb{E}\|\nabla_y^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) \mathbf{v}_{t+1} - \nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) \mathbf{v}_{t+1} \\ & \quad + \nabla_y f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_y f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\ & \leq 2\mathbb{E}\|(\nabla_y^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})) \mathbf{v}_{t+1}\|^2 \\ & \quad + 2\mathbb{E}\|\nabla_y f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_y f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\ & \leq 2\mathbb{E}\|\nabla_y^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 p^2 \\ & \quad + 2\mathbb{E}\|\nabla_y f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_y f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2. \end{aligned}$$

1235 This completes the proof. \square

1236
1237
1238
1239 As demonstrated in Lemma C.5, the gradient estimation error e_{t+1}^v for the linear system consists of four key components: (1)
1240 an iteratively refined error term $(1 - \lambda_{t+1})^2 (1 + 72\ell_{g,1}^2 \delta_t^2) \mathbb{E}\|e_t^v\|^2$, which depends on the stepsize δ_t ; (2) the error arising
1241 from the variation in the Hessian of the lower-level objective; (3) the error resulting from the variation in the gradient of the
1242 upper-level objective, and (4) an approximation error term of order $\mathcal{O}(\delta_t^2 \theta_t^v)$ associated with solving the linear system.
1243

1244 **Lemma C.6.** Suppose Assumptions 3.2 and 3.3 hold. Then, for the sequence $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ generated by Algorithm 1,
1245 we have

$$\begin{aligned} & \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \leq (1 + \hat{c}) \left(1 - 2\delta_t \frac{(\ell_{g,1} + \ell_{g,1}^3) \mu_g}{\mu_g + \ell_{g,1}} + \delta_t^2 \ell_{g,1}^2 \right) \mathbb{E}\|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \\ & \quad + (1 + \frac{1}{\hat{c}}) \delta_t^2 \mathbb{E}\|e_t^v\|^2, \end{aligned}$$

1252 where e_t^v defined in (55) and for any $\hat{c} > 0$.

1253
1254
1255
1256 *Proof.* From the update rules in Algorithm 1, we have the following:

$$\begin{aligned} & \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 = \mathbb{E}\|\mathbf{v}_t - \delta_t \mathbf{d}_t^v - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \\ & \leq (1 + \hat{c}) \mathbb{E}\|\mathbf{v}_t - \delta_t \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t) - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \\ & \quad + (1 + \frac{1}{\hat{c}}) \delta_t^2 \mathbb{E}\|\mathbf{d}_t^v - \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2, \end{aligned} \quad (60)$$

1263 where $\nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t) := \nabla_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t + \nabla_y f_t(\mathbf{x}_t, \mathbf{y}_t)$.
1264

1265 For the first term of the above eq. (60), we have

$$\begin{aligned}
 & \mathbb{E} \|\mathbf{v}_t - \delta_t \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t) - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \\
 &= \mathbb{E} \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 - 2\delta_t \mathbb{E} \langle \mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t), \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t) \rangle + \delta_t^2 \mathbb{E} \|\nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2 \\
 &\leq \left(1 - 2\delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}\right) \mathbb{E} \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 - (2\delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} - \delta_t^2) \mathbb{E} \|\nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2 \\
 &\leq \left(1 - 2\delta_t \frac{(\ell_{g,1} + \ell_{g,1}^3) \mu_g}{\mu_g + \ell_{g,1}} + \delta_t^2 \ell_{g,1}^2\right) \mathbb{E} \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2,
 \end{aligned} \tag{61}$$

1275 where the first inequality follows from the strong convexity of P_t function (in eq. (4b)) that

$$\mathbb{E} \langle \mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t), \nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t) \rangle \geq \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \mathbb{E} \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 + \frac{1}{\mu_g + \ell_{g,1}} \mathbb{E} \|\nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2.$$

1279 The second inequality is derived from the following inequality.

$$\begin{aligned}
 \mathbb{E} \|\nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2 &= \mathbb{E} \|\nabla_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t + \nabla_y f_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 &= \mathbb{E} \|\nabla_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t)(\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t))\|^2 \leq \ell_{g,1}^2 \mathbb{E} \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2.
 \end{aligned} \tag{62}$$

1284 Combining (60) and (61), we get the desired result. \square

1286 **Lemma C.7.** Suppose Assumptions 3.2 and 3.3 hold. Then, we have

$$\|\mathbf{v}_t^*(\mathbf{x}_t) - \mathbf{v}_{t+1}^*(\mathbf{x}_{t+1})\|^2 \leq 2 \frac{\nu^2}{\mu_g^2} \left(\|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \right),$$

1290 where $\nu := \ell_{f,1} + \frac{\ell_{g,2} \ell_{f,0}}{\mu_g}$, and $\mathbf{v}_t^*(\mathbf{x})$ is a solution of Subproblem (4b).

1293 *Proof.* Based on (4b), we have that

$$\begin{aligned}
 & \|\mathbf{v}_t^*(\mathbf{x}_t) - \mathbf{v}_{t+1}^*(\mathbf{x}_{t+1})\|^2 \\
 &= \|(\nabla_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))^{-1} \nabla_y f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \\
 &\quad - (\nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})))^{-1} \nabla_y f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}))\|^2 \\
 &\leq 2 \left\| \left((\nabla_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))^{-1} - (\nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})))^{-1} \right) \nabla_y f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \right\|^2 \tag{63a}
 \end{aligned}$$

$$\begin{aligned}
 &+ 2 \left\| (\nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})))^{-1} (\nabla_y f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - \nabla_y f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}))) \right\|^2. \tag{63b}
 \end{aligned}$$

1304 In the following steps, we bound the terms (63a) and (63b), respectively.

1305 For (63a), we have:

$$\begin{aligned}
 & \left\| (\nabla_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))^{-1} - (\nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})))^{-1} \right\|^2 \\
 &= \|(\nabla_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))^{-1} (\nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}))) \\
 &\quad - \nabla_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) (\nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})))^{-1}\|^2 \\
 &\leq \frac{1}{\mu_g^2} \|\nabla_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - \nabla_y^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}))\|^2 \\
 &\leq \frac{\ell_{g,2}}{\mu_g^2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\
 &\leq \frac{\ell_{g,2}}{\mu_g^2} \left(\|\mathbf{y}_t^*(\mathbf{x}_t) - \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})\|^2 + \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \right),
 \end{aligned} \tag{64}$$

where the first equality holds since for any invertible matrix \mathbf{A} and \mathbf{B} we have $\|\mathbf{A}^{-1} - \mathbf{B}^{-1}\| = \|\mathbf{A}^{-1}(\mathbf{B} - \mathbf{A})\mathbf{B}^{-1}\|$, and the second inequality is obtained from Assumption 3.3.

Thus, from (64) and Assumption 3.3, we get

$$(63a) \leq \frac{\ell_{f,0}\ell_{g,2}}{\mu_g^2} \left(\|\mathbf{y}_t^*(\mathbf{x}_t) - \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1})\|^2 + \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \right). \quad (65)$$

For (63b), we have

$$\begin{aligned} (63b) &\leq \frac{1}{\mu_g} \|\nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}))\|^2 \\ &\leq \frac{\ell_{f,1}}{\mu_g} \|(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - (\mathbf{x}_{t+1}, \mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}))\|^2 \\ &\leq \frac{\ell_{f,1}}{\mu_g} (\|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2). \end{aligned} \quad (66)$$

Combining (65) and (66), we have

$$\|\mathbf{v}_t^*(\mathbf{x}_t) - \mathbf{v}_{t+1}^*(\mathbf{x}_{t+1})\|^2 \leq \frac{1}{\mu_g} \left(\frac{\ell_{f,0}\ell_{g,2}}{\mu_g} + \ell_{f,1} \right) \left(\|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \right).$$

By raising both sides of the above inequality to the power 2 and using $(a + b)^2 \leq 2a^2 + 2b^2$, we complete the proof. \square

Lemma C.8. Suppose Assumptions 3.2 and 3.3 hold. Let $\theta_t^{\mathbf{v}}$ be defined in (47). Then, for any positive choice of step sizes as

$$\delta_t \leq \frac{L_{\mu_g}}{\ell_{g,1}^2}, \quad \text{where } L_{\mu_g} = \frac{(\ell_{g,1} + \ell_{g,1}^3)\mu_g}{(\mu_g + \ell_{g,1})},$$

for all $t \in [T]$, the sequence $\{\mathbf{v}_t\}_{t=1}^T$ generated by Algorithm 1 satisfy

$$\begin{aligned} \sum_{t=1}^T (\mathbb{E}[\theta_{t+1}^{\mathbf{v}}] - \mathbb{E}[\theta_t^{\mathbf{v}}]) &\leq -\frac{\delta_t L_{\mu_g}}{4} \sum_{t=1}^T \mathbb{E}[\theta_t^{\mathbf{v}}] + \frac{4}{L_{\mu_g}} \delta_t \sum_{t=1}^T \mathbb{E}\|\mathbf{e}_t^{\mathbf{v}}\|^2 \\ &\quad + \frac{16\nu^2}{L_{\mu_g} \mu_g^2 \delta_t} \sum_{t=1}^T \mathbb{E}\|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 \\ &\quad + \frac{8\nu^2}{L_{\mu_g} \mu_g^2 \delta_t} (1 + 2L_y^2) \sum_{t=1}^T \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2, \end{aligned} \quad (67)$$

where $\mathbf{e}_t^{\mathbf{v}}$ is defined in (55).

Proof. By Lemma B.5, for any $a > 0$, we have

$$\begin{aligned} \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_{t+1}^*(\mathbf{x}_{t+1})\|^2 &= \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t^*(\mathbf{x}_t) + \mathbf{v}_t^*(\mathbf{x}_t) - \mathbf{v}_{t+1}^*(\mathbf{x}_{t+1})\|^2 \\ &\leq (1 + a) \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \\ &\quad + \left(1 + \frac{1}{a}\right) \mathbb{E}\|\mathbf{v}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{v}_t^*(\mathbf{x}_t)\|^2. \end{aligned} \quad (68)$$

From Lemma C.6, we have for any $\hat{c} > 0$:

$$\begin{aligned} \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 &\leq (1 + \hat{c}) \left(1 - 2\delta_t \frac{(\ell_{g,1} + \ell_{g,1}^3)\mu_g}{\mu_g + \ell_{g,1}} + \delta_t^2 \ell_{g,1}^2 \right) \mathbb{E}\|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \\ &\quad + (1 + \frac{1}{\hat{c}}) \delta_t^2 \mathbb{E}\|\mathbf{e}_t^{\mathbf{v}}\|^2. \end{aligned} \quad (69)$$

1375 Substituting (69) into (68), we get

$$\begin{aligned}
 1377 \quad \mathbb{E} \|\mathbf{v}_{t+1} - \mathbf{v}_{t+1}^*(\mathbf{x}_{t+1})\|^2 &\leq (1+a)(1+\hat{c}) \left(1 - 2\delta_t \frac{(\ell_{g,1} + \ell_{g,1}^3)\mu_g}{\mu_g + \ell_{g,1}} + \delta_t^2 \ell_{g,1}^2 \right) \mathbb{E} \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \\
 1378 \quad &+ (1+a)(1+\frac{1}{\hat{c}})\delta_t^2 \mathbb{E} \|\mathbf{e}_t^\mathbf{v}\|^2 \\
 1380 \quad &+ \left(1 + \frac{1}{a} \right) \mathbb{E} \|\mathbf{v}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{v}_t^*(\mathbf{x}_t)\|^2. \\
 1381 \quad & \\
 1382 \quad & \\
 1383 \quad &
 \end{aligned} \tag{70}$$

1384 In the following, we provide a bound for the third term on the right-hand side of (70). To this end, we have from Lemma C.7:

$$\begin{aligned}
 1386 \quad \mathbb{E} \|\mathbf{v}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 &\leq 2 \frac{\nu^2}{\mu_g^2} \left(\mathbb{E} \|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \right) \\
 1387 \quad &\leq 2 \frac{\nu^2}{\mu_g^2} \left(2 \mathbb{E} \|\mathbf{y}_{t+1}^*(\mathbf{x}_{t+1}) - \mathbf{y}_{t+1}^*(\mathbf{x}_t)\|^2 \right. \\
 1388 \quad &\quad \left. + 2 \mathbb{E} \|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \right) \\
 1389 \quad &\leq 2 \frac{\nu^2}{\mu_g^2} \left((1+2L_y^2) \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2 \mathbb{E} \|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 \right), \\
 1390 \quad & \\
 1391 \quad & \\
 1392 \quad & \\
 1393 \quad & \\
 1394 \quad & \\
 1395 \quad & \\
 1396 \quad & where the last inequality follows from Lemma C.1. \\
 1397 \quad Combining this result with (70) gives \\
 1398 \quad & \\
 1399 \quad \mathbb{E} \|\mathbf{v}_{t+1} - \mathbf{v}_{t+1}^*(\mathbf{x}_{t+1})\|^2 &\leq (1+a)(1+\hat{c}) \left(1 - 2\delta_t \frac{(\ell_{g,1} + \ell_{g,1}^3)\mu_g}{\mu_g + \ell_{g,1}} + \delta_t^2 \ell_{g,1}^2 \right) \mathbb{E} \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \\
 1400 \quad &+ (1+a)(1+\frac{1}{\hat{c}})\delta_t^2 \mathbb{E} \|\mathbf{e}_t^\mathbf{v}\|^2 + 4 \left(1 + \frac{1}{a} \right) \frac{\nu^2}{\mu_g^2} \mathbb{E} \|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 \\
 1401 \quad &+ 2 \left(1 + \frac{1}{a} \right) \frac{\nu^2}{\mu_g^2} (1+2L_y^2) \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2. \\
 1402 \quad & \\
 1403 \quad & \\
 1404 \quad & \\
 1405 \quad & \\
 1406 \quad & \\
 1407 \quad Let L_{\mu_g} := \frac{(\ell_{g,1} + \ell_{g,1}^3)\mu_g}{\mu_g + \ell_{g,1}}, then we have \\
 1408 \quad & \\
 1409 \quad 1 - 2\delta_t \frac{(\ell_{g,1} + \ell_{g,1}^3)\mu_g}{\mu_g + \ell_{g,1}} + \delta_t^2 \ell_{g,1}^2 &= 1 - 2\delta_t L_{\mu_g} + \delta_t^2 \ell_{g,1}^2 \\
 1410 \quad &\leq 1 - \delta_t L_{\mu_g}, \\
 1411 \quad & \\
 1412 \quad & \\
 1413 \quad where the last inequality follows from \delta_t \leq \frac{L_{\mu_g}}{\ell_{g,1}^2}. \\
 1414 \quad & \\
 1415 \quad Choose a = \frac{\delta_t L_{\mu_g}/4}{1 - \frac{\delta_t L_{\mu_g}}{2}} and \hat{c} = \frac{\delta_t L_{\mu_g}/2}{1 - \delta_t L_{\mu_g}}. Then, from (72), we have \\
 1416 \quad & \\
 1417 \quad & \\
 1418 \quad (1+a)(1+\hat{c}) \left(1 - 2\delta_t \frac{(\ell_{g,1} + \ell_{g,1}^3)\mu_g}{\mu_g + \ell_{g,1}} + \delta_t^2 \ell_{g,1}^2 \right) & \\
 1419 \quad &\leq (1+a)(1+\hat{c})(1 - \delta_t L_{\mu_g}) = 1 - \frac{\delta_t L_{\mu_g}}{4}, \\
 1420 \quad & \\
 1421 \quad (1+a) \left(1 + \frac{1}{\hat{c}} \right) \leq \frac{4}{\delta_t L_{\mu_g}}, \\
 1422 \quad & \\
 1423 \quad 1 + \frac{1}{\hat{c}} \leq \frac{2}{\delta_t L_{\mu_g}}, \quad 1 + \frac{1}{a} \leq \frac{4}{\delta_t L_{\mu_g}}, \\
 1424 \quad & \\
 1425 \quad & \\
 1426 \quad & \\
 1427 \quad & \\
 1428 \quad where L_{\mu_g} := \frac{(\ell_{g,1} + \ell_{g,1}^3)\mu_g}{\mu_g + \ell_{g,1}}. \\
 1429 \quad &
 \end{aligned} \tag{71}$$

1430 Thus, from (71) and (73) we have

$$\begin{aligned} \mathbb{E} \|\mathbf{v}_{t+1} - \mathbf{v}_{t+1}^*(\mathbf{x}_{t+1})\|^2 &\leq \left(1 - \frac{\delta_t L_{\mu_g}}{4}\right) \mathbb{E} \|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2 \\ &\quad + \frac{4}{L_{\mu_g}} \delta_t \mathbb{E} \|e_t^{\mathbf{y}}\|^2 + \frac{16\nu^2}{L_{\mu_g} \mu_g^2 \delta_t} \mathbb{E} \|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 \\ &\quad + \frac{8\nu^2}{L_{\mu_g} \mu_g^2 \delta_t} (1 + 2L_{\mathbf{y}}^2) \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2. \end{aligned}$$

1439 Rearranging the terms and summing from $t = 1$ to T , gives the desired result. \square

C.4. Bounds on the Outer Objective and its Projected Gradient

1442 **Lemma C.9.** Suppose Assumptions B2., B3., C3. and C5. hold. Let $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ be generated according to Algorithm 1. For e_t^f defined as

$$1445 \quad e_t^f := \mathbf{d}_t^{\mathbf{x}} - \tilde{\mathbf{d}}_t(\mathbf{z}_t), \quad \text{where } \tilde{\mathbf{d}}_t(\mathbf{z}_t) = \nabla_{\mathbf{x}} f_t(\mathbf{z}_t) + \nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{z}_t) \mathbf{v}_t, \quad (74)$$

1447 we have:

$$\begin{aligned} 1449 \quad \mathbb{E} \|e_{t+1}^f\|^2 &\leq (1 - \eta_{t+1})^2 \mathbb{E} \|e_t^f\|^2 + 4\eta_{t+1}^2 \left(\frac{\sigma_{g_{\mathbf{xy}}}^2}{\bar{b}} p^2 + \frac{\sigma_{f_{\mathbf{x}}}^2}{b} \right) \\ 1450 &\quad + 12p^2(1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\ 1451 &\quad + 12(1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\ 1452 &\quad + 72(1 - \eta_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E} \|e_t^g\|^2 + 2\beta_t^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\ 1453 &\quad + 72\ell_{g,1}^2 (1 - \eta_{t+1})^2 \delta_t^2 \mathbb{E} \|e_t^{\mathbf{y}}\|^2 + 72(1 - \eta_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \mathbb{E} [\theta_t^{\mathbf{y}}], \end{aligned} \quad (75)$$

1457 for all $t \in [T]$, and $\theta_t^{\mathbf{y}}$ are defined in (47).

1466 *Proof.* Note that

$$1468 \quad e_{t+1}^f = \mathbf{d}_{t+1}^{\mathbf{x}} - \tilde{\mathbf{d}}_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}, \mathbf{v}_{t+1}),$$

1469 where

$$1471 \quad \tilde{\mathbf{d}}_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}, \mathbf{v}_{t+1}) = \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) + \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) \mathbf{v}_{t+1}. \quad (76)$$

1473 From Algorithm 1, we have

$$1474 \quad \mathbf{d}_{t+1}^{\mathbf{x}} = \mathbf{d}_{t+1}^{\mathbf{xx}}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1}) + (1 - \eta_{t+1})(\mathbf{d}_{t+1}^{\mathbf{x}} - \mathbf{d}_{t+1}^{\mathbf{xx}}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1})),$$

1476 where $\mathbf{d}_{t+1}^{\mathbf{xx}}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1}) = \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1}) + \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1}) \mathbf{v}_{t+1}$.

1477 Let $\mathbf{u} = [\mathbf{x}; \mathbf{y}; \mathbf{v}]$. Then, we have

$$\begin{aligned} 1479 \quad \mathbb{E} \|e_{t+1}^f\|^2 &= \mathbb{E} \|\mathbf{d}_{t+1}^{\mathbf{x}} - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1})\|^2 \\ 1480 &= \mathbb{E} \|\tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) + (1 - \eta_{t+1})(\mathbf{d}_t^{\mathbf{x}} - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1})) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1})\|^2 \\ 1481 &= \mathbb{E} \|(1 - \eta_{t+1})e_t^f + \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}) \\ 1482 &\quad - (1 - \eta_{t+1})(\tilde{\mathbf{d}}_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1}) - \tilde{\mathbf{d}}_t(\mathbf{u}_t))\|^2, \end{aligned}$$

which implies that

$$\begin{aligned}
 & \mathbb{E}\|e_{t+1}^f\|^2 \\
 &= (1 - \eta_{t+1})^2 \mathbb{E}\|e_t^f\|^2 + \mathbb{E}\|\eta_{t+1}(\tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1})) \\
 &\quad - (1 - \eta_{t+1})(\tilde{\mathbf{d}}_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) + \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}) - \tilde{\mathbf{d}}_t(\mathbf{u}_t))\|^2 \\
 &\leq (1 - \eta_{t+1})^2 \mathbb{E}\|e_t^f\|^2 + 2\eta_{t+1}^2 \mathbb{E}\|\tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1})\|^2 \\
 &\quad + 2(1 - \eta_{t+1})^2 \mathbb{E}\|\tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1}) + \tilde{\mathbf{d}}_t(\mathbf{u}_t)\|^2,
 \end{aligned} \tag{77}$$

where the inequality follows from $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$.

Let us bound the second term in the right-hand side of (77). Based on (76), we have

$$\begin{aligned}
 & \mathbb{E}\|\tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1})\|^2 \\
 &= \mathbb{E}\|(\nabla_{\mathbf{x}\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})) \mathbf{v}_{t+1} \\
 &\quad + \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1}) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 &\leq 2\mathbb{E}\|(\nabla_{\mathbf{x}\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})) \mathbf{v}_{t+1}\|^2 \\
 &\quad + 2\mathbb{E}\|\nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \mathcal{B}_{t+1}) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 &\leq 2\left(\frac{\sigma_{g_{\mathbf{x}\mathbf{y}}}^2}{b} p^2 + \frac{\sigma_{f_{\mathbf{x}}}^2}{b}\right),
 \end{aligned}$$

where the first inequality is by and $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$; the second inequality follows from Assumptions C3., C5. and (8).

Substituting the above inequality into (77) and using $\|a + b + c\|^2 \leq 3(\|a\|^2 + \|b\|^2 + \|c\|^2)$, we obtain

$$\begin{aligned}
 \mathbb{E}\|e_{t+1}^f\|^2 &\leq (1 - \eta_{t+1})^2 \mathbb{E}\|e_t^f\|^2 + 4\lambda_{t+1}^2 \left(\frac{\sigma_{g_{\mathbf{x}\mathbf{y}}}^2}{b} p^2 + \frac{\sigma_{f_{\mathbf{x}}}^2}{b}\right) \\
 &\quad + 6(1 - \eta_{t+1})^2 \mathbb{E}\|\tilde{\mathbf{d}}_t(\mathbf{u}_t) - \tilde{\mathbf{d}}_t(\mathbf{u}_{t+1})\|^2 \\
 &\quad + 6(1 - \eta_{t+1})^2 \mathbb{E}\|\tilde{\mathbf{d}}_t(\mathbf{u}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1})\|^2 \\
 &\quad + 6(1 - \eta_{t+1})^2 \mathbb{E}\|\tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1})\|^2.
 \end{aligned} \tag{78}$$

Moreover, from $\|a + b + c\|^2 \leq 3(\|a\|^2 + \|b\|^2 + \|c\|^2)$, we have

$$\begin{aligned}
 & \mathbb{E}\|\tilde{\mathbf{d}}_t(\mathbf{u}_{t+1}) - \tilde{\mathbf{d}}_t(\mathbf{u}_t)\|^2 \\
 &\leq 3\mathbb{E}\|\tilde{\mathbf{d}}_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}, \mathbf{v}_{t+1}) - \tilde{\mathbf{d}}_t(\mathbf{x}_t, \mathbf{y}_{t+1}, \mathbf{v}_{t+1})\|^2 \\
 &\quad + 3\mathbb{E}\|\tilde{\mathbf{d}}_t(\mathbf{x}_t, \mathbf{y}_{t+1}, \mathbf{v}_{t+1}) - \tilde{\mathbf{d}}_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_{t+1})\|^2 \\
 &\quad + 3\mathbb{E}\|\tilde{\mathbf{d}}_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_{t+1}) - \tilde{\mathbf{d}}_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2 \\
 &\stackrel{(i)}{\leq} 3\mathbb{E}\|(\nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_{t+1})) \mathbf{v}_{t+1} + \nabla_{\mathbf{x}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_{t+1})\|^2 \\
 &\quad + 3\mathbb{E}\|(\nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)) \mathbf{v}_{t+1} + \nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 &\quad + 3\mathbb{E}\|\tilde{\mathbf{d}}_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_{t+1}) - \tilde{\mathbf{d}}_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2 \\
 &\stackrel{(ii)}{\leq} 6(\ell_{g,2}^2 \mathbb{E}\|\mathbf{v}_{t+1}\|^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2) + 3\ell_{g,1}^2 \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 \\
 &\stackrel{(iii)}{\leq} 6(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \beta_t^2 \mathbb{E}\|\mathbf{d}^y\|^2) + 3\ell_{g,1}^2 \delta_t^2 \mathbb{E}\|\mathbf{d}^y\|^2 \\
 &\stackrel{(iv)}{\leq} 6(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\
 &\quad + 6\ell_{g,1}^2 \delta_t^2 (\mathbb{E}\|e_t^y\|^2 + \mathbb{E}\|\nabla P_t(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2) \\
 &\stackrel{(vi)}{\leq} 6(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\
 &\quad + 6\ell_{g,1}^2 \delta_t^2 (\mathbb{E}\|e_t^y\|^2 + \ell_{g,1}^2 \mathbb{E}\|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2),
 \end{aligned} \tag{79}$$

where the (i) follows from (76); (ii) follows from Assumptions B2., B3. and B4.; (iii) follows from (8); (iv) follows from (43) and (55); (vi) follows from (62).

Similarly, we have

$$\begin{aligned} & \mathbb{E}\|\tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1}; \mathcal{B}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_t; \mathcal{B}_{t+1})\|^2 \\ & \leq 6(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\ & \quad + 6\ell_{g,1}^2 \delta_t^2 (\mathbb{E}\|e_t^v\|^2 + \ell_{g,1}^2 \mathbb{E}\|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2). \end{aligned} \quad (80)$$

Substituting (80) and (79) into (78), we have

$$\begin{aligned} \mathbb{E}\|e_{t+1}^f\|^2 & \leq (1 - \eta_{t+1})^2 \mathbb{E}\|e_t^f\|^2 + 4\eta_{t+1}^2 \left(\frac{\sigma_{g_{\mathbf{y}\mathbf{y}}}^2}{b} p^2 + \frac{\sigma_{f_{\mathbf{y}}}^2}{b} \right) \\ & \quad + 6(1 - \eta_{t+1})^2 \mathbb{E}\|\tilde{\mathbf{d}}_t(\mathbf{u}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1})\|^2 \\ & \quad + 72(1 - \eta_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\ & \quad + 72\ell_{g,1}^2 (1 - \eta_{t+1})^2 \delta_t^2 \mathbb{E}\|e_t^v\|^2 + 72(1 - \eta_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \mathbb{E}\|\mathbf{v}_t - \mathbf{v}_t^*(\mathbf{x}_t)\|^2. \end{aligned}$$

From $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$ and (8), we have

$$\begin{aligned} \mathbb{E}\|\tilde{\mathbf{d}}_t(\mathbf{u}_{t+1}) - \tilde{\mathbf{d}}_{t+1}(\mathbf{u}_{t+1})\|^2 & = \mathbb{E}\|\nabla_{\mathbf{xy}}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) \mathbf{v}_{t+1} - \nabla_{\mathbf{xy}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) \mathbf{v}_{t+1} \\ & \quad + \nabla_{\mathbf{x}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\ & \leq 2\mathbb{E}\|(\nabla_{\mathbf{xy}}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{xy}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})) \mathbf{v}_{t+1}\|^2 \\ & \quad + 2\mathbb{E}\|\nabla_{\mathbf{x}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\ & \leq 2\mathbb{E}\|\nabla_{\mathbf{xy}}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{xy}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 p^2 \\ & \quad + 2\mathbb{E}\|\nabla_{\mathbf{x}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2. \end{aligned}$$

This completes the proof. \square

As demonstrated in Lemma C.9, the hypergradient estimator error e_{t+1}^f comprises five key components: (1) the term $(1 - \eta_{t+1})^2 \mathbb{E}\|e_t^f\|^2$, representing the per-iteration improvement achieved by the momentum-based update; (2) the error arising from the variation in the Jacobian of the lower-level objective; (3) the error caused by the variation in the gradient of the upper-level objective ; (4) the error term $\mathcal{O}(2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2)$, which is due to solving the lower-level problem; and (5) the error term $\mathcal{O}(\delta_t^2 \mathbb{E}\|e_t^v\|^2 + 72(1 - \eta_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \theta_t^v)$, which is introduced by the one-step momentum update in solving the linear system problem.

Lemma C.10. *Let Assumption 3.4 holds. Then, for the sequence of functions $\{f_t\}_{t=1}^T$, we have*

$$\sum_{t=1}^T (f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1}))) \leq 2M + V_T,$$

where M is defined in Assumption 3.4; V_T is defined in (10).

1582

1583 *Proof.* Note that, we have

$$\begin{aligned} & \sum_{t=1}^T (f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1}))) \\ & = f_1(\mathbf{x}_1, \mathbf{y}_1^*(\mathbf{x}_1)) - f_T(\mathbf{x}_{T+1}, \mathbf{y}_T^*(\mathbf{x}_{T+1})) \\ & \quad + \sum_{t=2}^T (f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_{t-1}(\mathbf{x}_t, \mathbf{y}_{t-1}^*(\mathbf{x}_t))) \\ & \leq 2M + V_T, \end{aligned}$$

where the inequality follows from Assumption 3.4. \square

1595 **Lemma C.11.** Let $\{f_t\}_{t=1}^T$ denote the sequence of functions presented to Algorithm 1, satisfying Assumptions 3.2, 3.3 and
 1596 3.4. Let $\mathcal{P}_{\mathcal{X}, \alpha_t}$ be defined as in Definition B.1. For any positive step size α_t such that $\alpha_t \leq \frac{1}{L_f}$ for all $t \in [T]$, Algorithm 1
 1597 ensures the following bound:

$$\begin{aligned} & \sum_{t=1}^T (\alpha_t - L_f \alpha_t^2) \mathbb{E} \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \\ & \leq 8M + 4V_T + 2M_f^2 \sum_{t=1}^T (2\alpha_t - L_f \alpha_t^2) (\mathbb{E}[\theta_t^y] + \mathbb{E}[\theta_t^v]) \\ & \quad + 2 \sum_{t=1}^T (2\alpha_t - L_f \alpha_t^2) \mathbb{E} \|e_t^f\|^2. \end{aligned} \quad (81)$$

1608 Here, θ_t^y and θ_t^v are defined in (47); V_T is defined in (10), M is given in Assumption 3.4; and M_f is defined in (40).

1613 *Proof.* It follows from Lemma C.1 that

$$\begin{aligned} & f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1})) - f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \\ & \leq \langle \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \frac{L_f}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & = -\alpha_t \langle \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)), \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^x) \rangle + \frac{L_f \alpha_t^2}{2} \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^x)\|^2. \end{aligned} \quad (82)$$

1621 For the first term on the right hand side of (82), we have that

$$\begin{aligned} & -\langle \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)), \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^x) \rangle \\ & = -\langle \mathbf{d}_t^x, \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^x) \rangle - \langle \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - \mathbf{d}_t^x, \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^x) \rangle \\ & \leq -\frac{1}{2} \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^x)\|^2 + \frac{1}{2} \|\mathbf{d}_t^x - \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\|^2, \end{aligned}$$

1627 where the inequality follows from Lemma B.8.

1628 Let $\tilde{\mathbf{d}}_t(\mathbf{z}_t) = \nabla_{\mathbf{x}} f_t(\mathbf{z}_t) + \nabla_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{z}_t) \mathbf{v}_t$. Then, from Lemma C.1, we have

$$\begin{aligned} \|\mathbf{d}_t^x - \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\|^2 &= \left\| \mathbf{d}_t^x - \tilde{\mathbf{d}}_t(\mathbf{z}_t) + \tilde{\mathbf{d}}_t(\mathbf{z}_t) - \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \right\|^2 \\ &\leq 2 \|\mathbf{d}_t^x - \tilde{\mathbf{d}}_t(\mathbf{z}_t)\|^2 + 2 \|\tilde{\mathbf{d}}_t(\mathbf{z}_t) - \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\|^2 \\ &\leq 2 \|e_t^f\|^2 + 2 \|\tilde{\mathbf{d}}_t(\mathbf{z}_t) - \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\|^2 \\ &\leq 2 \|e_t^f\|^2 + M_f^2 (\theta_t^y + \theta_t^v), \end{aligned} \quad (83)$$

1639 where $e_t^f := \mathbf{d}_t^x - \tilde{\mathbf{d}}_t(\mathbf{z}_t)$. This implies that

$$\begin{aligned} & -\langle \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)), \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^x) \rangle \\ & \leq -\frac{1}{2} \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^x)\|^2 + 2 \|e_t^f\|^2 + M_f^2 (\theta_t^y + \theta_t^v), \end{aligned} \quad (84)$$

1644 Plugging the bound (84) into (82), we have that

$$\begin{aligned} & f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1})) - f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \\ & \leq \frac{(L_f \alpha_t^2 - \alpha_t)}{2} \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^x)\|^2 + 2\alpha_t \|e_t^f\|^2 + M_f^2 (\theta_t^y + \theta_t^v) \alpha_t, \end{aligned}$$

which can be rearranged into

$$\begin{aligned} & (\alpha_t - L_f \alpha_t^2) \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^{\mathbf{x}})\|^2 \\ & \leq 2f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1})) + 4\alpha_t \|e_t^f\|^2 + 2M_f^2 (\theta_t^{\mathbf{y}} + \theta_t^{\mathbf{v}}) \alpha_t. \end{aligned} \quad (85)$$

In addition, we have

$$\begin{aligned} & \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \\ & \leq 2 \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^{\mathbf{x}}) - \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 + 2 \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^{\mathbf{x}})\|^2 \\ & \leq 2 \|\mathbf{d}_t^{\mathbf{x}} - \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\|^2 + 2 \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^{\mathbf{x}})\|^2 \\ & \leq 4 \|e_t^f\|^2 + 4M_f^2 (\theta_t^{\mathbf{y}} + \theta_t^{\mathbf{v}}) + 4 \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^{\mathbf{x}})\|^2, \end{aligned} \quad (86)$$

where the second inequality follows from non-expansiveness of the projection operator and the last inequality follows from (83).

Combining (85) and (86), we have

$$\begin{aligned} & \sum_{t=1}^T (\alpha_t - L_f \alpha_t^2) \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \\ & \leq 4 \sum_{t=1}^T (f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1}))) \\ & \quad + 2M_f^2 \sum_{t=1}^T (2\alpha_t - L_f \alpha_t^2) (\theta_t^{\mathbf{y}} + \theta_t^{\mathbf{v}}) + 2 \sum_{t=1}^T (2\alpha_t - L_f \alpha_t^2) \|e_t^f\|^2 \\ & \leq 8M + 4V_T \\ & \quad + 2M_f^2 \sum_{t=1}^T (2\alpha_t - L_f \alpha_t^2) (\theta_t^{\mathbf{y}} + \theta_t^{\mathbf{v}}) + 2 \sum_{t=1}^T (2\alpha_t - L_f \alpha_t^2) \|e_t^f\|^2, \end{aligned}$$

where the second inequality is due to Lemma C.10. \square

Lemma C.12. Let Assumptions 3.3 and 3.4 hold. Let $\{\mathbf{x}_t\}_{t=1}^T$ be generated according to Algorithm 1. Then, we have

$$\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \leq 2\alpha_t^2 \left(\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 + M_f^2 (\theta_t^{\mathbf{y}} + \theta_t^{\mathbf{v}}) \right),$$

where $\theta_t^{\mathbf{y}}$ and $\theta_t^{\mathbf{v}}$ are defined in (47).

Proof. From the update rule of Algorithm 1, we have

$$\begin{aligned} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 &= \alpha_t^2 \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^{\mathbf{x}})\|^2 \\ &\leq 2\alpha_t^2 \left(\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right. \\ &\quad \left. + \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \mathbf{d}_t^{\mathbf{x}}) - \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right) \\ &\leq 2\alpha_t^2 \left(\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right. \\ &\quad \left. + \|\mathbf{d}_t^{\mathbf{x}} - \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))\|^2 \right) \\ &\leq 2\alpha_t^2 \left(\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 + M_f^2 (\theta_t^{\mathbf{y}} + \theta_t^{\mathbf{v}}) \right), \end{aligned} \quad (87)$$

where the first inequality is by $(a+b)^2 \leq 2a^2 + 2b^2$; the second inequality follows from non-expansiveness of the projection operator; and the last inequality follows from Eq. (37a) in Lemma C.1. \square

C.5. Proof of Theorem 3.6

Proof. **Bounding $\mathbb{E}\|e_t^f\|^2$ in (75).** From (75), we have

$$\begin{aligned}
 & \frac{\mathbb{E}\|e_{t+1}^f\|^2}{\alpha_t} - \frac{\mathbb{E}\|e_t^f\|^2}{\alpha_{t-1}} \leq \left(\frac{(1-\eta_{t+1})^2}{\alpha_t} - \frac{1}{\alpha_{t-1}} \right) \mathbb{E}\|e_t^f\|^2 + \frac{4\eta_{t+1}^2}{\alpha_t} \left(\frac{\sigma_{g_{xy}}^2}{b} p^2 + \frac{\sigma_{f_x}^2}{b} \right) \\
 & + \frac{12p^2}{\alpha_t} (1-\eta_{t+1})^2 \mathbb{E}\|\nabla_{xy}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{xy}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 & + \frac{12}{\alpha_t} (1-\eta_{t+1})^2 \mathbb{E}\|\nabla_{\mathbf{x}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 & + \frac{72}{\alpha_t} (1-\eta_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\
 & + \frac{72}{\alpha_t} \ell_{g,1}^2 (1-\eta_{t+1})^2 \delta_t^2 \mathbb{E}\|e_t^v\|^2 + \frac{72}{\alpha_t} (1-\eta_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \mathbb{E}[\theta_t^v].
 \end{aligned} \tag{88}$$

With respect to the coefficient of the first term on the right-hand side of equation (88), it is important to note that we have:

$$\frac{(1-\eta_{t+1})^2}{\alpha_t} - \frac{1}{\alpha_{t-1}} \leq \frac{1}{\alpha_t} - \frac{\eta_{t+1}}{\alpha_t} - \frac{1}{\alpha_{t-1}}. \tag{89}$$

Using the definition of α_t in (15), we have

$$\begin{aligned}
 \frac{1}{\alpha_t} - \frac{1}{\alpha_{t-1}} &= (c+t)^{1/3} - (c+t-1)^{1/3} \stackrel{(i)}{\leq} \frac{1}{3(c+t-1)^{2/3}} \stackrel{(ii)}{\leq} \frac{1}{3(\frac{c}{2}+t)^{2/3}} \\
 &= \frac{2^{2/3}}{3(c+2t)^{2/3}} \stackrel{(iii)}{\leq} \frac{2^{2/3}}{3(c+t)^{2/3}} \stackrel{(iv)}{\leq} \frac{2^{2/3}}{3} \alpha_t^2 \stackrel{(vi)}{\leq} \frac{\alpha_t}{6L_f},
 \end{aligned} \tag{90}$$

where the (i) follows from $(a+b)^{1/3} - a^{1/3} \leq b/(3a^{2/3})$; (ii) follows from $c \geq 2$ in (104); (iii) follows from (15); (iv) follows from $\alpha_t \leq 1/4L_f$ in (104).

Substituting (90) into (89) and using $\delta_t = c_\delta \alpha_t$ and $\eta_{t+1} = c_\eta \alpha_t^2$, we have

$$\frac{(1-\eta_{t+1})^2}{\alpha_t} - \frac{1}{\alpha_{t-1}} \leq \frac{\alpha_t}{6L_f} - \frac{\eta_{t+1}}{\alpha_t} = \frac{\alpha_t}{6L_f} - c_\eta \alpha_t \leq -5\Omega \alpha_t, \tag{91}$$

where the inequalities follow from $c_\eta = \frac{1}{6L_f} + 5\Omega$ with Ω in (103).

Then, substituting (91) into (88) yields

$$\begin{aligned}
 & \frac{1}{\Omega} \mathbb{E} \left(\frac{\|e_{t+1}^f\|^2}{\alpha_t} - \frac{\|e_t^f\|^2}{\alpha_{t-1}} \right) \leq -5\alpha_t \mathbb{E}\|e_t^f\|^2 + \frac{4\eta_{t+1}^2}{\Omega \alpha_t} \left(\frac{\sigma_{g_{xy}}^2}{b} p^2 + \frac{\sigma_{f_x}^2}{b} \right) \\
 & + \frac{12p^2}{\Omega \alpha_t} (1-\eta_{t+1})^2 \mathbb{E}\|\nabla_{xy}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{xy}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 & + \frac{12}{\Omega \alpha_t} (1-\eta_{t+1})^2 \mathbb{E}\|\nabla_{\mathbf{x}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 & + \frac{72}{\Omega \alpha_t} (1-\eta_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\
 & + \frac{72}{\Omega \alpha_t} \ell_{g,1}^2 (1-\eta_{t+1})^2 \delta_t^2 \mathbb{E}\|e_t^v\|^2 + \frac{72}{\Omega \alpha_t} (1-\eta_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \mathbb{E}[\theta_t^v].
 \end{aligned} \tag{92}$$

Bounding $\mathbb{E}\|e_t^g\|^2$ in (44).

From (44), we have

$$\begin{aligned}
 \frac{\mathbb{E}\|e_{t+1}^g\|^2}{\alpha_t} - \frac{\mathbb{E}\|e_t^g\|^2}{\alpha_{t-1}} &\leq \left(\frac{1}{\alpha_t} (1 - \gamma_{t+1})^2 (1 + 48\ell_{g,1}^2 \beta_t^2) - \frac{1}{\alpha_{t-1}} \right) \mathbb{E}\|e_t^g\|^2 \\
 &\quad + 2 \frac{\gamma_{t+1}^2 \sigma_{gy}^2}{\alpha_t} + \frac{24}{\alpha_t} (1 - \gamma_{t+1})^2 \ell_{g,1}^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\
 &\quad + \frac{6}{\alpha_t} (1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_y g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 &\quad + 48(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \frac{\beta_t^2}{\alpha_t} \mathbb{E}\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2. \tag{93}
 \end{aligned}$$

Let us examine the coefficient of the first term on the right-hand side of Eq. (93). Specifically, for $\gamma_{t+1} = c_\gamma \alpha_t^2$ and $\beta_t = c_\beta \alpha_t$, we have:

$$\begin{aligned}
 \frac{1}{\alpha_t} (1 - \gamma_{t+1})^2 (1 + 48\ell_{g,1}^2 \beta_t^2) - \frac{1}{\alpha_{t-1}} &\leq \frac{1}{\alpha_t} (1 - \gamma_{t+1}) (1 + 48\ell_{g,1}^2 \beta_t^2) - \frac{1}{\alpha_{t-1}} \\
 &= \frac{1}{\alpha_t} - \frac{1}{\alpha_{t-1}} - \frac{\gamma_{t+1}}{\alpha_t} + \frac{1 - \gamma_{t+1}}{\alpha_t} 48\ell_{g,1}^2 \beta_t^2 \\
 &= \frac{1}{\alpha_t} - \frac{1}{\alpha_{t-1}} - c_\gamma \alpha_t + \left(\frac{1}{\alpha_t} - c_\gamma \alpha_t \right) 48\ell_{g,1}^2 c_\beta^2 \alpha_t^2 \\
 &\leq \frac{\alpha_t}{6L_f} + 48\ell_{g,1}^2 c_\beta^2 \alpha_t - c_\gamma \alpha_t, \tag{94}
 \end{aligned}$$

where the last inequality follows from (90).

Recalling Φ from (103) that we selected, we obtain

$$c_\gamma = \frac{1}{6L_f} + 48\ell_{g,1}^2 c_\beta^2 + \hbar \Phi, \quad \text{where } \hbar := 25 \frac{M_f^2}{L_{\mu_g}^2},$$

which, when combined with Eq. (94), results in

$$\frac{1}{\alpha_t} (1 - \gamma_{t+1})^2 (1 + 48\ell_{g,1}^2 \beta_t^2) - \frac{1}{\alpha_{t-1}} \leq -\hbar \Phi \alpha_t. \tag{95}$$

Substituting eq. (95) into eq. (93) yields

$$\begin{aligned}
 \frac{1}{\Phi} \left(\frac{\mathbb{E}\|e_{t+1}^g\|^2}{\alpha_t} - \frac{\mathbb{E}\|e_t^g\|^2}{\alpha_{t-1}} \right) &\leq -\hbar \alpha_t \mathbb{E}\|e_t^g\|^2 \\
 &\quad + 2 \frac{\gamma_{t+1}^2 \sigma_{gy}^2}{\Phi \alpha_t} + \frac{24}{\Phi \alpha_t} (1 - \gamma_{t+1})^2 \ell_{g,1}^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\
 &\quad + \frac{6}{\Phi \alpha_t} (1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_y g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 &\quad + 48(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \frac{\beta_t^2}{\Phi \alpha_t} \mathbb{E}\|\nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2. \tag{96}
 \end{aligned}$$

Bounding $\mathbb{E}\|e_t^v\|^2$ in (56).

1815 From (56), we get

$$\begin{aligned}
 & \frac{\mathbb{E}\|e_{t+1}^{\mathbf{v}}\|^2 - \mathbb{E}\|e_t^{\mathbf{v}}\|^2}{\alpha_t} \leq \left(\frac{1}{\alpha_t}(1-\lambda_{t+1})^2(1+72\ell_{g,1}^2\delta_t^2) - \frac{1}{\alpha_{t-1}} \right) \mathbb{E}\|e_t^{\mathbf{v}}\|^2 \\
 & + 4\frac{\lambda_{t+1}^2}{\alpha_t} \left(\frac{\sigma_{g_{\mathbf{y}\mathbf{y}}}^2 p^2}{b} + \frac{\sigma_{f_{\mathbf{y}}}^2}{b} \right) + \frac{12p^2}{\alpha_t} (1-\lambda_{t+1})^2 \mathbb{E}\|\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 & + \frac{12}{\alpha_t} (1-\lambda_{t+1})^2 \mathbb{E}\|\nabla_{\mathbf{y}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 & + \frac{72}{\alpha_t} (1-\lambda_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\
 & + \frac{72}{\alpha_t} (1-\lambda_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \mathbb{E}[\theta_t^{\mathbf{v}}].
 \end{aligned} \tag{97}$$

1828 Let us examine the coefficient of the first term on the right-hand side of equation (97). Specifically, for $\lambda_{t+1} = c_{\lambda} \alpha_t^2$ and
 1829 $\delta_t = c_{\delta} \alpha_t$, we have:

$$\begin{aligned}
 & \frac{1}{\alpha_t} (1-\lambda_{t+1})^2 (1+72\ell_{g,1}^2\delta_t^2) - \frac{1}{\alpha_{t-1}} \leq \frac{1}{\alpha_t} (1-\lambda_{t+1})(1+72\ell_{g,1}^2\delta_t^2) - \frac{1}{\alpha_{t-1}} \\
 & = \frac{1}{\alpha_t} - \frac{1}{\alpha_{t-1}} - \frac{\lambda_{t+1}}{\alpha_t} + \frac{1-\lambda_{t+1}}{\alpha_t} 72\ell_{g,1}^2\delta_t^2 \\
 & = \frac{1}{\alpha_t} - \frac{1}{\alpha_{t-1}} - c_{\lambda} \alpha_t + \left(\frac{1}{\alpha_t} - c_{\lambda} \alpha_t \right) 72\ell_{g,1}^2 c_{\delta}^2 \alpha_t^2 \\
 & \leq \frac{\alpha_t}{6L_f} + 72\ell_{g,1}^2 c_{\delta}^2 \alpha_t - c_{\lambda} \alpha_t,
 \end{aligned} \tag{98}$$

1840 where the last inequality follows from (90).

1841 Recalling Ψ from (103) that we selected, we obtain

$$c_{\lambda} = \frac{1}{6L_f} + 72\ell_{g,1}^2 c_{\delta}^2 + \jmath\Psi, \quad \text{where } \jmath = 90 \frac{M_f^2}{L_{\mu_g}^2},$$

1846 which, when combined with Eq. (98), results in

$$\frac{1}{\alpha_t} (1-\lambda_{t+1})^2 (1+72\ell_{g,1}^2\delta_t^2) - \frac{1}{\alpha_{t-1}} \leq -\jmath\Psi \alpha_t. \tag{99}$$

1850 Substituting eq. (99) into eq. (97) yields

$$\begin{aligned}
 & \frac{1}{\Psi} \left(\frac{\mathbb{E}\|e_{t+1}^{\mathbf{v}}\|^2}{\alpha_t} - \frac{\mathbb{E}\|e_t^{\mathbf{v}}\|^2}{\alpha_{t-1}} \right) \leq -\jmath\Psi \mathbb{E}\|e_t^{\mathbf{v}}\|^2 \\
 & + 4\frac{\lambda_{t+1}^2}{\Psi\alpha_t} \left(\frac{\sigma_{g_{\mathbf{y}\mathbf{y}}}^2 p^2}{b} + \frac{\sigma_{f_{\mathbf{y}}}^2}{b} \right) + \frac{12p^2}{\Psi\alpha_t} (1-\lambda_{t+1})^2 \mathbb{E}\|\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 & + \frac{12}{\Psi\alpha_t} (1-\lambda_{t+1})^2 \mathbb{E}\|\nabla_{\mathbf{y}} f_t(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})\|^2 \\
 & + \frac{72}{\Psi\alpha_t} (1-\lambda_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) (\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2\beta_t^2 \mathbb{E}\|e_t^g\|^2 + 2\beta_t^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2) \\
 & + \frac{72}{\Psi\alpha_t} (1-\lambda_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \mathbb{E}[\theta_t^{\mathbf{v}}].
 \end{aligned} \tag{100}$$

1863 **Combining the outcomes.** We recall from Lemma C.12 that we have

$$\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \leq 2\alpha_t^2 \left(\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 + M_f^2 (\theta_t^{\mathbf{y}} + \theta_t^{\mathbf{v}}) \right). \tag{101}$$

1870 Let

$$\begin{aligned} \Lambda := & \Gamma \sum_{t=1}^T (\mathbb{E}[\theta_{t+1}^{\mathbf{y}}] - \mathbb{E}[\theta_t^{\mathbf{y}}]) + \Upsilon \sum_{t=1}^T (\mathbb{E}[\theta_{t+1}^{\mathbf{x}}] - \mathbb{E}[\theta_t^{\mathbf{x}}]) + \frac{1}{\Phi} \sum_{t=1}^T \left(\frac{\mathbb{E}\|e_{t+1}^g\|^2}{\alpha_t} - \frac{\mathbb{E}\|e_t^g\|^2}{\alpha_{t-1}} \right) \\ & + \frac{1}{\Psi} \sum_{t=1}^T \left(\frac{\mathbb{E}\|e_{t+1}^v\|^2}{\alpha_t} - \frac{\mathbb{E}\|e_t^v\|^2}{\alpha_{t-1}} \right) + \frac{1}{\Omega} \sum_{t=1}^T \left(\frac{\mathbb{E}\|e_{t+1}^f\|^2}{\alpha_t} - \frac{\mathbb{E}\|e_t^f\|^2}{\alpha_{t-1}} \right). \end{aligned} \quad (102)$$

1877 Here

$$\begin{aligned} \Gamma &= \frac{11M_f^2}{L_{\mu_g} c_\beta}, \quad \Upsilon = \frac{22M_f^2}{L_{\mu_g} c_\delta}, \quad \Phi = 480\ell_{g,1}^2, \\ \Psi &= \max \left\{ 144(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \left(10 + \frac{L_{\mu_g}^2 c_\beta^2}{11M_f^2} \right), \frac{288\ell_{g,1}^4}{M_f^2} c_\delta^2 \right\}, \\ \Omega &= \max \left\{ 144(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \left(10 + \frac{L_{\mu_g}^2 c_\beta^2}{11M_f^2} \right), \frac{288\ell_{g,1}^4}{M_f^2} c_\delta^2 \right\}. \end{aligned} \quad (103)$$

1888 Here, from (15), we have

$$\begin{aligned} c &\geq \max \{4L_f, c_\beta(\mu_g + \ell_{g,1}), 2\}, \\ c_\beta &= \sqrt{880 \frac{L_y^2 M_f^2}{L_{\mu_g}^2}}, \\ c_\delta &= \sqrt{3520 \frac{\nu^2 M_f^2}{L_{\mu_g}^2 \mu_g^2} (1 + 2L_y^2)}, \\ c_\gamma &= \frac{2}{3L_f} + 48\ell_{g,1}^2 c_\beta^2 + \hbar\Phi, \quad \text{where } \hbar := 25 \frac{M_f^2}{L_{\mu_g}^2}, \\ c_\eta &= \frac{2}{3L_f} + 5\Omega, \\ c_\lambda &= \frac{2}{3L_f} + 72\ell_{g,1}^2 c_\delta^2 + \jmath\Psi, \quad \text{where } \jmath = 90 \frac{M_f^2}{L_{\mu_g}^2}. \end{aligned} \quad (104)$$

1904 Using (100), (96), (92), (81), (67), and (48), along with (101) and the fact that α_t decreases with respect to t , we obtain:

$$\begin{aligned} & \sum_{t=1}^T A(\alpha_t, \beta_t, \delta_t) \mathbb{E} \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 + \Lambda \\ & \leq 8M + 4V_T + \sum_{t=1}^T B(\alpha_t, \beta_t, \delta_t) \mathbb{E}[\theta_t^v] + \sum_{t=1}^T C(\alpha_t, \beta_t) \mathbb{E}[\theta_t^y] \end{aligned} \quad (105a)$$

$$+ \sum_{t=1}^T D(\alpha_t) \mathbb{E}\|e_t^f\|^2 + \sum_{t=1}^T F(\beta_t, \delta_t) \mathbb{E}\|e_t^g\|^2 + \sum_{t=1}^T I(\alpha_t, \beta_t, \delta_t) \mathbb{E}\|e_t^v\|^2 \quad (105b)$$

$$+ \sum_{t=1}^T L(\beta_t) \mathbb{E}\|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 + \sum_{t=2}^T N(\beta_t, \delta_t) \sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{y}_{t-1}^*(\mathbf{x}) - \mathbf{y}_t^*(\mathbf{x})\|^2 \quad (105c)$$

$$+ \frac{\sigma_{gy}^2}{b} \frac{2}{\Phi} \sum_{t=1}^T \frac{\gamma_{t+1}^2}{\alpha_t} + \frac{4}{\Psi} \left(\frac{\sigma_{gxy}^2}{b} p^2 + \frac{\sigma_{fx}^2}{b} \right) \sum_{t=1}^T \frac{\lambda_{t+1}^2}{\alpha_t} + \frac{4}{\Omega} \left(\frac{\sigma_{gxy}^2}{b} p^2 + \frac{\sigma_{fx}^2}{b} \right) \sum_{t=1}^T \frac{\eta_{t+1}^2}{\alpha_t} \quad (105d)$$

$$+ \frac{6}{\Phi\alpha_T} G_{y,T} + \frac{12p^2}{\Omega\alpha_T} G_{xy,T} + \frac{12p^2}{\Psi\alpha_T} G_{yy,T} + \frac{12\ell_{f,1}^2}{\Psi\alpha_T} D_{y,T} + \frac{12\ell_{f,1}^2}{\Omega\alpha_T} D_{x,T}. \quad (105e)$$

1923 Here, M is defined in Assumption 3.4, V_T and $H_{2,T}$ are defined in (10). Moreover, $G_{y,T}$, $G_{xy,T}$, and $G_{yy,T}$ are defined in
1924

1925 (12). Let

$$\begin{aligned}
 1927 \quad E(\beta_t, \delta_t) &:= \frac{4L_y^2}{L_{\mu_g}\beta_t} \Gamma + \frac{8\nu^2}{L_{\mu_g}\mu_g^2\delta_t} (1+2L_y^2)\Upsilon + 72(1-\eta_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{1}{\Omega\alpha_t} \\
 1928 \quad &\quad + 24(1-\gamma_{t+1})^2\ell_{g,1}^2 \frac{1}{\Phi\alpha_t} + 72(1-\lambda_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{1}{\Psi\alpha_t}, \\
 1929 \quad A(\alpha_t, \beta_t, \delta_t) &:= \alpha_t - (L_f + 2E(\beta_t, \delta_t))\alpha_t^2, \\
 1930 \quad B(\alpha_t, \beta_t, \delta_t) &:= -\frac{L_{\mu_g}\Upsilon}{4}\delta_t + 4M_f^2\alpha_t - 2M_f^2L_f\alpha_t^2 + 2M_f^2E(\beta_t, \delta_t)\alpha_t^2 \\
 1931 \quad &\quad + 72(1-\lambda_{t+1})^2\ell_{g,1}^4\delta_t^2 \frac{1}{\Psi\alpha_t} + 72(1-\eta_{t+1})^2\ell_{g,1}^4\delta_t^2 \frac{1}{\Omega\alpha_t}, \\
 1932 \quad C(\alpha_t, \beta_t) &:= -\frac{L_{\mu_g}\Gamma}{2}\beta_t + 4M_f^2\alpha_t - 2L_fM_f^2\alpha_t^2 + 2M_f^2E(\beta_t, \delta_t)\alpha_t^2, \\
 1933 \quad D(\alpha_t) &:= 2(2\alpha_t - L_f\alpha_t^2) - 5\alpha_t, \\
 1934 \quad F(\alpha_t, \beta_t, \delta_t) &:= \frac{2\Gamma}{L_{\mu_g}}\beta_t - \hbar\alpha_t + 72(1-\lambda_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2)2\frac{\beta_t^2}{\Psi\alpha_t} \\
 1935 \quad &\quad + 72(1-\eta_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2)2\frac{\beta_t^2}{\Omega\alpha_t}, \\
 1936 \quad I(\alpha_t, \beta_t, \delta_t) &:= \frac{4\Upsilon}{L_{\mu_g}}\delta_t - \jmath\alpha_t + 72\ell_{g,1}^2(1-\eta_{t+1})^2\frac{\delta_t^2}{\Omega\alpha_t}.
 \end{aligned} \tag{106}$$

1948 Moreover, we have

$$\begin{aligned}
 1949 \quad L(\beta_t) &:= -\frac{2\Gamma}{\mu_g + \ell_{g,1}}\beta_t + \Gamma\beta_t^2 + 48(1-\gamma_{t+1})^2\ell_{g,1}^2\frac{\beta_t^2}{\Phi\alpha_t} \\
 1950 \quad &\quad + 72(1-\lambda_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2)2\frac{\beta_t^2}{\Psi\alpha_t} + 72(1-\eta_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2)2\frac{\beta_t^2}{\Omega\alpha_t}, \\
 1951 \quad N(\beta_t, \delta_t) &:= \frac{4}{L_{\mu_g}\beta_t}\Gamma + \frac{16\nu^2}{L_{\mu_g}\mu_g^2\delta_t}\Upsilon.
 \end{aligned} \tag{107}$$

1957 Note that, we have

$$\begin{aligned}
 1958 \quad E(\beta_t, \delta_t) &= \frac{4L_y^2}{L_{\mu_g}\beta_t}\Gamma + \frac{8\nu^2}{L_{\mu_g}\mu_g^2\delta_t}(1+2L_y^2)\Upsilon + 72(1-\eta_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{1}{\Omega\alpha_t} \\
 1959 \quad &\quad + 24(1-\gamma_{t+1})^2\ell_{g,1}^2 \frac{1}{\Phi\alpha_t} + 72(1-\lambda_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{1}{\Psi\alpha_t},
 \end{aligned}$$

1964 which together with $\beta_t = c_\beta\alpha_t$, $\delta_t = c_\delta\alpha_t$, we have

$$\begin{aligned}
 1965 \quad \alpha_t^2 E(\beta_t, \delta_t) &= \frac{4L_y^2}{L_{\mu_g}}\Gamma\frac{\alpha_t^2}{\beta_t} + \frac{8\nu^2}{L_{\mu_g}\mu_g^2}(1+2L_y^2)\Upsilon\frac{\alpha_t^2}{\delta_t} + 72(1-\eta_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{\alpha_t}{\Omega} \\
 1966 \quad &\quad + 24(1-\gamma_{t+1})^2\ell_{g,1}^2 \frac{\alpha_t}{\Phi} + 72(1-\lambda_{t+1})^2(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{\alpha_t}{\Psi} \\
 1967 \quad &\leq \frac{44L_y^2}{L_{\mu_g}^2}M_f^2\frac{\alpha_t}{c_\beta^2} + \frac{176\nu^2}{L_{\mu_g}^2\mu_g^2}(1+2L_y^2)M_f^2\frac{\alpha_t}{c_\delta^2} \\
 1968 \quad &\quad + 24\ell_{g,1}^2 \frac{\alpha_t}{\Phi} + 72(\ell_{g,2}^2 p^2 + \ell_{f,1}^2)(\frac{1}{\Omega} + \frac{1}{\Psi})\alpha_t \\
 1969 \quad &\leq \frac{\alpha_t}{4},
 \end{aligned} \tag{108}$$

1977 where the first inequality follows from $\Gamma = \frac{11M_f^2}{L_{\mu_g}c_\beta}$ and $\Upsilon = \frac{22M_f^2}{L_{\mu_g}c_\delta}$ in (103); the last inequality follows from $c_\beta \geq$
 1978
 1979

1980 $\sqrt{880 \frac{L_y^2 M_f^2}{L_{\mu_g}^2}}, c_\delta \geq \sqrt{3520 \frac{\nu^2 M_f^2}{L_{\mu_g}^2 \mu_g^2} (1 + 2L_y^2)}$, in (104) and $\Phi \geq 480 \ell_{g,1}^2$, and $\Omega, \Psi \geq 1440(\ell_{g,2}^2 p^2 + \ell_{f,1}^2)$ in (103).
 1981 Moreover, we have
 1982

$$\begin{aligned}
 1983 \quad A(\alpha_t, \beta_t, \delta_t) &:= \alpha_t - L_f \alpha_t^2 - 2E(\beta_t, \delta_t) \alpha_t^2 \\
 1984 \quad &\geq \alpha_t - L_f \alpha_t^2 - \frac{\alpha_t}{2} \\
 1985 \quad &\geq \frac{\alpha_t}{4}, \\
 1986 \\
 1987 \\
 1988
 \end{aligned} \tag{109}$$

where the last inequality is by $\alpha_t \leq 1/4L_f$ in (104).

Bounding (105a).

From (106), we have

$$\begin{aligned}
 1992 \quad B(\alpha_t, \beta_t, \delta_t) &= -\frac{L_{\mu_g} \Upsilon}{4} \delta_t + 4M_f^2 \alpha_t - 2M_f^2 L_f \alpha_t^2 + 2M_f^2 E(\beta_t, \delta_t) \alpha_t^2 \\
 1993 \quad &+ 72(1 - \lambda_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \frac{1}{\Psi \alpha_t} + 72(1 - \eta_{t+1})^2 \ell_{g,1}^4 \delta_t^2 \frac{1}{\Omega \alpha_t} \\
 1994 \quad &\leq -\frac{L_{\mu_g} \Upsilon}{4} \delta_t + 4M_f^2 \alpha_t - 2M_f^2 L_f \alpha_t^2 + \frac{M_f^2}{2} \alpha_t + 72 \ell_{g,1}^4 \left(\frac{1}{\Psi} + \frac{1}{\Omega} \right) \frac{\delta_t^2}{\alpha_t} \\
 1995 \quad &= \left(-\frac{L_{\mu_g} \Upsilon}{4} c_\delta + \frac{9}{2} M_f^2 + 72 \ell_{g,1}^4 \left(\frac{1}{\Psi} + \frac{1}{\Omega} \right) c_\delta^2 \right) \alpha_t \\
 1996 \quad &\leq -\frac{1}{2} M_f^2 \alpha_t, \\
 1997 \\
 1998 \\
 1999 \\
 2000 \\
 2001 \\
 2002 \\
 2003
 \end{aligned} \tag{110}$$

where the first inequality follows from $\beta_t = c_\beta \alpha_t$, $\delta_t = c_\delta \alpha_t$, and (108); the second inequality is by $\Upsilon = \frac{22M_f^2}{L_{\mu_g} c_\delta}$, and $\Psi, \Omega \geq \frac{288\ell_{g,1}^4}{M_f^2} c_\delta^2$ in (103); the last inequality follows from in (103).

Moreover, from (106), and $\beta_t = c_\beta \alpha_t$, we have

$$\begin{aligned}
 2004 \quad C(\alpha_t, \beta_t) &= -\frac{L_{\mu_g} \Gamma}{2} \beta_t + 4M_f^2 \alpha_t - 2L_f M_f^2 \alpha_t^2 + 2M_f^2 E(\beta_t, \delta_t) \alpha_t^2 \\
 2005 \quad &\leq -\frac{L_{\mu_g} \Gamma}{2} c_\beta \alpha_t + \frac{9}{2} M_f^2 \alpha_t \\
 2006 \quad &= -M_f^2 \alpha_t, \\
 2007 \\
 2008 \\
 2009 \\
 2010 \\
 2011 \\
 2012 \\
 2013 \\
 2014
 \end{aligned} \tag{111}$$

where the first inequality follows from (108); the last equality follows from $\Gamma = \frac{11M_f^2}{L_{\mu_g} c_\beta}$ in (103).

Thus, from (110) and (111), we get

$$(105a) \leq \mathcal{O}(V_T). \tag{112}$$

Bounding (105b).

From (106), we also obtain

$$D(\alpha_t) = 4\alpha_t - 2L_f \alpha_t^2 - 5\alpha_t \leq 0.$$

From $\beta_t = c_\beta \alpha_t$, $\Gamma = \frac{11M_f^2}{L_{\mu_g} c_\beta}$ and (106), we obtain

$$\begin{aligned}
 F(\alpha_t, \beta_t, \delta_t) &= \frac{2\Gamma}{L_{\mu_g}} \beta_t - \hbar \alpha_t + 144(1 - \lambda_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{\beta_t^2}{\Psi \alpha_t} \\
 &\quad + 144(1 - \eta_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{\beta_t^2}{\Omega \alpha_t} \\
 &\leq \frac{22M_f^2}{L_{\mu_g}^2} \alpha_t - \hbar \alpha_t + 144(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \left(\frac{1}{\Psi} + \frac{1}{\Omega} \right) c_\beta^2 \alpha_t \\
 &\leq 24 \frac{M_f^2}{L_{\mu_g}^2} \alpha_t - \hbar \alpha_t \\
 &= -\frac{M_f^2}{L_{\mu_g}^2} \alpha_t,
 \end{aligned}$$

where the second inequality follows from $\Omega, \Psi \geq 144(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{L_{\mu_g}^2 c_\beta^2}{M_f^2}$ in (103); and the last equality is by $\hbar := 25 \frac{M_f^2}{L_{\mu_g}^2}$.

From $\delta_t = c_\delta \alpha_t$, we obtain

$$\begin{aligned}
 I(\alpha_t, \beta_t, \delta_t) &= \frac{4\Upsilon}{L_{\mu_g}} \delta_t - \jmath \alpha_t + 72\ell_{g,1}^2 (1 - \eta_{t+1})^2 \frac{\delta_t^2}{\Omega \alpha_t} \\
 &\leq \frac{4\Upsilon}{L_{\mu_g}} c_\delta \alpha_t - \jmath \alpha_t + 72\ell_{g,1}^2 \frac{c_\delta^2 \alpha_t}{\Omega} \\
 &\leq \frac{89M_f^2}{L_{\mu_g}^2} \alpha_t - \jmath \alpha_t \\
 &= -\frac{M_f^2}{L_{\mu_g}^2} \alpha_t,
 \end{aligned}$$

where the second inequality follows from $\Upsilon = \frac{22M_f^2}{L_{\mu_g} c_\delta}$ and $\Omega \geq \frac{72\ell_{g,1}^2 L_{\mu_g}^2}{M_f^2} c_\delta^2$; the last equality follows from $\jmath = 90 \frac{M_f^2}{L_{\mu_g}^2}$.

Thus, we get

$$(105b) \leq 0. \quad (113)$$

Bounding (105c).

From $\beta_t = c_\beta \alpha_t$ and (107), we have

$$\begin{aligned}
 L(\beta_t) &= -\frac{2\Gamma \beta_t}{\mu_g + \ell_{g,1}} + \Gamma \beta_t^2 + 48(1 - \gamma_{t+1})^2 \ell_{g,1}^2 \frac{\beta_t^2}{\Phi \alpha_t} \\
 &\quad + 72(1 - \lambda_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) 2 \frac{\beta_t^2}{\Psi \alpha_t} + 72(1 - \eta_{t+1})^2 (\ell_{g,2}^2 p^2 + \ell_{f,1}^2) 2 \frac{\beta_t^2}{\Omega \alpha_t} \\
 &\leq -\frac{2\Gamma c_\beta \alpha_t}{\mu_g + \ell_{g,1}} + \Gamma c_\beta^2 \alpha_t^2 + 48\ell_{g,1}^2 c_\beta^2 \frac{\alpha_t}{\Phi} + 144(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \left(\frac{1}{\Psi} + \frac{1}{\Omega} \right) c_\beta^2 \alpha_t \\
 &\leq -\frac{2\Gamma c_\beta \alpha_t}{\mu_g + \ell_{g,1}} + \Gamma c_\beta^2 \alpha_t^2 + \frac{3\Gamma c_\beta \alpha_t}{4(\mu_g + \ell_{g,1})} \\
 &\leq -\frac{\Gamma c_\beta \alpha_t}{4(\mu_g + \ell_{g,1})},
 \end{aligned}$$

where the second inequality is by $\Phi \geq 192\ell_{g,1}^2 \frac{(\mu_g + \ell_{g,1})}{\Gamma} c_\beta$, and $\Omega, \Psi \geq 576(\ell_{g,2}^2 p^2 + \ell_{f,1}^2) \frac{(\mu_g + \ell_{g,1})}{\Gamma} c_\beta$ in (103); the last inequality follows from $\alpha_t \leq \frac{1}{c_\beta(\mu_g + \ell_{g,1})}$ in (104).

From $\beta_t = c_\beta \alpha_t$, $\delta_t = c_\delta \alpha_t$ and (107), we obtain

$$\begin{aligned} N(\beta_t, \delta_t) &= \frac{4}{L_{\mu_g} \beta_t} \Gamma + \frac{16\nu^2}{L_{\mu_g} \mu_g^2 \delta_t} \Upsilon \\ &= \frac{4}{L_{\mu_g} c_\beta \alpha_t} \Gamma + \frac{16\nu^2}{L_{\mu_g} \mu_g^2 c_\delta \alpha_t} \Upsilon. \end{aligned}$$

Thus, we get

$$\begin{aligned} (105c) &= \sum_{t=1}^T L(\beta_t) \mathbb{E} \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 + \sum_{t=2}^T N(\beta_t, \delta_t) \sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{y}_{t-1}^*(\mathbf{x}) - \mathbf{y}_t^*(\mathbf{x})\|^2 \\ &\leq \mathcal{O}\left(\frac{H_{2,T}}{\alpha_T}\right). \end{aligned} \tag{114}$$

Bounding (105d).

From $\eta_{t+1} = c_\eta \alpha_t^2$, $\gamma_{t+1} = c_\gamma \alpha_t^2$, $\lambda_{t+1} = c_\lambda \alpha_t^2$, we obtain

$$\begin{aligned} (105d) &= \frac{\sigma_{g_y}^2}{b} \frac{2}{\Phi} \sum_{t=1}^T \frac{\gamma_{t+1}^2}{\alpha_t} + \frac{4}{\Psi} \left(\frac{\sigma_{g_{yy}}^2}{b} p^2 + \frac{\sigma_{f_y}^2}{b} \right) \sum_{t=1}^T \frac{\lambda_{t+1}^2}{\alpha_t} + \frac{4}{\Omega} \left(\frac{\sigma_{g_{xy}}^2}{b} p^2 + \frac{\sigma_{f_x}^2}{b} \right) \sum_{t=1}^T \frac{\eta_{t+1}^2}{\alpha_t} \\ &\leq \mathcal{O}\left((\frac{\sigma_{g_y}^2}{b} + \frac{\sigma_{g_{yy}}^2}{b} + \frac{\sigma_{f_y}^2}{b} + \frac{\sigma_{g_{xy}}^2}{b} + \frac{\sigma_{f_x}^2}{b}) \sum_{t=1}^T \alpha_t^3\right). \end{aligned} \tag{115}$$

Bounding (105e).

We have

$$\begin{aligned} (105e) &= \frac{6}{\Phi \alpha_T} G_{y,T} + \frac{12p^2}{\Omega \alpha_T} G_{xy,T} + \frac{12p^2}{\Psi \alpha_T} G_{yy,T} + \frac{12\ell_{f,1}^2}{\Psi \alpha_T} D_{y,T} + \frac{12\ell_{f,1}^2}{\Omega \alpha_T} D_{x,T} \\ &\leq \mathcal{O}\left(\frac{1}{\alpha_T} (G_{y,T} + G_{xy,T} + G_{yy,T} + D_{y,T} + D_{x,T})\right). \end{aligned} \tag{116}$$

Let

$$\begin{aligned} G_T &:= G_{y,T} + G_{xy,T} + G_{yy,T}, \\ D_T &:= D_{y,T} + D_{x,T}, \\ \sigma^2 &:= \sigma_{g_y}^2 + \sigma_{g_{yy}}^2 + \sigma_{f_y}^2 + \sigma_{g_{xy}}^2 + \sigma_{f_x}^2, \\ b &= \bar{b} = 1. \end{aligned}$$

By inequalities (109), (112), (113), (114), (115), (116), we have

$$\begin{aligned} &\sum_{t=1}^T \frac{\alpha_t}{2} \mathbb{E} \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 + \Lambda \\ &\leq \mathcal{O}\left(V_T + \frac{H_{2,T}}{\alpha_T} + \frac{\sigma^2}{b} \sum_{t=1}^T \alpha_t^3 + \frac{G_T}{\alpha_T} + \frac{D_T}{\alpha_T}\right). \end{aligned} \tag{117}$$

2145 From the definition of Λ in (102), we have

$$\begin{aligned}
 -\Lambda &= \Gamma \sum_{t=1}^T (\mathbb{E}[\theta_t^y] - \mathbb{E}[\theta_{t+1}^y]) + \Upsilon \sum_{t=1}^T (\mathbb{E}[\theta_t^v] - \mathbb{E}[\theta_{t+1}^v]) + \frac{1}{\Phi} \sum_{t=1}^T \left(\frac{\mathbb{E}\|e_t^g\|^2}{\alpha_{t-1}} - \frac{\mathbb{E}\|e_{t+1}^g\|^2}{\alpha_t} \right) \\
 &\quad + \frac{1}{\Psi} \sum_{t=1}^T \left(\frac{\mathbb{E}\|e_t^v\|^2}{\alpha_{t-1}} - \frac{\mathbb{E}\|e_{t+1}^v\|^2}{\alpha_t} \right) + \frac{1}{\Omega} \sum_{t=1}^T \left(\frac{\mathbb{E}\|e_t^f\|^2}{\alpha_{t-1}} - \frac{\mathbb{E}\|e_{t+1}^f\|^2}{\alpha_t} \right) \\
 &\leq \Gamma \theta_1^y + \Upsilon \theta_1^v + \frac{\sigma_{gy}^2}{\Phi \alpha_0} + \frac{\sigma_{gyy}^2 + \sigma_{fy}^2}{\Psi \alpha_0} + \frac{\sigma_{gxy}^2 + \sigma_{fx}^2}{\Omega \alpha_0}.
 \end{aligned} \tag{118}$$

2155 Using (118), we get

$$\begin{aligned}
 &\sum_{t=1}^T \frac{\alpha_t}{2} \mathbb{E} \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \\
 &\leq \mathcal{O} \left(V_T + \frac{H_{2,T}}{\alpha_T} + \frac{\sigma^2}{b} \sum_{t=1}^T \alpha_t^3 + \frac{G_T}{\alpha_T} + \frac{D_T}{\alpha_T} - \Lambda \right) \\
 &\leq \mathcal{O} \left(V_T + \theta_1^y + \theta_1^v + \frac{\sigma^2}{b} \sum_{t=1}^T \alpha_t^3 + \frac{H_{2,T}}{\alpha_T} + \frac{G_T}{\alpha_T} + \frac{D_T}{\alpha_T} + \frac{\sigma^2}{\alpha_0} \right).
 \end{aligned}$$

2166 Since $\alpha_t = 1/(c+t)^{1/3}$, we get

$$\sum_{t=1}^T \alpha_t^3 = \sum_{t=1}^T \frac{1}{c+t} \leq \sum_{t=1}^T \frac{1}{1+t} \leq \log(T+1),$$

2171 which, combined with the fact that α_t decreases with respect to t and by multiplying both sides by $2/\alpha_T$, results in Thus,
2172 we have

$$\begin{aligned}
 \text{BL-Reg}_T &= \sum_{t=1}^T \mathbb{E} \|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \\
 &\leq \mathcal{O} \left(\frac{1}{\alpha_T} (V_T + \|\mathbf{y}_1 - \mathbf{y}_1^*(\mathbf{x}_1)\|^2 + \|\mathbf{v}_1 - \mathbf{v}_1^*(\mathbf{x}_1)\|^2 + \sigma^2 \log(T+1) + \frac{\sigma^2}{\alpha_0}) \right. \\
 &\quad \left. + \frac{1}{\alpha_T^2} (H_{2,T} + G_T + D_T) \right).
 \end{aligned}$$

2181 This completes the proof. \square

D. Proof of Regret Bounds for Zeroth Order SOGD (ZO-SOGD)

2185 **Proof Roadmap.** We provide Lemma D.9, which quantifies the error between the approximated direction of the momentum-based gradient estimator, $\hat{\mathbf{d}}_t^y$ and the true direction, $\nabla_y g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)$, at each iteration. Lemma D.11 assesses the convergence of the iterative solutions $\{\mathbf{y}_t\}_{t=1}^T$, specifically the gap $\mathbb{E}[\|\mathbf{y}_{t+1} - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2]$, while accounting for the error introduced in Lemma D.9. To establish Lemma D.15, which quantifies the error between the approximated direction of the momentum-based gradient estimator, $\hat{\mathbf{d}}_t^y$, and the true direction, $\nabla_y f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) + \nabla_y^2 g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t$, we need to present Lemma D.13. This lemma quantifies the error between $\hat{\mathbf{d}}_t^y$ and $\nabla_y f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) + \frac{1}{2\rho_v} (\nabla_y g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_y g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t))$. Then, Lemma D.17 captures the error of the system solution of Problem (17), i.e., gap $\mathbb{E}[\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2]$, based on these errors. To establish Lemma D.21, which quantifies the error between the approximated direction of the momentum-based hypergradient estimator, $\hat{\mathbf{d}}_t^x$, and the true direction, $\nabla_x f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) + \nabla_x^2 g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t$, it is necessary to introduce Lemma D.19. This lemma quantifies the error between $\hat{\mathbf{d}}_t^x$ and $\nabla_x f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) + \frac{1}{2\rho_v} (\nabla_x g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_x g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t))$. Then, Lemma D.22 bounds the projection mapping based on these errors. By combining these lemmas and setting parameters, we achieve the desired result.

D.1. Auxiliary Lemmas for Proof of Theorem 4.2

To solve the online bilevel problem in (17), we need to obtain the hyper-gradient of $f_{t,\rho}$ in (17) at (\mathbf{x}, \mathbf{y}) as

$$\begin{aligned} \nabla f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) &:= \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) + \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) \mathbf{v}_t^*(\mathbf{x}), \quad \text{where} \\ \mathbf{v}_t^*(\mathbf{x}) &:= -[\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))]^{-1} \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})). \end{aligned}$$

Obtaining $\hat{\mathbf{y}}_t^*(\mathbf{x})$ in closed-form is usually a challenging task, so it is natural to use the following gradient surrogate. At any (\mathbf{x}, \mathbf{y}) , define:

$$\tilde{\nabla} f_{t,\rho}(\mathbf{x}, \mathbf{y}) := \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) + \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) \hat{\mathbf{v}}_t^*(\mathbf{x}), \quad \text{where} \quad (119a)$$

$$\hat{\mathbf{v}}_t^*(\mathbf{x}) := -[\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y})]^{-1} \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}). \quad (119b)$$

To do so, we also introduce $\mathbf{d}_{t,\rho}^y$, $\mathbf{d}_{t,\rho}^v$ and $\mathbf{d}_{t,\rho}^x$ as follows:

$$\mathbf{d}_{t,\rho}^y(\mathbf{x}, \mathbf{y}) = \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}, \mathbf{y}), \quad (120a)$$

$$\mathbf{d}_{t,\rho}^v(\mathbf{x}, \mathbf{y}, \mathbf{v}) = \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) + \nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) \mathbf{v}, \quad (120b)$$

$$\mathbf{d}_{t,\rho}^x(\mathbf{x}, \mathbf{y}, \mathbf{v}) = \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) + \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) \mathbf{v}. \quad (120c)$$

To approximate these directions, we use (19)-(21). It can be shown that $\hat{\nabla}_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y}; \xi)$ and $\hat{\nabla}_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{y}; \xi)$ are unbiased estimators of the true gradients $\nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y})$ and $\nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \mathbf{y})$ with respect to \mathbf{y} and \mathbf{x} (Flaxman et al., 2004), respectively, i.e.,

$$\mathbb{E}_{\mathbf{r}} [\hat{\nabla}_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y}; \xi)] = \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}), \quad \mathbb{E}_{\mathbf{z}} [\hat{\nabla}_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{y}; \xi)] = \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \mathbf{y}),$$

and,

$$\mathbb{E}_{\mathbf{r}} [\hat{\nabla}_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}; \zeta)] = \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}, \mathbf{y}), \quad \mathbb{E}_{\mathbf{z}} [\hat{\nabla}_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}; \zeta)] = \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}, \mathbf{y}). \quad (121)$$

Similarly,

$$\mathbb{E}_{\mathbf{r}} [\hat{\nabla}_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y}; \mathcal{B})] = \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}), \quad \mathbb{E}_{\mathbf{z}} [\hat{\nabla}_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{y}; \mathcal{B})] = \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \mathbf{y}),$$

and,

$$\mathbb{E}_{\mathbf{r}} [\hat{\nabla}_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}})] = \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}, \mathbf{y}), \quad \mathbb{E}_{\mathbf{z}} [\hat{\nabla}_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}})] = \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}, \mathbf{y}). \quad (122)$$

Lemma D.1. (Allen-Zhu & Li (2018, Lemma A.1.) Suppose Assumption B4. holds. Then, for any $\mathbf{x}, \mathbf{v} \in \mathcal{X}$, we have:

$$\|\nabla g_t(\mathbf{x} + \mathbf{v}, \mathbf{y} + \mathbf{v}) - \nabla g_t(\mathbf{x}, \mathbf{y}) - \nabla^2 g_t(\mathbf{x}, \mathbf{y}) \mathbf{v}\| \leq \ell_{g,2} \|\mathbf{v}\|^2.$$

Lemma D.2. Suppose that Assumptions 3.2 and 3.3 hold for all $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, and $t \in [T]$, and that $\mathbf{d}_{t,\rho}^x$ and $\mathbf{d}_{t,\rho}^y$ are defined in (120). Then, we have

$$\|\mathbf{d}_{t,\rho}^x - \nabla f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))\|^2 \leq M_f^2 (\|\mathbf{y} - \hat{\mathbf{y}}_t^*(\mathbf{x})\|^2 + \|\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x})\|^2), \quad (123a)$$

$$\|\mathbf{d}_{t,\rho}^y\|^2 \leq M_v^2 (\|\mathbf{y} - \hat{\mathbf{y}}_t^*(\mathbf{x})\|^2 + \|\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x})\|^2), \quad (123b)$$

$$\|\nabla f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) - \nabla f_{t,\rho}(\mathbf{x}', \hat{\mathbf{y}}_t^*(\mathbf{x}'))\| \leq L_f \|\mathbf{x} - \mathbf{x}'\|, \quad (123c)$$

$$\|\hat{\mathbf{y}}_t^*(\mathbf{x}) - \hat{\mathbf{y}}_t^*(\mathbf{x}')\| \leq L_y \|\mathbf{x} - \mathbf{x}'\|, \quad (123d)$$

$$\|\hat{\mathbf{v}}_t^*(\mathbf{x}) - \hat{\mathbf{v}}_t^*(\mathbf{x}')\| \leq L_v \|\mathbf{x} - \mathbf{x}'\|. \quad (123e)$$

Here, $\hat{\mathbf{v}}_t^*(\mathbf{x})$ and $f_{t,\rho}, \hat{\mathbf{y}}_t^*(\mathbf{x})$ are defined in (119b) and (17), respectively. Moreover, the constants M_f , M_v , and (L_y, L_v, L_f) are defined as in (40), (41), and (42), respectively.

Proof. We first show Eq. (123a).

Using Assumptions 3.2 and B1., we have $\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) \succeq \mu_g$, and

$$\|\hat{\mathbf{v}}_t^*(\mathbf{x})\| = \|(\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})))^{-1} \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))\| \leq \frac{\ell_{f,0}}{\mu_g}. \quad (124)$$

Observe that we have

$$\begin{aligned} \|\mathbf{d}_{t,\rho}^{\mathbf{x}} - \nabla f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))\| &\leq \|\nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))\| \\ &\quad + \|\mathbf{v} \nabla_{\mathbf{xy}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \hat{\mathbf{v}}_t^*(\mathbf{x}) \nabla_{\mathbf{xy}}^2 g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))\| \\ &\leq \|\nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))\| \\ &\quad + \|\nabla_{\mathbf{xy}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y})\| \|\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x})\| \\ &\quad + \|\hat{\mathbf{v}}_t^*(\mathbf{x})\| \|\nabla_{\mathbf{xy}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{xy}}^2 g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))\| \\ &\leq \left(\ell_{f,1} + \frac{\ell_{g,2}\ell_{f,0}}{\mu_g} \right) \|\mathbf{y} - \hat{\mathbf{y}}_t^*(\mathbf{x})\| + \ell_{g,1} \|\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x})\| \\ &\leq M_f^2 (\|\mathbf{y} - \hat{\mathbf{y}}_t^*(\mathbf{x})\| + \|\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x})\|), \end{aligned} \quad (125)$$

where M_f is defined as in (40); the third inequality is by Assumption 3.3 and the last inequality is by Eq. (124).

We now show Eq. (123b).

Since $\mathbf{d}_{t,\rho}^{\mathbf{y}^*} := \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) + \nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) \hat{\mathbf{v}}_t^*(\mathbf{x}) = 0$, we have

$$\begin{aligned} \|\mathbf{d}_{t,\rho}^{\mathbf{y}}\| &= \|\mathbf{d}_{t,\rho}^{\mathbf{y}} - \mathbf{d}_{t,\rho}^{\mathbf{y}^*}\| \\ &= \|\mathbf{v}_t \nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) + \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) \\ &\quad - (\hat{\mathbf{v}}_t^*(\mathbf{x}) \nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) + \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})))\| \\ &\leq \|(\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))) \hat{\mathbf{v}}_t^*(\mathbf{x})\| \\ &\quad + \|\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) (\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x}))\| \\ &\quad + \|\nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))\|. \end{aligned}$$

Then, from Assumption 3.3 and Eq. (124), we have

$$\begin{aligned} \|\mathbf{d}_{t,\rho}^{\mathbf{y}}\| &\leq \ell_{g,2} \|\mathbf{y} - \hat{\mathbf{y}}_t^*(\mathbf{x})\| \|\hat{\mathbf{v}}_t^*(\mathbf{x})\| + \ell_{g,1} \|\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x})\| + \ell_{f,1} \|\mathbf{y} - \hat{\mathbf{y}}_t^*(\mathbf{x})\| \\ &\leq \left(\frac{\ell_{g,2}\ell_{f,0}}{\mu_g} + \ell_{f,1} \right) \|\mathbf{y} - \hat{\mathbf{y}}_t^*(\mathbf{x})\| + \ell_{g,1} \|\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x})\| \\ &\leq M_{\mathbf{v}} (\|\mathbf{y} - \hat{\mathbf{y}}_t^*(\mathbf{x})\| + \|\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x})\|), \end{aligned}$$

where $M_{\mathbf{v}}$ is defined as in (41).

The proofs of Eqs. (123c)-(123e) follow from Tarzanagh et al. (2024, Lemma 17) by setting $(L_{\mathbf{y}}, L_{\mathbf{v}}, L_f)$ as in (42). \square

D.2. Perturbation Bounds for OBO Objectives and Their Smoothing Variants

Lemma D.3. Given $\rho = (\rho_s, \rho_r)$ as positive smoothing parameters, let $g_{t,\rho}(\mathbf{x}, \mathbf{y})$ and $f_{t,\rho}(\mathbf{x}, \mathbf{y})$ be the functions defined by (17).

(a) Suppose Assumption B3. holds. Then, we have

$$|g_{t,\rho}(\mathbf{x}, \mathbf{y}) - g_t(\mathbf{x}, \mathbf{y})| \leq \frac{\ell_{g,1}(\rho_s^2 + \rho_r^2)}{2}. \quad (126)$$

(b) Suppose Assumption B2. holds. Then, we have

$$|f_{t,\rho}(\mathbf{x}, \mathbf{y}) - f_t(\mathbf{x}, \mathbf{y})| \leq \frac{\ell_{f,1}(\rho_s^2 + \rho_r^2)}{2}. \quad (127)$$

Proof. Let B_1 and B_2 be the unit ball in \mathbb{R}^{d_1} and \mathbb{R}^{d_2} , respectively. Let $\mathcal{V}(d_1)$ and $\mathcal{V}(d_2)$ be volume of the unit ball in \mathbb{R}^{d_1}

2310 and \mathbb{R}^{d_2} , respectively. Then, we have

$$\begin{aligned} & |g_{t,\rho}(\mathbf{x}, \mathbf{y}) - g_t(\mathbf{x}, \mathbf{y})| \\ &= \left| \frac{1}{\mathcal{V}(d_1)\mathcal{V}(d_2)} \int_{B_1} \int_{B_2} (g_t(\mathbf{x} + \rho_s \mathbf{s}, \mathbf{y} + \rho_r \mathbf{r}) - g_t(\mathbf{x}, \mathbf{y})) d\mathbf{s} d\mathbf{r} \right| \\ &= \left| \frac{1}{\mathcal{V}(d_1)\mathcal{V}(d_2)} \int_{B_1} \int_{B_2} (g_t(\mathbf{x} + \rho_s \mathbf{s}, \mathbf{y} + \rho_r \mathbf{r}) - g_t(\mathbf{x}, \mathbf{y}) - \langle \nabla g_t(\mathbf{x}, \mathbf{y}), (\rho_s \mathbf{s}, \rho_r \mathbf{r}) \rangle) d\mathbf{s} d\mathbf{r} \right|. \end{aligned}$$

2318 Thus, we get

$$\begin{aligned} & |g_{t,\rho}(\mathbf{x}, \mathbf{y}) - g_t(\mathbf{x}, \mathbf{y})| \\ &\leq \int_{B_1} \int_{B_2} |g_t(\mathbf{x} + \rho_s \mathbf{s}, \mathbf{y} + \rho_r \mathbf{r}) - g_t(\mathbf{x}, \mathbf{y}) - \langle \nabla g_t(\mathbf{x}, \mathbf{y}), (\rho_s \mathbf{s}, \rho_r \mathbf{r}) \rangle| d\mathbf{s} d\mathbf{r} \\ &\leq \int_{B_1} \int_{B_2} \frac{\ell_{g,1}}{2} (\rho_s^2 \|\mathbf{s}\|^2 + \rho_r^2 \|\mathbf{r}\|^2) d\mathbf{s} d\mathbf{r} \\ &= \frac{\ell_{g,1} \rho_s^2}{2} \int_{B_1} \|\mathbf{s}\|^2 d\mathbf{s} + \frac{\ell_{g,1} \rho_r^2}{2} \int_{B_2} \|\mathbf{r}\|^2 d\mathbf{r} \\ &= \frac{\ell_{g,1} \rho_s^2}{2} \frac{d_1}{d_1 + 2} + \frac{\ell_{g,1} \rho_r^2}{2} \frac{d_2}{d_2 + 2} \\ &\leq \frac{\ell_{g,1} (\rho_s^2 + \rho_r^2)}{2}, \end{aligned}$$

2333 where the last equality follows since $\frac{1}{\mathcal{V}(d)} \int_{s \in B} \|s\|^p ds = \frac{d}{d+p}$.

2334 The proof of part (b) follows using similar arguments. \square

2347 **Lemma D.4.** Given $\rho = (\rho_s, \rho_r)$ as positive smoothing parameters, let $g_{t,\rho}(\mathbf{x}, \mathbf{y})$ and $f_{t,\rho}(\mathbf{x}, \mathbf{y})$ be the functions defined by (17).

2353 (a) Suppose Assumption B3. holds. Then, we have

$$2355 \quad \|\nabla g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla g_t(\mathbf{x}, \mathbf{y})\| \leq \frac{\ell_{g,1}(\rho_s d_1 + \rho_r d_2)}{2}. \quad (128)$$

2361 (b) Suppose Assumption B2. holds. Then, we have

$$2363 \quad \|\nabla f_t(\mathbf{x}, \mathbf{y}) - \nabla f_{t,\rho}(\mathbf{x}, \mathbf{y})\| \leq \frac{\ell_{f,1}(\rho_s d_1 + \rho_r d_2)}{2}. \quad (129)$$

2365 *Proof.* Let $S(d_1)$ be the surface area of the unit sphere in \mathbb{R}^{d_1} . Moreover, let U_{B_1} be the unit sphere.
 2366

$$\begin{aligned}
 & \| \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}) \| \\
 &= \left\| \frac{1}{S(d_1)} \left(\frac{d_1}{\rho_s} \int_{U_{B_1}} g_t(\mathbf{x} + \rho_s \mathbf{s}, \mathbf{y}) s ds \right) - \nabla_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}) \right\| \\
 &= \left\| \frac{1}{S(d_1)} \left(\frac{d_1}{\rho_s} \int_{U_{B_1}} g_t(\mathbf{x} + \rho_s \mathbf{s}, \mathbf{y}) s ds - \int_{U_{B_1}} \frac{d_1}{\rho_s} g_t(\mathbf{x}, \mathbf{y}) s ds \right. \right. \\
 &\quad \left. \left. - \int_{U_{B_1}} \frac{d_1}{\rho_s} \langle \nabla_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}), \rho_s \mathbf{s} \rangle s ds \right) \right\| \\
 &\leq \frac{d_1}{S(d_1) \rho_s} \int_{U_{B_1}} |g_t(\mathbf{x}_t + \rho_s \mathbf{s}, \mathbf{y}) - g_t(\mathbf{x}, \mathbf{y}) - \langle \nabla_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}), \rho_s \mathbf{s} \rangle| \|s\| ds \\
 &\leq \frac{d_1}{S(d_1) \rho_s} \cdot \frac{\ell_{g,1} \rho_s^2}{2} \int_{U_{B_1}} \|s\|^3 ds \\
 &= \frac{\rho_s d_1 \ell_{g,1}}{2},
 \end{aligned} \tag{130}$$

2384 where the second equality follows from $\int_{U_{B_1}} \mathbf{s} \mathbf{s}^\top ds = \frac{S(d_1)}{d_1} \mathbf{I}$.
 2385

2386 Similarly, let $S(d_2)$ be the surface area of the unit sphere in \mathbb{R}^{d_2} . Moreover, let U_{B_2} be the unit sphere.
 2387

$$\begin{aligned}
 & \| \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}) \| \\
 &= \left\| \frac{1}{S(d_2)} \left(\frac{d_2}{\rho_r} \int_{U_{B_2}} g_t(\mathbf{x}, \mathbf{y} + \rho_r \mathbf{r}) r dr \right) - \nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}) \right\| \\
 &= \left\| \frac{1}{S(d_2)} \left(\frac{d_2}{\rho_r} \int_{U_{B_2}} g_t(\mathbf{x}, \mathbf{y} + \rho_r \mathbf{r}) r dr - \int_{U_{B_2}} \frac{d_2}{\rho_r} g_t(\mathbf{x}, \mathbf{y}) r dr \right. \right. \\
 &\quad \left. \left. - \int_{U_{B_2}} \frac{d_2}{\rho_r} \langle \nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}), \rho_r \mathbf{r} \rangle r dr \right) \right\| \\
 &\leq \frac{d_2}{S(d_2) \rho_r} \int_{U_{B_2}} |g_t(\mathbf{x}_t, \mathbf{y} + \rho_r \mathbf{r}) - g_t(\mathbf{x}, \mathbf{y}) - \langle \nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}), \rho_r \mathbf{r} \rangle| \|r\| dr \\
 &\leq \frac{d_2}{S(d_2) \rho_r} \cdot \frac{\ell_{g,1} \rho_r^2}{2} \int_{U_{B_2}} \|r\|^3 dr \\
 &= \frac{\rho_r d_2 \ell_{g,1}}{2},
 \end{aligned} \tag{131}$$

2405 where the second equality follows from $\int_{U_{B_2}} \mathbf{r} \mathbf{r}^\top dr = \frac{S(d_2)}{d_2} \mathbf{I}$.
 2406

2407 Thus, we get

$$\begin{aligned}
 & \| \nabla g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla g_t(\mathbf{x}, \mathbf{y}) \| \\
 &\leq \| \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}) \| + \| \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}) \| \\
 &\leq \frac{\rho_s d_1 \ell_{g,1}}{2} + \frac{\rho_r d_2 \ell_{g,1}}{2}.
 \end{aligned}$$

2413 Finally, by a similar argument as in Part (a), we obtain
 2414

$$\| \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{x}} f_t(\mathbf{x}, \mathbf{y}) \| \leq \frac{\rho_s d_1 \ell_{f,1}}{2}, \tag{132}$$

2420 and

$$2422 \quad \|\nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\| \leq \frac{\rho_r d_2 \ell_{f,1}}{2}, \quad (133)$$

2424 which implies

$$2426 \quad \|\nabla f_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla f_t(\mathbf{x}, \mathbf{y})\| \leq \frac{(\rho_s d_1 + \rho_r d_2) \ell_{f,1}}{2}. \quad \square$$

2430 **Lemma D.5.** Suppose Assumption B4. holds. Given $\rho = (\rho_s, \rho_r)$ as positive smoothing parameters, let $g_{t,\rho}(\mathbf{x}, \mathbf{y})$ be the
2431 function defined in (17). Then, we have

$$2433 \quad \|\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y})\|^2 \leq \frac{d_2^2 \ell_{g,2}^2}{4} \rho_r^2.$$

2436 *Proof.* Similary, let $S(d_2)$ be the surface area of the unit sphere in \mathbb{R}^{d_2} . Moreover, let U_{B_2} be the unit sphere.

$$\begin{aligned} & \|\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y})\| \\ &= \left\| \frac{1}{S(d_2)} \left(\frac{d_2}{\rho_r} \int_{U_{B_2}} \nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y} + \rho_r \mathbf{r}) \mathbf{r} d\mathbf{r} \right) - \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y}) \right\| \\ &= \left\| \frac{1}{S(d_2)} \left(\frac{d_2}{\rho_r} \int_{U_{B_2}} \nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y} + \rho_r \mathbf{r}) \mathbf{r} d\mathbf{r} - \int_{U_{B_2}} \frac{d_2}{\rho_r} \nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}) \mathbf{r} d\mathbf{r} \right. \right. \\ &\quad \left. \left. - \int_{U_{B_2}} \frac{d_2}{\rho_r} \langle \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y}), \rho_r \mathbf{r} \rangle \mathbf{r} d\mathbf{r} \right) \right\| \\ &\leq \frac{d_2}{S(d_2) \rho_r} \int_{U_{B_2}} |\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y} + \rho_r \mathbf{r}) - \nabla_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}) - \langle \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y}), \rho_r \mathbf{r} \rangle| \|\mathbf{r}\| d\mathbf{r} \\ &\leq \frac{d_2}{S(d_2) \rho_r} \cdot \frac{\ell_{g,2} \rho_r^2}{2} \int_{U_{B_2}} \|\mathbf{r}\|^3 d\mathbf{r} \\ &= \frac{\rho_r d_2 \ell_{g,2}}{2}, \end{aligned}$$

2455 where the second equality follows from $\int_{U_{B_2}} \mathbf{r} \mathbf{r}^\top d\mathbf{r} = \frac{S(d_2)}{d_2} \mathbf{I}$. \square

2457 **Lemma D.6.** Suppose Assumption B3. holds. Let $\hat{\nabla}_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}})$ and $\hat{\nabla}_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}})$ be defined as in (19a) and (19b),
2458 respectively. Then, for any $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_1} \times \mathbb{R}^{d_2}$ and $\rho_r, \rho_s \geq 0$, we have

$$2460 \quad \mathbb{E}_{\mathbf{r}} \left[\|\hat{\nabla}_{\mathbf{y}} g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}}) - \hat{\nabla}_{\mathbf{y}} g_t(\mathbf{x}, \hat{\mathbf{y}}; \bar{\mathcal{B}})\|^2 \right] \leq 3d_2 \ell_{g,1}^2 \|\mathbf{y} - \hat{\mathbf{y}}\|^2 + \frac{3\ell_{g,1}^2 d_2^2 \rho_r^2}{2} \quad \forall \hat{\mathbf{y}} \in \mathbb{R}^{d_2}, \quad (134a)$$

$$2463 \quad \mathbb{E}_{\mathbf{z}} \left[\|\hat{\nabla}_{\mathbf{x}} g_t(\mathbf{x}, \mathbf{y}; \bar{\mathcal{B}}) - \hat{\nabla}_{\mathbf{x}} g_t(\hat{\mathbf{x}}, \mathbf{y}; \bar{\mathcal{B}})\|^2 \right] \leq 3d_1 \ell_{g,1}^2 \|\mathbf{x} - \hat{\mathbf{x}}\|^2 + \frac{3\ell_{g,1}^2 d_1^2 \rho_s^2}{2} \quad \forall \hat{\mathbf{x}} \in \mathbb{R}^{d_1}. \quad (134b)$$

2465 *Proof.* The proof is similar to that of Lemma 5 in (Ji et al., 2019). \square

2467 **Lemma D.7.** Suppose Assumptions 3.2 and B3. hold. Let $\rho = (\rho_s, \rho_r)$ be positive smoothing parameters. Let $\mathbf{y}_t^*(\mathbf{x})$ and
2468 $\hat{\mathbf{y}}_t^*(\mathbf{x})$ be defined in (1) and (18), respectively. Then, we have

$$2470 \quad \mathbb{E} \left[\|\hat{\mathbf{y}}_t^*(\mathbf{x}) - \mathbf{y}_t^*(\mathbf{x})\|^2 \right] \leq \frac{\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g}. \quad (135)$$

2473 *Proof.* From (1), we have $\mathbf{y}_t^*(\mathbf{x}) \in \arg \min_{\mathbf{y} \in \mathbb{R}^{d_2}} g_t(\mathbf{x}, \mathbf{y})$. Since by Assumption 3.2, $g_t(\mathbf{x}, \mathbf{y})$ is μ_g -strongly convex in
2474

2475 terms of \mathbf{y} . Then, by Lemma B.2, we get

$$2477 \quad \|\mathbf{y} - \mathbf{y}_t^*(\mathbf{x})\|^2 \leq \frac{2}{\mu_g} (g_t(\mathbf{x}, \mathbf{y}) - g_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))).$$

2479 By setting $\mathbf{y} = \hat{\mathbf{y}}_t^*(\mathbf{x})$, we have

$$2481 \quad \|\hat{\mathbf{y}}_t^*(\mathbf{x}) - \mathbf{y}_t^*(\mathbf{x})\|^2 \leq \frac{2}{\mu_g} (g_t(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) - g_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))). \quad (136)$$

2483 Similarly, from (18), we have

$$2485 \quad \hat{\mathbf{y}}_t^*(\mathbf{x}) \in \arg \min_{\mathbf{y} \in \mathbb{R}^{d_2}} \{g_{t,\rho}(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{(\mathbf{z}, \mathbf{r})} [g_t(\mathbf{x} + \rho_s \mathbf{s}, \mathbf{y} + \rho_r \mathbf{r}; \zeta)]\}.$$

2488 Since by Assumption 3.2, $g_{t,\rho}(\mathbf{x}, \mathbf{y})$ is μ_g -strongly convex in terms of \mathbf{y} . Then, from Lemma B.2, we have

$$2489 \quad \|\mathbf{y} - \hat{\mathbf{y}}_t^*(\mathbf{x})\|^2 \leq \frac{2}{\mu_g} (g_{t,\rho}(\mathbf{x}, \mathbf{y}) - g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))).$$

2492 By setting $\mathbf{y} = \mathbf{y}_t^*(\mathbf{x})$, we have

$$2494 \quad \|\mathbf{y}_t^*(\mathbf{x}) - \hat{\mathbf{y}}_t^*(\mathbf{x})\|^2 \leq \frac{2}{\mu_g} (g_{t,\rho}(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) - g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))). \quad (137)$$

2496 Summing up (136) and (137), we get

$$2498 \quad \|\mathbf{y}_t^*(\mathbf{x}) - \hat{\mathbf{y}}_t^*(\mathbf{x})\|^2 \leq \frac{1}{\mu_g} (g_{t,\rho}(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) - g_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))) \\ 2500 \quad + \frac{1}{\mu_g} (g_t(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) - g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))),$$

2502 which implies

$$2505 \quad \|\mathbf{y}_t^*(\mathbf{x}) - \hat{\mathbf{y}}_t^*(\mathbf{x})\|^2 \leq \frac{1}{\mu_g} |g_{t,\rho}(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x})) - g_t(\mathbf{x}, \mathbf{y}_t^*(\mathbf{x}))| \\ 2506 \quad + \frac{1}{\mu_g} |g_t(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x})) - g_{t,\rho}(\mathbf{x}, \hat{\mathbf{y}}_t^*(\mathbf{x}))| \\ 2508 \quad \leq \frac{\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g},$$

2512 where the last inequality is by Eq. (126). \square

2513
2514
2515 **Lemma D.8.** Suppose Assumptions 3.2 and 3.3 hold. Let $\mathbf{v}_t^*(\mathbf{x})$ and $\hat{\mathbf{v}}_t^*(\mathbf{x})$ be defined in (4b) and (119b), respectively.
2516 Then, we have

$$2517 \quad \mathbb{E} [\|\hat{\mathbf{v}}_t^*(\mathbf{x}) - \mathbf{v}_t^*(\mathbf{x})\|^2] \leq \frac{d_2^2}{2\mu_g^4} (\ell_{f,1}^2 \mu_g^2 + \ell_{f,2}^2 \ell_{f,0}^2) \rho_r^2. \quad (138)$$

2521 *Proof.* From (4b) and (119b), we have

$$2523 \quad \begin{aligned} & \|\hat{\mathbf{v}}_t^*(\mathbf{x}) - \mathbf{v}_t^*(\mathbf{x})\|^2 \\ &= \|[\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y})]^{-1} \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) - [\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y})]^{-1} \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\|^2 \\ &\leq 2 \|[\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y})]^{-1} \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) - [\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y})]^{-1} \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\|^2 \\ &\quad + 2 \|[\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y})]^{-1} \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y}) - [\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y})]^{-1} \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\|^2. \end{aligned} \quad (139a)$$

$$2528 \quad (139b)$$

2530 Next, we separately bound each of the above terms, (139a) and (139b).

$$\begin{aligned}
 2532 \quad (139a) &\leq 2 \|[\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y})]^{-1}\|^2 \|\nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\|^2 \\
 2533 &\leq \frac{2}{\mu_g^2} \|\nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\|^2 \\
 2534 &\leq \frac{2}{\mu_g^2} \frac{\rho_r^2 d_2^2 \ell_{f,1}^2}{4}, \\
 2535 \\
 2536 \\
 2537
 \end{aligned} \tag{140}$$

2538 where the second inequality holds due to the Assumption 3.2, the third inequality is by (133).

2539 To bound (139b), note that for any invertible matrices \mathbf{A}_1 and \mathbf{A}_2 , we have

$$2541 \quad \|\mathbf{A}_2^{-1} - \mathbf{A}_1^{-1}\| = \|\mathbf{A}_1^{-1}(\mathbf{A}_1 - \mathbf{A}_2)\mathbf{A}_2^{-1}\| \leq \|\mathbf{A}_1^{-1}\| \|\mathbf{A}_2^{-1}\| \|\mathbf{A}_1 - \mathbf{A}_2\|,$$

2543 which implies

$$\begin{aligned}
 2545 \quad (139b) &\leq 2 \|\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y})]^{-1} - [\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y})]^{-1}\|^2 \|\nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\|^2 \\
 2546 &\leq 2 \|\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y})]^{-1}\|^2 \|[\nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y})]^{-1}\|^2 \\
 2547 &\quad \|\nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}}^2 g_t(\mathbf{x}, \mathbf{y})\|^2 \|\nabla_{\mathbf{y}} f_t(\mathbf{x}, \mathbf{y})\|^2 \\
 2548 &\leq \frac{2}{\mu_g^4} \frac{\rho_r^2 d_2^2 \ell_{g,2}^2}{4} \ell_{f,0}^2, \\
 2549 \\
 2550 \\
 2551
 \end{aligned} \tag{141}$$

2552 where the last inequality follows from Lemma D.5.

2553 Using (139)–(140), we obtain the desired result. \square

D.3. Bounds on the Zeroth-Order Inner Solution

2554 **Lemma D.9.** Suppose that Assumptions B3. and D1. hold. Consider the sequence $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ generated by Algorithm 2, and define

$$2560 \quad e_t^{g\rho} := \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{d}}_t^{\mathbf{y}}. \tag{142}$$

2561 Then, we have

$$\begin{aligned}
 2563 \quad \mathbb{E}\|e_{t+1}^{g\rho}\|^2 &\leq (1 - \gamma_{t+1})^2 \mathbb{E}\|e_t^{g\rho}\|^2 + 12(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_{\mathbf{y}} g_{t-1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 2564 &\quad + 9d_2^2 \ell_{g,1}^2 (1 - \gamma_{t+1})^2 \rho_r^2 + 24d_2 \ell_{g,1}^2 (1 - \gamma_{t+1})^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\
 2565 &\quad + 24d_2 \ell_{g,1}^2 (1 - \gamma_{t+1})^2 \mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 + 2\frac{\hat{\sigma}_{gy}^2}{b} \gamma_{t+1}^2. \\
 2566 \\
 2567 \\
 2568 \\
 2569 \\
 2570 \\
 2571 \\
 2572 \\
 2573 \\
 2574 \\
 2575 \\
 2576 \\
 2577 \\
 2578 \\
 2579
 \end{aligned} \tag{143}$$

From the definition of $\hat{\mathbf{d}}_{t+1}^{\mathbf{y}}$ in Algorithm 2, we have

$$\begin{aligned}
 2581 \quad \hat{\mathbf{d}}_{t+1}^{\mathbf{y}} - \hat{\mathbf{d}}_t^{\mathbf{y}} &= -\gamma_{t+1} \hat{\mathbf{d}}_t^{\mathbf{y}} + \gamma_{t+1} \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) \\
 2582 &\quad + (1 - \gamma_{t+1}) \left(\hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1}) \right). \\
 2583 \\
 2584
 \end{aligned}$$

2585 Then, we have

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{y}}\|^2 \\
 &= \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\mathbf{d}}_t^{\mathbf{y}} - (\hat{\mathbf{d}}_{t+1}^{\mathbf{y}} - \hat{\mathbf{d}}_t^{\mathbf{y}})\|^2 \\
 &= \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\mathbf{d}}_t^{\mathbf{y}} + \gamma_{t+1} \hat{\mathbf{d}}_t^{\mathbf{y}} - \gamma_{t+1} \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) \\
 &\quad - (1 - \gamma_{t+1}) (\hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1}))\|^2 \\
 &= \mathbb{E} \|(1 - \gamma_{t+1}) (\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{d}}_t^{\mathbf{y}}) \\
 &\quad + \gamma_{t+1} (\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1})) \\
 &\quad + (1 - \gamma_{t+1}) (\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) \\
 &\quad + \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_t, \mathbf{y}_t) \\
 &\quad - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) + \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1}))\|^2.
 \end{aligned}$$

2600 From (122), we have

$$\begin{aligned}
 & \mathbb{E} [\hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1})] = \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}), \\
 & \mathbb{E} [\hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1})] \\
 &= \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_t, \mathbf{y}_t),
 \end{aligned}$$

2607 then, we have

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{y}}\|^2 \\
 &= (1 - \gamma_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{d}}_t^{\mathbf{y}}\|^2 \\
 &\quad + \mathbb{E} \|\gamma_{t+1} (\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1})) \\
 &\quad + (1 - \gamma_{t+1}) (\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) + \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_t, \mathbf{y}_t) \\
 &\quad - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) + \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1}))\|^2 \\
 &\leq (1 - \gamma_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{d}}_t^{\mathbf{y}}\|^2 \\
 &\quad + 2(1 - \gamma_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) + \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_t, \mathbf{y}_t) \\
 &\quad - \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) + \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 &\quad + 2\gamma_{t+1}^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1})\|^2,
 \end{aligned}$$

2622 where the second inequality holds by Cauchy-Schwarz inequality.

2623 Then, from $\mathbb{E} \|a - \mathbb{E}[a]\|^2 = \mathbb{E}\|a\|^2 - \|\mathbb{E}[a]\|^2$ and Assumption 4.1, we have

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{y}}\|^2 \\
 &\leq (1 - \gamma_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{d}}_t^{\mathbf{y}}\|^2 \\
 &\quad + 4(1 - \gamma_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 &\quad + 4(1 - \gamma_{t+1})^2 \mathbb{E} \|\hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t; \bar{\mathcal{B}}_{t+1})\|^2 + 2\gamma_{t+1}^2 \frac{\hat{\sigma}_{g_{\mathbf{y}}}^2}{b} \\
 &\leq (1 - \gamma_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{d}}_t^{\mathbf{y}}\|^2 \\
 &\quad + 4(1 - \gamma_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 &\quad + 12(1 - \gamma_{t+1})^2 d_2 \ell_{g, 1}^2 \mathbb{E} \|(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - (\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 &\quad + 3(1 - \gamma_{t+1})^2 \ell_{g, 1}^2 d_2^2 \rho_r^2 + 2\gamma_{t+1}^2 \frac{\hat{\sigma}_{g_{\mathbf{y}}}^2}{b},
 \end{aligned}$$

2637 where the second inequality follows from Young's inequality and Lemma D.6.

From Eq. (131), we have

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 & \leq 3\mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 & + 3\mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 & + 3\mathbb{E} \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 & \leq 3\mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 + \frac{3\rho_r^2 d_2^2 \ell_{g,1}^2}{2}.
 \end{aligned}$$

Finally, we get

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1,\rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\mathbf{d}}_t^{\mathbf{y}}\|^2 \leq (1 - \gamma_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{d}}_t^{\mathbf{y}}\|^2 \\
 & + 12(1 - \gamma_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 + 6(1 - \gamma_{t+1})^2 \rho_r^2 d_2^2 \ell_{g,1}^2 \\
 & + 12(1 - \gamma_{t+1})^2 d_2 \ell_{g,1}^2 \mathbb{E} \|(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - (\mathbf{x}_t, \mathbf{y}_t)\|^2 + 3(1 - \gamma_{t+1})^2 \ell_{g,1}^2 d_2^2 \rho_r^2 + 2\gamma_{t+1}^2 \frac{\hat{\sigma}_{g_y}^2}{b}.
 \end{aligned}$$

□

Lemma D.10. Suppose Assumptions 3.2 and B3. hold. Then, for the sequence $\{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^T$ generated by Algorithm 2, we have

$$\begin{aligned}
 \mathbb{E} [\|\mathbf{y}_{t+1} - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] & \leq (1 + a) \left(1 - 2\beta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \right) \mathbb{E} [\|\mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\
 & + \left(-(1 + a) \left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2 \right) \right) \mathbb{E} [\|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\
 & + (1 + \frac{1}{a}) \beta_t^2 \mathbb{E} [\|e_t^{g,\rho}\|^2],
 \end{aligned}$$

where $a > 0$ is a constant, $e_t^{g,\rho}$ is defined in (142), and $\hat{\mathbf{y}}_t^*(\mathbf{x}_t)$ is defined in (18).

Proof. From Lemma B.5, we have

$$\begin{aligned}
 \mathbb{E} [\|\mathbf{y}_{t+1} - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] & = \mathbb{E} [\|\mathbf{y}_t - \beta_t \hat{\mathbf{d}}_t^{\mathbf{y}} - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\
 & \leq (1 + a) \mathbb{E} [\|\mathbf{y}_t - \beta_t \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\
 & + (1 + \frac{1}{a}) \beta_t^2 \mathbb{E} [\|\hat{\mathbf{d}}_t^{\mathbf{y}} - \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2].
 \end{aligned} \tag{144}$$

Next, we will separately bound the first term on the RHS of the above inequality.

We have

$$\begin{aligned}
 \mathbb{E} [\|\mathbf{y}_t - \beta_t \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] & = \mathbb{E} [\|\mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] + \beta_t^2 \mathbb{E} [\|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\
 & - 2\beta_t \mathbb{E} [\langle \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t), \mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t) \rangle] \\
 & \leq \left(1 - 2\beta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \right) \mathbb{E} [\|\mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\
 & - \left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2 \right) \mathbb{E} [\|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2],
 \end{aligned} \tag{145}$$

where the inequality results from the strong convexity of $g_{t,\rho}$ by Assumption 3.2, which implies

$$\langle \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t), \mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t) \rangle \geq \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \|\mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2 + \frac{1}{\mu_g + \ell_{g,1}} \|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2.$$

Substituting (145) into (144), gives the desired result.

□

2695 For notational brevity in the analysis, we define

$$2697 \hat{\theta}_t^y := \|\mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2, \quad \hat{\theta}_t^v := \|\mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2, \quad (146)$$

2698 where $\hat{\mathbf{y}}_t^*(\mathbf{x})$ and $\hat{\mathbf{v}}_t^*(\mathbf{x})$ are defined in (18) and (119b), respectively.

Lemma D.11. Suppose Assumptions 3.2, B2, and B3, hold. Let $\hat{\theta}_t^y$ be defined in (146). Then, for the sequence $\{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^T$ generated by Algorithm 2 guarantees the following bound:

$$\begin{aligned} & \sum_{t=1}^T \left(\mathbb{E}[\hat{\theta}_{t+1}^y] - \mathbb{E}[\hat{\theta}_t^y] \right) \\ & \leq \left(-\frac{L_{\mu_g}}{2} \sum_{t=1}^T \mathbb{E}[\hat{\theta}_t^y] + \frac{2}{L_{\mu_g}} \sum_{t=1}^T \mathbb{E}[\|e_t^{g,\rho}\|^2] \right) \beta_t + \frac{4L_y^2}{L_{\mu_g}} \sum_{t=1}^T \mathbb{E}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \frac{1}{\beta_t} \\ & \quad + \sum_{t=1}^T \left(\frac{24\ell_{g,1}}{L_{\mu_g}\mu_g} (\rho_s^2 + \rho_r^2) + \frac{12}{L_{\mu_g}} \sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{y}_{t-1}^*(\mathbf{x}) - \mathbf{y}_t^*(\mathbf{x})\|^2 \right) \frac{1}{\beta_t} \\ & \quad + \sum_{t=1}^T \left(-\frac{2\beta_t}{\mu_g + \ell_{g,1}} + \beta_t^2 \right) \mathbb{E}[\|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2], \end{aligned} \quad (147)$$

2715 where $L_{\mu_g} = \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}$, and $L_y = \frac{\ell_{g,1}}{\mu_g}$ is defined as in (42).

2723 *Proof.* From Lemma B.5, we have for any $c > 0$

$$\begin{aligned} 2725 \mathbb{E}[\|\mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1})\|^2] &= \mathbb{E}[\|\mathbf{y}_{t+1} - \hat{\mathbf{y}}_t^*(\mathbf{x}_t) + \hat{\mathbf{y}}_t^*(\mathbf{x}_t) - \hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1})\|^2] \\ 2726 &\leq (1+c) \mathbb{E}[\|\mathbf{y}_{t+1} - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\ 2727 &\quad + \left(1 + \frac{1}{c}\right) \mathbb{E}[\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2]. \end{aligned} \quad (148)$$

2730 From Lemma D.10, we have for any $a > 0$

$$\begin{aligned} 2732 \mathbb{E}[\|\mathbf{y}_{t+1} - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] &\leq (1+a) \left(1 - 2\beta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}\right) \mathbb{E}[\|\mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\ 2733 &\quad + \left(-(1+a) \left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2\right)\right) \mathbb{E}[\|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\ 2734 &\quad + \left(1 + \frac{1}{a}\right) \beta_t^2 \mathbb{E}[\|e_t^{g,\rho}\|^2]. \end{aligned} \quad (149)$$

2739 Substituting (149) into (148), we get

$$\begin{aligned} 2741 \mathbb{E}[\|\mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1})\|^2] &\leq (1+c)(1+a) \left(1 - 2\beta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}\right) \mathbb{E}[\|\mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\ 2743 &\quad + \left(-(1+c)(1+a) \left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2\right)\right) \mathbb{E}[\|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\ 2745 &\quad + (1+c)(1+\frac{1}{a}) \beta_t^2 \mathbb{E}[\|e_t^{g,\rho}\|^2] + \left(1 + \frac{1}{c}\right) \mathbb{E}[\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2]. \end{aligned} \quad (150)$$

Choose $c = \frac{\beta_t L_{\mu_g}/2}{1-\beta_t L_{\mu_g}}$ and $a = \frac{\beta_t L_{\mu_g}}{1-2\beta_t L_{\mu_g}}$. Then, the following equations and inequalities are satisfied.

$$\begin{aligned} (1+c)(1+a)(1-2\beta_t L_{\mu_g}) &= 1 - \frac{\beta_t L_{\mu_g}}{2}, \\ (1+a)(1-2\beta_t L_{\mu_g}) &= 1 - \beta_t L_{\mu_g}, \\ (1+c)(1-\beta_t L_{\mu_g}) &= 1 - \frac{\beta_t L_{\mu_g}}{2}, \\ 1 + \frac{1}{a} &\leq \frac{1}{\beta_t L_{\mu_g}}, \quad 1 + \frac{1}{c} \leq \frac{2}{\beta_t L_{\mu_g}}, \end{aligned} \tag{151}$$

where $L_{\mu_g} = \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}$. Based on (150) and (151), we get

$$\begin{aligned} &\mathbb{E} [\|\mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1})\|^2] - \mathbb{E} [\|\mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\ &\leq -\frac{\beta_t L_{\mu_g}}{2} \mathbb{E} [\|\mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] + \left(-\left(\frac{2\beta_t}{\mu_g + \ell_{g,1}} - \beta_t^2 \right) \right) \mathbb{E} [\|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\ &\quad + \frac{2}{\beta_t L_{\mu_g}} \beta_t^2 \mathbb{E} [\|e_t^{g,\rho}\|^2] + \frac{2}{\beta_t L_{\mu_g}} \mathbb{E} [\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2]. \end{aligned} \tag{152}$$

Next, we upper-bound the last term of the above inequality.

$$\begin{aligned} &\mathbb{E} [\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\ &\leq 2 (\mathbb{E} [\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t)\|^2] + \mathbb{E} [\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2]) \\ &\leq 2 (L_y^2 \mathbb{E} [\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2] + \|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2), \end{aligned} \tag{153}$$

where the second inequality is by Lemma D.2.

Moreover, from Lemma D.7, we get

$$\begin{aligned} \mathbb{E} [\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] &\leq 3\mathbb{E} [\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_{t+1}^*(\mathbf{x}_t)\|^2] \\ &\quad + 3\mathbb{E} [\|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] + 3\mathbb{E} [\|\mathbf{y}_t^*(\mathbf{x}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\ &\leq 3\mathbb{E} [\|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] + \frac{6\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g}. \end{aligned} \tag{154}$$

Combining (153) and (154) yields

$$\begin{aligned} &\mathbb{E} [\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2] \\ &\leq 2 \left(L_y^2 \mathbb{E} [\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2] + 3\mathbb{E} [\|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2] + \frac{6\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g} \right). \end{aligned} \tag{155}$$

Substituting (155) into (152) and summing over $t \in [T]$, give the desired result. \square

D.4. Bounds on the Zeroth-Order System Solution

Lemma D.12. Suppose Assumptions B2. and B3. hold. Then, for the sequence $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ generated by Algorithm 2, we have

$$\begin{aligned} &\mathbb{E} \|\hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\ &\leq (12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_v^2}) d_2 \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + (12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_v^2}) d_2 \mathbb{E} \|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\ &\quad + \frac{9}{2} d_2 \ell_{g,1}^2 \mathbb{E} \|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + (3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_r^2}) d_2^2 \rho_r^2, \end{aligned}$$

where $\hat{\nabla}_{\mathbf{y}} f_t$ and $\hat{\nabla}_{\mathbf{y}}^2 g_t$ are defined in (20a) and (21a), respectively.

2805 Proof. From Lemma D.6, we have

$$\begin{aligned}
 & \|\hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1})\|^2 \\
 & \leq 3d_2 \ell_{f,1}^2 \|\mathbf{z}_{t+1} - \mathbf{z}_t\|^2 + \frac{3}{2} \ell_{f,1}^2 d_2^2 \rho_{\mathbf{r}}^2 \\
 & \leq 6d_2 \ell_{f,1}^2 \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6d_2 \ell_{f,1}^2 \|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 + \frac{3}{2} \ell_{f,1}^2 d_2^2 \rho_{\mathbf{r}}^2.
 \end{aligned} \tag{156}$$

2812 Moreover, from (21a), we have

$$\begin{aligned}
 & \|\hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 & = \frac{1}{4\rho_{\mathbf{v}}^2} \|\hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} + \rho_{\mathbf{v}} \mathbf{v}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 & \leq \frac{3}{4\rho_{\mathbf{v}}^2} d_2 \ell_{g,1}^2 \|(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} + \rho_{\mathbf{v}} \mathbf{v}_{t+1}) - (\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 + \frac{3}{8\rho_{\mathbf{v}}^2} \ell_{g,1}^2 d_2^2 \rho_{\mathbf{r}}^2 \\
 & \leq \frac{9}{4\rho_{\mathbf{v}}^2} d_2 \ell_{g,1}^2 \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \frac{9}{4\rho_{\mathbf{v}}^2} d_2 \ell_{g,1}^2 \|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 & + \frac{9}{4} d_2 \ell_{g,1}^2 \|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + \frac{3}{8\rho_{\mathbf{v}}^2} \ell_{g,1}^2 d_2^2 \rho_{\mathbf{r}}^2,
 \end{aligned} \tag{157}$$

2824 where the first inequality follows from Lemma D.6.

2825 From $\|a + b\|^2 \leq 2(\|a\|^2 + \|b\|^2)$, we get

$$\begin{aligned}
 & \|\hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 & \leq 2\|\hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 & + 2\|\hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1})\|^2 \\
 & \leq (12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2})d_2 \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + (12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2})d_2 \|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 & + \frac{9}{2} d_2 \ell_{g,1}^2 \|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + (3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_{\mathbf{v}}^2})d_2^2 \rho_{\mathbf{r}}^2,
 \end{aligned}$$

2837 where the second inequality follows from (156) and (157). □

2851 **Lemma D.13.** Suppose Assumptions B2., B3., D1., and D3. hold. Consider the sequence $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ generated by Algorithm 2, and define

$$e_{t+1}^M := \nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{v}}, \quad \text{where} \tag{158}$$

$$\begin{aligned}
 \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) &= \frac{1}{2\rho_{\mathbf{v}}} (\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} + \rho_{\mathbf{v}} \mathbf{v}_{t+1}) \\
 &\quad - \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} - \rho_{\mathbf{v}} \mathbf{v}_{t+1})). \tag{159}
 \end{aligned}$$

2860 Then, we have

$$\begin{aligned}
 \mathbb{E}\|e_{t+1}^M\|^2 &\leq (1 - \lambda_{t+1})^2 \mathbb{E}\|e_t^M\|^2 + 36\mathbb{E}\|\nabla_y f_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_y f_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 &\quad + \left(18d_2^2\ell_{f,1}^2 + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_v^2})d_2^2\right)\rho_r^2 + 18d_2^2\ell_{g,1}^2\frac{\rho_r^2}{\rho_v^2} \\
 &\quad + \frac{18}{\rho_v^2}\mathbb{E}\|\nabla_y g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \\
 &\quad + \frac{18}{\rho_v^2}\mathbb{E}\|\nabla_y g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t) - \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t)\|^2 \\
 &\quad + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_v^2})d_2\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_v^2})d_2\mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 &\quad + 27d_2\ell_{g,1}^2\mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + 3(\frac{\hat{\sigma}_{g,y}^2}{b\rho_v^2} + \frac{\hat{\sigma}_{f,y}^2}{b})\lambda_{t+1}^2. \tag{160}
 \end{aligned}$$

2881 *Proof.* According to the definition of $\hat{\mathbf{d}}_t^v$ in Algorithm 2, we have

$$\begin{aligned}
 \hat{\mathbf{d}}_{t+1}^v - \hat{\mathbf{d}}_t^v &= -\lambda_{t+1}\hat{\mathbf{d}}_t^v + \lambda_{t+1}(\hat{\nabla}_y f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})) \\
 &\quad + (1 - \lambda_{t+1})\left(\hat{\nabla}_y f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})\right. \\
 &\quad \left.- \hat{\nabla}_y f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\right).
 \end{aligned}$$

2888 Then we have

$$\begin{aligned}
 \mathbb{E}\|\nabla_y f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^v\|^2 &= \mathbb{E}\|\nabla_y f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}) - (\hat{\mathbf{d}}_{t+1}^v - \hat{\mathbf{d}}_t^v)\|^2 \\
 &= \mathbb{E}\|\nabla_y f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_t^v + \lambda_{t+1}\hat{\mathbf{d}}_t^v \\
 &\quad - \lambda_{t+1}\left(\hat{\nabla}_y f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})\right) \\
 &\quad - (1 - \lambda_{t+1})\left(\hat{\nabla}_y f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})\right. \\
 &\quad \left.- \hat{\nabla}_y f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\right)\|^2 \\
 &= \mathbb{E}\|(1 - \lambda_{t+1})(\nabla_y f_{t,\rho}(\mathbf{z}_t) + \tilde{\nabla}_y^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^v) \\
 &\quad + \lambda_{t+1}(\nabla_y f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\nabla}_y f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})) \\
 &\quad + (1 - \lambda_{t+1})\left(\nabla_y f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}) - \nabla_y f_{t,\rho}(\mathbf{z}_t) - \tilde{\nabla}_y^2 g_t(\mathbf{z}_t)\right. \\
 &\quad \left.+ \nabla_y f_{t+1,\rho}(\mathbf{z}_t) + \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_t) - \nabla_y f_{t+1,\rho}(\mathbf{z}_t) - \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_t)\right. \\
 &\quad \left.- \hat{\nabla}_y f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) + \hat{\nabla}_y f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) + \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\right)\|^2.
 \end{aligned}$$

2907 Since

$$\begin{aligned}
 \mathbb{E}\left[\hat{\nabla}_y f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})\right] &= \nabla_y f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}), \\
 \mathbb{E}\left[\hat{\nabla}_y f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_y f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_y^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\right] &= \nabla_y f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_{t+1}) - \nabla_y f_{t+1,\rho}(\mathbf{z}_t) - \tilde{\nabla}_y^2 g_{t+1}(\mathbf{z}_t),
 \end{aligned}$$

2915 then, we have

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{v}}\|^2 \\
 &= (1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{v}}\|^2 \\
 &+ \|\lambda_{t+1} (\nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})) \\
 &+ (1 - \lambda_{t+1}) (\nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t) \\
 &+ \nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t) - \nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t) \\
 &- \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) + \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1}))\|^2 \\
 &\leq (1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{v}}\|^2 \\
 &+ 3(1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t) \\
 &+ \nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t) - \nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t) \\
 &- \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) + \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 &+ 3\lambda_{t+1}^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_{t+1}) - \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1})\|^2 \\
 &+ 3\lambda_{t+1}^2 \mathbb{E} \|\tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})\|^2,
 \end{aligned} \tag{161}$$

2935 where the second inequality holds by Cauchy-Schwarz inequality.

2936 Note that, for the last term on the right-hand side of (161), from (21a) and (159), we have

$$\begin{aligned}
 & \|\tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})\|^2 \\
 &\leq 2 \left\| \frac{1}{2\rho_{\mathbf{v}}} (\nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} + \rho_{\mathbf{v}} \mathbf{v}_{t+1}) - \hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} + \rho_{\mathbf{v}} \mathbf{v}_{t+1}; \bar{\mathcal{B}}_{t+1})) \right\|^2 \\
 &+ 2 \left\| \frac{1}{2\rho_{\mathbf{v}}} (\hat{\nabla}_{\mathbf{y}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} - \rho_{\mathbf{v}} \mathbf{v}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_{\mathbf{y}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} - \rho_{\mathbf{v}} \mathbf{v}_{t+1})) \right\|^2 \\
 &\leq \frac{\hat{\sigma}_{g_{\mathbf{y}}}^2}{b\rho_{\mathbf{v}}^2},
 \end{aligned}$$

2946 where the last inequality follows from Assumption 4.1.

2947 Then, from $\mathbb{E} \|a - \mathbb{E}[a]\|^2 = \mathbb{E}\|a\|^2 - \|\mathbb{E}[a]\|^2$ and Assumption 4.1, we have

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{v}}\|^2 \\
 &\leq (1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{v}}\|^2 \\
 &+ 6(1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t) - \nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t)\|^2 \\
 &+ 6(1 - \lambda_{t+1})^2 \mathbb{E} \|\hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) \\
 &- \hat{\nabla}_{\mathbf{y}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 + 3\lambda_{t+1}^2 \left(\frac{\hat{\sigma}_{g_{\mathbf{y}}}^2}{b\rho_{\mathbf{v}}^2} + \frac{\hat{\sigma}_{f_{\mathbf{y}}}^2}{b} \right).
 \end{aligned}$$

2970 Then, from Young's inequality and Lemma D.12, we obtain
 2971

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{y}} f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{v}}\|^2 \\
 & \leq (1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{v}}\|^2 \\
 & + 12(1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_{t+1,\rho}(\mathbf{z}_t) - \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{z}_t)\|^2 \\
 & + 12(1 - \lambda_{t+1})^2 \mathbb{E} \|\tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t)\|^2 \\
 & + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2}) d_2 \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2}) d_2 \|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 & + 27d_2 \ell_{g,1}^2 \|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_{\mathbf{v}}^2}) d_2^2 \rho_{\mathbf{r}}^2 + 3\lambda_{t+1}^2 (\frac{\hat{\sigma}_{g_{\mathbf{y}}}^2}{b\rho_{\mathbf{v}}^2} + \frac{\hat{\sigma}_{f_{\mathbf{y}}}^2}{b}),
 \end{aligned} \tag{162}$$

2982 For the third term on the right-hand side of (162), based on (159), we have
 2983

$$\begin{aligned}
 & \|\tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 & \leq \frac{1}{2\rho_{\mathbf{v}}^2} \|\nabla_{\mathbf{y}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t)\|^2
 \end{aligned} \tag{163a}$$

$$+ \frac{1}{2\rho_{\mathbf{v}}^2} \|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)\|^2. \tag{163b}$$

2990 For (163a), we get
 2991

$$\begin{aligned}
 & \|\nabla_{\mathbf{y}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 \\
 & \leq 3 \|\nabla_{\mathbf{y}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 \\
 & + 3 \|\nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 \\
 & + 3 \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 \\
 & \leq 3 \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 + \frac{3\rho_{\mathbf{r}}^2 d_2^2 \ell_{g,1}^2}{2},
 \end{aligned}$$

3000 where the last inequality follows from Eq. (131).
 3001 Similary, for (163b), we have
 3002

$$\begin{aligned}
 & \|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 \\
 & \leq 3 \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 + \frac{3\rho_{\mathbf{r}}^2 d_2^2 \ell_{g,1}^2}{2}.
 \end{aligned}$$

3007 Substituting the above inequalities in (163), we have
 3008

$$\begin{aligned}
 & \|\tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \leq \frac{3}{2\rho_{\mathbf{v}}^2} \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 \\
 & + \frac{3}{2\rho_{\mathbf{v}}^2} \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 + \frac{3\rho_{\mathbf{r}}^2 d_2^2 \ell_{g,1}^2}{2\rho_{\mathbf{v}}^2}.
 \end{aligned} \tag{164}$$

3014 For the second term on the right-hand side of (162), we have
 3015

$$\begin{aligned}
 & \|\nabla_{\mathbf{y}} f_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 & \leq 3 \|\nabla_{\mathbf{y}} f_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 & + 3 \|\nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 & + 3 \|\nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 & \leq 3 \|\nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_t, \mathbf{y}_t)\|^2 + \frac{3\rho_{\mathbf{r}}^2 d_2^2 \ell_{f,1}^2}{2},
 \end{aligned} \tag{165}$$

3023 where the last inequality follows from Eq. (133).
 3024

3025 From (164), (165) and (162), we get

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{y}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{y}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{v}}\|^2 \\
 & \leq (1 - \lambda_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{v}}\|^2 \\
 & + 36 \|\nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} f_{t+1}(\mathbf{x}_t, \mathbf{y}_t)\|^2 + 18\rho_{\mathbf{r}}^2 d_2^2 \ell_{f,1}^2 \\
 & + \frac{18}{\rho_{\mathbf{v}}^2} \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 \\
 & + \frac{18}{\rho_{\mathbf{v}}^2} \|\nabla_{\mathbf{y}} g_t(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)\|^2 + \frac{18\rho_{\mathbf{r}}^2 d_2^2 \ell_{g,1}^2}{\rho_{\mathbf{v}}^2} \\
 & + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2}) d_2 \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2}) d_2 \|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 & + 27d_2 \ell_{g,1}^2 \|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_{\mathbf{v}}^2}) d_2^2 \rho_{\mathbf{r}}^2 + 3\lambda_{t+1}^2 (\frac{\hat{\sigma}_{g_y}^2}{b\rho_{\mathbf{v}}^2} + \frac{\hat{\sigma}_{f_y}^2}{b}).
 \end{aligned}$$

□

3041 **Lemma D.14.** Suppose Assumptions B3. and B4. hold. Let

$$e_t^H := \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}}^2 g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t, \quad (166a)$$

$$e_t^J := \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{xy}}^2 g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t, \quad (166b)$$

3042 where

$$\tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) = \frac{1}{2\rho_{\mathbf{v}}} (\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)),$$

$$\tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) = \frac{1}{2\rho_{\mathbf{v}}} (\nabla_{\mathbf{xy}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{xy}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)).$$

3050 Then, for $(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)$ presented to Algorithm 2, we have

3057 (a)

$$\mathbb{E} [\|e_t^H\|^2] \leq \ell_{g,2}^2 \rho_{\mathbf{v}}^2 p^4. \quad (167a)$$

3061 (b)

$$\mathbb{E} [\|e_t^J\|^2] \leq \ell_{g,2}^2 \rho_{\mathbf{v}}^2 p^4. \quad (167b)$$

3067 **Proof. For part (a):** From Lemma D.1, We have

$$\begin{aligned}
 \mathbb{E} [\|e_t^H\|] &= \mathbb{E} [\|\tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}}^2 g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t\|] \\
 &\leq \frac{1}{2\rho_{\mathbf{v}}} \mathbb{E} [\|\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}}^2 g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) \rho_{\mathbf{v}} \mathbf{v}_t\|] \\
 &+ \frac{1}{2\rho_{\mathbf{v}}} \mathbb{E} [\|\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}}^2 g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) \rho_{\mathbf{v}} \mathbf{v}_t\|] \\
 &\leq \ell_{g,2} \rho_{\mathbf{v}} \mathbb{E} [\|\mathbf{v}_t\|^2] \\
 &\leq \ell_{g,2} \rho_{\mathbf{v}} p^2,
 \end{aligned} \quad (168)$$

3078 where the last inequality follows from (8).

3080 **For part (b):** From Lemma D.1, We have

$$\begin{aligned}
 3082 \quad \mathbb{E} [\|e_t^J\|] &= \mathbb{E} [\|\tilde{\nabla}_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t\|] \\
 3083 &\leq \frac{1}{2\rho_{\mathbf{v}}} \mathbb{E} [\|\nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) \rho_{\mathbf{v}} \mathbf{v}_t\|] \\
 3084 &+ \frac{1}{2\rho_{\mathbf{v}}} \mathbb{E} [\|\nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{x}\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) \rho_{\mathbf{v}} \mathbf{v}_t\|] \\
 3085 &\leq \ell_{g,2} \rho_{\mathbf{v}} \mathbb{E} [\|\mathbf{v}_t\|^2] \\
 3086 &\leq \ell_{g,2} \rho_{\mathbf{v}} p^2,
 \end{aligned} \tag{169}$$

3091 where the last inequality follows from (8). \square

3094 **Lemma D.15.** Suppose Assumptions B2., B3. and B4. hold. Then, for directions $\hat{\mathbf{d}}_t^{\mathbf{y}}$ and $\hat{\mathbf{d}}_t^{\mathbf{x}}$ presented to Algorithm 2, and

3096 (a) $\mathbf{d}_{t,\rho}^{\mathbf{y}}$ defined in (120b), we have

$$\mathbb{E} [\|\hat{\mathbf{d}}_t^{\mathbf{y}} - \mathbf{d}_{t,\rho}^{\mathbf{y}}\|^2] \leq 2\mathbb{E} [\|e_t^M\|^2] + 2\ell_{g,2}^2 \rho_{\mathbf{v}}^2 p^4 := B_t, \tag{170a}$$

3101 where $e_t^M := \hat{\mathbf{d}}_t^{\mathbf{y}} - \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)$, with $\tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)$ is defined in (171).

3103 (b) $\mathbf{d}_{t,\rho}^{\mathbf{x}}$ defined in (120c), we have

$$\mathbb{E} [\|\hat{\mathbf{d}}_t^{\mathbf{x}} - \mathbf{d}_{t,\rho}^{\mathbf{x}}\|^2] \leq 2\mathbb{E} [\|e_t^L\|^2] + 2\ell_{g,2}^2 \rho_{\mathbf{v}}^2 p^4, \tag{170b}$$

3108 where $e_t^L := \hat{\mathbf{d}}_t^{\mathbf{x}} - \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)$ with $\tilde{\nabla}_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)$ is defined in (176).

3112 *Proof.* **For part (a):** Let

$$\tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) = \frac{1}{2\rho_{\mathbf{v}}} (\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)). \tag{171}$$

3116 According to the definition of $\mathbf{d}_{t,\rho}^{\mathbf{y}}$ in (120b), we have

$$\begin{aligned}
 3118 \quad \mathbb{E} [\|\hat{\mathbf{d}}_t^{\mathbf{y}} - \mathbf{d}_{t,\rho}^{\mathbf{y}}\|^2] &= \mathbb{E} [\|\hat{\mathbf{d}}_t^{\mathbf{y}} - \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}\|^2] \\
 3119 &\leq 2\mathbb{E} [\|\hat{\mathbf{d}}_t^{\mathbf{y}} - \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2] \\
 3120 &+ 2\mathbb{E} [\|\tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}\|^2].
 \end{aligned} \tag{172a}$$

$$\tag{172b}$$

3125 Next, we separately bound (172a)–(172b) on the RHS of the above inequality.

3126 **Bounding (172a).** We have

$$3128 \quad 2\mathbb{E} [\|\hat{\mathbf{d}}_t^{\mathbf{y}} - \nabla_{\mathbf{y}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2] := 2\mathbb{E} [\|e_t^M\|^2]. \tag{173}$$

3131 **Bounding (172b).** From Lemmas D.1 and D.14, we have

$$\tag{172b} = \mathbb{E} [\|e_t^H\|^2] \leq 3\ell_{g,2}^2 \rho_{\mathbf{v}}^2 p^4. \tag{174}$$

Combining (173)–(174) yields

$$\mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^{\mathbf{v}} - \mathbf{d}_{t,\rho}^{\mathbf{v}} \right\|^2 \right] \leq 2\mathbb{E} \left[\|e_t^M\|^2 \right] + 2\ell_{g,2}^2 \rho_{\mathbf{v}}^2 p^4. \quad (175)$$

For part (b): Let

$$\tilde{\nabla}_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) = \frac{1}{2\rho_{\mathbf{v}}} (\nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}} \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}} \mathbf{v}_t)). \quad (176)$$

According to the definition of $\mathbf{d}_{t,\rho}^{\mathbf{x}}$ in (120c), we have

$$\begin{aligned} \mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^{\mathbf{x}} - \mathbf{d}_{t,\rho}^{\mathbf{x}} \right\|^2 \right] &= \mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^{\mathbf{x}} - \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v} \right\|^2 \right] \\ &\leq 2\mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^{\mathbf{x}} - \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) \right\|^2 \right] \end{aligned} \quad (177a)$$

$$+ 2\mathbb{E} \left[\left\| \tilde{\nabla}_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}\mathbf{y}}^2 g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) \mathbf{v}_t \right\|^2 \right]. \quad (177b)$$

Next, we separately bound (177a)–(177b) on the RHS of the above inequality.

Bounding (177a). We have

$$2\mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^{\mathbf{x}} - \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{x}\mathbf{y}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) \right\|^2 \right] := 2\mathbb{E} \left[\|e_t^L\|^2 \right]. \quad (178)$$

Bounding (177b). From Lemmas D.1 and D.14, we have

$$(177b) = \mathbb{E} \left[\|e_t^J\|^2 \right] \leq 2\ell_{g,2}^2 \rho_{\mathbf{v}}^2 p^4. \quad (179)$$

Combining (178)–(179) yields

$$\mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^{\mathbf{x}} - \mathbf{d}_{t,\rho}^{\mathbf{x}} \right\|^2 \right] \leq 2\mathbb{E} \left[\|e_t^L\|^2 \right] + 2\ell_{g,2}^2 \rho_{\mathbf{v}}^2 p^4.$$

□

Lemma D.16. Suppose Assumptions 3.2, B1, and B3 hold. Set the step size δ_t and the parameter p in (8), as

$$\delta_t \leq \left(2 + \frac{1}{\ell_{g,1}^2} \right) \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}, \quad \forall t \in [T], \quad \text{and} \quad p = \frac{\ell_{f,0}}{\mu_g}. \quad (180)$$

Then, for the sequence $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ generated by Algorithm 2, we have

$$\mathbb{E} \left[\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2 \right] \leq (1 + \hat{a}) \left(1 - \delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \right) \mathbb{E}[\hat{\theta}_t^{\mathbf{v}}] + \left(1 + \frac{1}{\hat{a}} \right) \delta_t^2 B_t,$$

for some $\hat{a} > 0$, where $\hat{\theta}_t^{\mathbf{v}}$ and B_t are defined in Eq. (146) and Lemma D.15, respectively.

3190 *Proof.* By setting the radius $p := \frac{\ell_{f,0}}{\mu_g}$ in (8), we have

$$\begin{aligned}
 3192 \quad & \mathbb{E} [\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2] = \mathbb{E} \left[\left\| \Pi_{\mathcal{Z}_p} [\mathbf{v}_t - \delta_t \hat{\mathbf{d}}_t^\mathbf{v}] - \Pi_{\mathcal{Z}_p} [\hat{\mathbf{v}}_t^*(\mathbf{x}_t)] \right\|^2 \right] \\
 3193 \quad & \leq \mathbb{E} [\|\mathbf{v}_t - \delta_t \hat{\mathbf{d}}_t^\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2] \\
 3194 \quad & \leq (1 + \hat{a}) \underbrace{\mathbb{E} [\|\mathbf{v}_t - \delta_t \mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t) - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2]}_{I_t} \\
 3195 \quad & + \left(1 + \frac{1}{\hat{a}}\right) \delta_t^2 \underbrace{\mathbb{E} [\|\hat{\mathbf{d}}_t^\mathbf{v} - \mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2]}_{K_t}, \tag{181}
 \end{aligned}$$

3202 where $\mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)$ is defined in (120b); the first inequality follows from non-expansiveness property of a projection
3203 operator.

3204 We next bound the I_t , and K_t terms in (181), respectively.

3205 **Bounding I_t .** We have

$$\begin{aligned}
 3208 \quad I_t &= \mathbb{E} [\|\mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2] - 2\delta_t \mathbb{E} [\langle \mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t), \mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t) \rangle] + \delta_t^2 \mathbb{E} [\|\mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2] \\
 3209 \quad &\leq \left(1 - 2\delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}\right) \mathbb{E} [\|\mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2] \\
 3210 \quad &- \left(2\delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} - \delta_t^2\right) \mathbb{E} [\|\mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2],
 \end{aligned}$$

3214 where the inequality holds since $\mathbf{d}_{t,\rho}^\mathbf{v}$ in (120) is the gradient of the strongly convex quadratic program $\frac{1}{2}\mathbf{v}^\top \nabla_y^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) \mathbf{v} +$
3215 $\mathbf{v}^\top \nabla_y f_{t,\rho}(\mathbf{x}, \mathbf{y})$.

3216 Thus, we have

$$\begin{aligned}
 3219 \quad & \mathbb{E} [\langle \mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t), \mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t) \rangle] \\
 3220 \quad &\geq \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \mathbb{E} [\|\mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2] + \frac{1}{\mu_g + \ell_{g,1}} \mathbb{E} [\|\mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2].
 \end{aligned}$$

3223 Since $\delta_t \leq \left(2 + \frac{1}{\ell_{g,1}^2}\right) \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}$, then we have

$$\begin{aligned}
 3225 \quad I_t &\leq \left(1 - 2\delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}\right) \mathbb{E} [\|\mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2] + \frac{1}{\ell_{g,1}^2} \left(\frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}\right) \delta_t \mathbb{E} [\|\mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2] \\
 3226 \quad &\leq \left(1 - \delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}\right) \mathbb{E} [\|\mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2], \tag{182}
 \end{aligned}$$

3230 where the second inequality holds since from (119b), we have

$$\begin{aligned}
 3232 \quad \mathbb{E} [\|\mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2] &= \mathbb{E} [\|\nabla_y f_{t,\rho}(\mathbf{x}, \mathbf{y}) + \nabla_y^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) \mathbf{v}\|^2] \\
 3233 \quad &= \mathbb{E} [\|\nabla_y^2 g_{t,\rho}(\mathbf{x}, \mathbf{y}) (\mathbf{v} - \hat{\mathbf{v}}_t^*(\mathbf{x}_t))\|^2] \\
 3234 \quad &\leq \ell_{g,1}^2 \mathbb{E} [\|\mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2].
 \end{aligned}$$

3236 **Bounding K_t .** Let

$$\tilde{\nabla}_y^2 g_t(\mathbf{x}_t, \mathbf{y}_t) = \frac{1}{2\rho_\mathbf{v}} (\nabla_y g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_\mathbf{v} \mathbf{v}_t) - \nabla_y g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_\mathbf{v} \mathbf{v}_t)).$$

3240 From Lemma D.15, we have

$$K_t = \mathbb{E} [\|\hat{\mathbf{d}}_t^\mathbf{v} - \mathbf{d}_{t,\rho}^\mathbf{v}(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\|^2] \leq B_t. \tag{183}$$

Putting (182), and (183) together with Eq. (181) yields the desired result.

$$\mathbb{E} \left[\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2 \right] \leq (1 + \hat{\alpha}) \left(1 - \delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \right) \mathbb{E} \left[\|\mathbf{v}_t - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2 \right] + \left(1 + \frac{1}{\hat{\alpha}} \right) \delta_t^2 B_t.$$

□

Lemma D.17. Suppose Assumptions 3.2 and 3.3 hold. Let $\hat{\theta}_t^\mathbf{v}$ be defined in (146). Set the parameter p in (8) as $p = \frac{\ell_{f,0}}{\mu_g}$. Then, for any positive choice of step sizes satisfying

$$\delta_t \leq \left(2 + \frac{1}{\ell_{g,1}^2} \right) \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}},$$

the sequence $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ generated by Algorithm 2 guarantees the following bound:

$$\begin{aligned} \sum_{t=1}^T \left(\mathbb{E}[\hat{\theta}_{t+1}^\mathbf{v}] - \mathbb{E}[\hat{\theta}_t^\mathbf{v}] \right) &\leq \sum_{t=1}^T \left(-\frac{L_{\mu_g}}{4} \mathbb{E}[\hat{\theta}_t^\mathbf{v}] + \frac{4}{L_{\mu_g}} B_t \right) \delta_t \\ &\quad + \frac{16\nu^2}{L_{\mu_g} \mu_g^2} (2L_y^2 + 1) \sum_{t=1}^T \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \frac{1}{\delta_t} \\ &\quad + \sum_{t=1}^T \left(\frac{96\ell_{g,1}\nu^2}{L_{\mu_g} \mu_g^3} (\rho_s^2 + \rho_r^2) + \frac{48\nu^2}{L_{\mu_g} \mu_g^2} \sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{y}_{t-1}^*(\mathbf{x}) - \mathbf{y}_t^*(\mathbf{x})\|^2 \right) \frac{1}{\delta_t}, \end{aligned} \quad (184)$$

where B_t and ν are defined in Lemmas D.15 and C.7, respectively.

Proof. From Lemma B.5, we have, for any $\hat{c} > 0$

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_{t+1}^*(\mathbf{x}_{t+1})\|^2 \right] &= \mathbb{E} \left[\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_t^*(\mathbf{x}_t) + \hat{\mathbf{v}}_t^*(\mathbf{x}_t) - \hat{\mathbf{v}}_{t+1}^*(\mathbf{x}_{t+1})\|^2 \right] \\ &\leq (1 + \hat{c}) \mathbb{E} [\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2] \\ &\quad + \left(1 + \frac{1}{\hat{c}} \right) \mathbb{E} \left[\|\hat{\mathbf{v}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2 \right]. \end{aligned} \quad (185)$$

From Lemma D.16, we have, for any $\hat{\alpha} > 0$

$$\mathbb{E} \left[\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2 \right] \leq (1 + \hat{\alpha}) \left(1 - \delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \right) \hat{\theta}_t^\mathbf{v} + \left(1 + \frac{1}{\hat{\alpha}} \right) \delta_t^2 B_t. \quad (186)$$

Substituting (186) into (185), we get

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_{t+1}^*(\mathbf{x}_{t+1})\|^2 \right] &\leq (1 + \hat{c})(1 + \hat{\alpha}) \left(1 - \delta_t \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}} \right) \hat{\theta}_t^\mathbf{v} \\ &\quad + (1 + \hat{c}) \left(1 + \frac{1}{\hat{\alpha}} \right) \delta_t^2 B_t \\ &\quad + \left(1 + \frac{1}{\hat{c}} \right) \mathbb{E} \left[\|\hat{\mathbf{v}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2 \right]. \end{aligned} \quad (187)$$

Choose $\hat{c} = \frac{\delta_t L_{\mu_g}/4}{1 - \frac{\delta_t L_{\mu_g}}{2}}$ and $\hat{a} = \frac{\delta_t L_{\mu_g}/2}{1 - \delta_t L_{\mu_g}}$. Then, the following equations and inequalities are satisfied.

$$\begin{aligned} (1 + \hat{c})(1 + \hat{a})(1 - \delta_t L_{\mu_g}) &= 1 - \frac{\delta_t L_{\mu_g}}{4}, \\ (1 + \hat{c})\left(1 + \frac{1}{\hat{a}}\right) &\leq \frac{4}{\delta_t L_{\mu_g}}, \\ 1 + \frac{1}{\hat{a}} &\leq \frac{2}{\delta_t L_{\mu_g}}, \quad 1 + \frac{1}{\hat{c}} \leq \frac{4}{\delta_t L_{\mu_g}}, \end{aligned} \tag{188}$$

where $L_{\mu_g} = \frac{\mu_g \ell_{g,1}}{\mu_g + \ell_{g,1}}$.

Thus, we have

$$\begin{aligned} \mathbb{E} [\|\mathbf{v}_{t+1} - \hat{\mathbf{v}}_{t+1}^*(\mathbf{x}_{t+1})\|^2] &\leq \left(1 - \frac{\delta_t L_{\mu_g}}{4}\right) \hat{\theta}_t^\mathbf{v} + \frac{4}{L_{\mu_g}} \delta_t B_t \\ &\quad + \frac{4}{L_{\mu_g}} \frac{1}{\delta_t} \mathbb{E} [\|\hat{\mathbf{v}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2]. \end{aligned} \tag{189}$$

We now bound the last term on the right-hand side of (189). By Lemma C.7, we have:

$$\begin{aligned} &\|\hat{\mathbf{v}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2 \\ &\leq 2 \frac{\nu^2}{\mu_g^2} \left(\|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \right) \\ &\leq 2 \frac{\nu^2}{\mu_g^2} \left(2 \|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t)\|^2 \right. \\ &\quad \left. + 2 \|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \right) \\ &\leq 2 \frac{\nu^2}{\mu_g^2} \left(2 L_y^2 \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 2 \|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \right), \end{aligned} \tag{190}$$

where the last inequality follows from Lemma D.2.

From (154), we have

$$\begin{aligned} \|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2 &\leq 3 \|\hat{\mathbf{y}}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_{t+1}^*(\mathbf{x}_t)\|^2 \\ &\quad + 3 \|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + 3 \|\mathbf{y}_t^*(\mathbf{x}_t) - \hat{\mathbf{y}}_t^*(\mathbf{x}_t)\|^2 \\ &\leq 3 \|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \frac{6\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g}. \end{aligned} \tag{191}$$

Plugging (191) into (190), we get

$$\begin{aligned} &\|\hat{\mathbf{v}}_{t+1}^*(\mathbf{x}_{t+1}) - \hat{\mathbf{v}}_t^*(\mathbf{x}_t)\|^2 \\ &\leq 4 \frac{\nu^2}{\mu_g^2} (2 L_y^2 + 1) \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ &\quad + 4 \frac{\nu^2}{\mu_g^2} \left(3 \|\mathbf{y}_{t+1}^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\|^2 + \frac{6\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g} \right). \end{aligned} \tag{192}$$

Then, substituting (192) into (189), rearranging the resulting inequality and summing over $t \in [T]$, we obtain the desired result. \square

D.5. Bounds on the Zeroth-Order Objective Function and its Projected Gradients

Lemma D.18. Suppose Assumptions B2., B3. hold. Then, for the sequence $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ generated by Algorithm 2, we have

$$\begin{aligned} & \|\hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\ & \leq (12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2})d_1\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + (12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2})d_1\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\ & + \frac{9}{2}d_1\ell_{g,1}^2\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + (3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_{\mathbf{v}}^2})d_1^2\rho_{\mathbf{s}}^2, \end{aligned}$$

where $\hat{\nabla}_{\mathbf{x}} f_{t+1}$ and $\hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}$ are defined in (20b) and (21b), respectively.

Proof. From Lemma D.6, we have

$$\begin{aligned} & \|\hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1})\|^2 \\ & \leq 3d_1\ell_{g,1}^2\|\mathbf{z}_{t+1} - \mathbf{z}_t\|^2 + \frac{3}{2}\ell_{f,1}^2d_1^2\rho_{\mathbf{s}}^2 \\ & \leq 6d_1\ell_{f,1}^2\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6d_1\ell_{f,1}^2\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 + \frac{3}{2}\ell_{f,1}^2d_1^2\rho_{\mathbf{s}}^2. \end{aligned} \quad (193)$$

Moreover, from (21a), we have

$$\begin{aligned} & \|\hat{\nabla}_{\mathbf{yy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{yy}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\ & = \frac{1}{4\rho_{\mathbf{v}}^2}\|\hat{\nabla}_{\mathbf{x}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} + \rho_{\mathbf{v}}\mathbf{v}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{x}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}}\mathbf{v}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\ & \leq \frac{3}{4\rho_{\mathbf{v}}^2}d_1\ell_{g,1}^2\|(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} + \rho_{\mathbf{v}}\mathbf{v}_{t+1}) - (\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}}\mathbf{v}_t)\|^2 + \frac{3}{8\rho_{\mathbf{v}}^2}\ell_{g,1}^2d_1^2\rho_{\mathbf{s}}^2 \\ & \leq \frac{9}{4\rho_{\mathbf{v}}^2}d_1\ell_{g,1}^2\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \frac{9}{4\rho_{\mathbf{v}}^2}d_1\ell_{g,1}^2\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\ & + \frac{9}{4}d_1\ell_{g,1}^2\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + \frac{3}{8\rho_{\mathbf{v}}^2}\ell_{g,1}^2d_1^2\rho_{\mathbf{s}}^2, \end{aligned} \quad (194)$$

where the first inequality follows from Lemma D.6.

From $\|a + b\|^2 \leq 2(\|a\|^2 + \|b\|^2)$, we get

$$\begin{aligned} & \|\hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\ & \leq 2\|\hat{\nabla}_{\mathbf{yy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{yy}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\ & + 2\|\hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1})\|^2 \\ & \leq (12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2})d_1\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + (12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2})d_1\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\ & + \frac{9}{2}d_1\ell_{g,1}^2\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + (3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_{\mathbf{v}}^2})d_1^2\rho_{\mathbf{s}}^2, \end{aligned}$$

where the second inequality follows from (193) and (194). \square

Lemma D.19. Suppose Assumptions B2., B3., D2., and D4. hold. Consider the sequence $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ generated by Algorithm 2, and define

$$e_t^L := \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t) + \hat{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) - \hat{\mathbf{d}}_t^x, \quad \text{where} \quad (195)$$

$$\tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t) = \frac{1}{2\rho_{\mathbf{v}}}(\nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}}\mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}}\mathbf{v}_t)). \quad (196)$$

3410 Then, we have

$$\begin{aligned}
 3411 \quad & \mathbb{E}\|e_{t+1}^L\|^2 \leq (1 - \eta_{t+1})^2 \mathbb{E}\|e_t^L\|^2 + 36\mathbb{E}\|\nabla_{\mathbf{x}}f_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}}f_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\
 3412 \quad & + \left(18d_1^2\ell_{f,1}^2 + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_{\mathbf{v}}^2})d_1^2\right)\rho_{\mathbf{s}}^2 + 18d_1^2\ell_{g,1}^2\frac{\rho_{\mathbf{s}}^2}{\rho_{\mathbf{v}}^2} \\
 3413 \quad & + \frac{18}{\rho_{\mathbf{v}}^2}\mathbb{E}\|\nabla_{\mathbf{x}}g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}}\mathbf{v}_t) - \nabla_{\mathbf{x}}g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_{\mathbf{v}}\mathbf{v}_t)\|^2 \\
 3414 \quad & + \frac{18}{\rho_{\mathbf{v}}^2}\mathbb{E}\|\nabla_{\mathbf{x}}g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}}\mathbf{v}_t) - \nabla_{\mathbf{x}}g_t(\mathbf{x}_t, \mathbf{y}_t - \rho_{\mathbf{v}}\mathbf{v}_t)\|^2 \\
 3415 \quad & + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2})d_1\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2})d_1\mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 3416 \quad & + 27d_1\ell_{g,1}^2\mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + 3(\frac{\hat{\sigma}_{g,\mathbf{x}}^2}{b\rho_{\mathbf{v}}^2} + \frac{\hat{\sigma}_{f,\mathbf{x}}^2}{b})\eta_{t+1}^2. \\
 3417 \quad & \\
 3418 \quad & \\
 3419 \quad & \\
 3420 \quad & \\
 3421 \quad & \\
 3422 \quad & \\
 3423 \quad & \\
 3424 \quad & \\
 3425 \quad & \\
 3426 \quad & \\
 3427 \quad & \\
 3428 \quad & \\
 3429 \quad & \\
 3430 \quad & \\
 3431 \quad \textit{Proof.} According to the definition of } \hat{\mathbf{d}}_t^{\mathbf{x}} \text{ in Algorithm 2, we have} \\
 3432 \quad & \\
 3433 \quad \hat{\mathbf{d}}_{t+1}^{\mathbf{x}} - \hat{\mathbf{d}}_t^{\mathbf{x}} = -\eta_{t+1}\hat{\mathbf{d}}_t^{\mathbf{x}} + \eta_{t+1}(\hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})) \\
 3434 \quad & + (1 - \eta_{t+1})(\hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) \\
 3435 \quad & - \hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})). \\
 3436 \quad & \\
 3437 \quad & \\
 3438 \quad \text{Then, we have} \\
 3439 \quad & \\
 3440 \quad \mathbb{E}\|\nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{x}}\|^2 \\
 3441 \quad = \mathbb{E}\|\nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}) - (\hat{\mathbf{d}}_{t+1}^{\mathbf{x}} - \hat{\mathbf{d}}_t^{\mathbf{x}})\|^2 \\
 3442 \quad = \mathbb{E}\|\nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_t^{\mathbf{x}} + \eta_{t+1}\hat{\mathbf{d}}_t^{\mathbf{x}} \\
 3443 \quad - \eta_{t+1}(\hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})) \\
 3444 \quad - (1 - \eta_{t+1})(\hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) \\
 3445 \quad - \hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1}))\|^2 \\
 3446 \quad = \mathbb{E}\|(1 - \eta_{t+1})(\nabla_{\mathbf{x}}f_{t,\rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{x}}) \\
 3447 \quad + \eta_{t+1}(\nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}) - \hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})) \\
 3448 \quad + (1 - \eta_{t+1})(\nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}) - \nabla_{\mathbf{x}}f_{t,\rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{xy}}^2g_t(\mathbf{z}_t) \\
 3449 \quad + \nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_t) - \nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_t) \\
 3450 \quad - \hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) + \hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1}))\|^2. \\
 3451 \quad & \\
 3452 \quad & \\
 3453 \quad & \\
 3454 \quad & \\
 3455 \quad & \\
 3456 \quad & \\
 3457 \quad & \\
 3458 \quad \text{Since} \\
 3459 \quad & \\
 3460 \quad \mathbb{E}\left[\hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})\right] = \nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}), \\
 3461 \quad & \\
 3462 \quad \mathbb{E}\left[\hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \hat{\nabla}_{\mathbf{x}}f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\right] \\
 3463 \quad = \nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_{t+1}) - \nabla_{\mathbf{x}}f_{t+1,\rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{xy}}^2g_{t+1}(\mathbf{z}_t), \\
 3464 \quad &
 \end{aligned} \tag{197}$$

3465 then, we have

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{x}}\|^2 \\
 &= (1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_{t, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{x}}\|^2 \\
 &+ \|\eta_{t+1} (\nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})) \\
 &+ (1 - \eta_{t+1}) (\nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \nabla_{\mathbf{x}} f_{t, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{z}_t) \\
 &+ \nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t) - \nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t) \\
 &- \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) + \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1}))\|^2 \\
 &\leq (1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_{t, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{x}}\|^2 \\
 &+ 3(1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \nabla_{\mathbf{x}} f_{t, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{z}_t) \\
 &+ \nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t) - \nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t) \\
 &- \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) + \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 \\
 &+ 3\eta_{t+1}^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_{t+1}) - \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1})\|^2 \\
 &+ 3\eta_{t+1}^2 \mathbb{E} \|\tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})\|^2,
 \end{aligned} \tag{198}$$

3485 where the second inequality holds by Cauchy-Schwarz inequality.

3486 Note that for the last term on the right-hand side of (198), using (196) and (21b), we have

$$\begin{aligned}
 & \|\tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1})\|^2 \\
 &\leq 2 \left\| \frac{1}{2\rho_{\mathbf{v}}} (\nabla_{\mathbf{x}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} + \rho_{\mathbf{v}} \mathbf{v}_{t+1}) - \hat{\nabla}_{\mathbf{x}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} + \rho_{\mathbf{v}} \mathbf{v}_{t+1}; \bar{\mathcal{B}}_{t+1})) \right\|^2 \\
 &+ 2 \left\| \frac{1}{2\rho_{\mathbf{v}}} (\hat{\nabla}_{\mathbf{x}} g_{t+1}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} - \rho_{\mathbf{v}} \mathbf{v}_{t+1}; \bar{\mathcal{B}}_{t+1}) - \nabla_{\mathbf{x}} g_{t+1, \rho}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1} - \rho_{\mathbf{v}} \mathbf{v}_{t+1})) \right\|^2 \\
 &\leq \frac{\hat{\sigma}_{g_{\mathbf{x}}}^2}{b\rho_{\mathbf{v}}^2},
 \end{aligned}$$

3496 where the last inequality follows from Assumption 4.1.

3497 Then, from $\mathbb{E}\|a - \mathbb{E}[a]\|^2 = \mathbb{E}\|a\|^2 - \|\mathbb{E}[a]\|^2$ and Assumption 4.1, we have

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{x}}\|^2 \\
 &\leq (1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_{t, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{x}}\|^2 \\
 &+ 6(1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t) - \nabla_{\mathbf{x}} f_{t, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{z}_t)\|^2 \\
 &+ 6(1 - \eta_{t+1})^2 \mathbb{E} \|\hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_{t+1}; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}; \bar{\mathcal{B}}_{t+1}) \\
 &+ \hat{\nabla}_{\mathbf{x}} f_{t+1}(\mathbf{z}_t; \mathcal{B}_{t+1}) + \hat{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t; \bar{\mathcal{B}}_{t+1})\|^2 + 3\eta_{t+1}^2 \left(\frac{\hat{\sigma}_{g_{\mathbf{x}}}^2}{b\rho_{\mathbf{v}}^2} + \frac{\hat{\sigma}_{f_{\mathbf{x}}}^2}{b} \right),
 \end{aligned} \tag{199}$$

3507 Then, from Young's inequality and Lemma D.18, we have

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{x}}\|^2 \\
 &\leq (1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_{t, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{x}}\|^2 \\
 &+ 12(1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_{t+1, \rho}(\mathbf{z}_t) - \nabla_{\mathbf{x}} f_{t, \rho}(\mathbf{z}_t)\|^2 + 12(1 - \eta_{t+1})^2 \mathbb{E} \|\tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{z}_t)\|^2 \\
 &+ 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2}) d_1 \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2}) d_1 \|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 &+ 27d_1 \ell_{g,1}^2 \|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_{\mathbf{v}}^2}) d_1^2 \rho_{\mathbf{s}}^2 + 3\eta_{t+1}^2 \left(\frac{\hat{\sigma}_{g_{\mathbf{x}}}^2}{b\rho_{\mathbf{v}}^2} + \frac{\hat{\sigma}_{f_{\mathbf{x}}}^2}{b} \right).
 \end{aligned} \tag{200}$$

3520 For the third term on the right-hand side of (199), we have

$$3522 \quad \|\tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\ 3523 \quad \leq \frac{1}{2\rho_v^2} \|\nabla_{\mathbf{x}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \quad (201a)$$

$$3524 \quad + \frac{1}{2\rho_v^2} \|\nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t)\|^2. \quad (201b)$$

3528 For (201a), we get

$$3529 \quad \|\nabla_{\mathbf{x}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \\ 3530 \quad \leq 3 \|\nabla_{\mathbf{x}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \\ 3531 \quad + 3 \|\nabla_{\mathbf{x}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \\ 3532 \quad + 3 \|\nabla_{\mathbf{x}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \\ 3533 \quad \leq 3 \|\nabla_{\mathbf{x}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 + \frac{3\rho_s^2 d_1^2 \ell_{g,1}^2}{2},$$

3537 where the last inequality follows from Lemma 131.

3538 Similary, for (163b), we have

$$3540 \quad \|\nabla_{\mathbf{x}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t)\|^2 \\ 3541 \quad \leq 3 \|\nabla_{\mathbf{x}} g_t(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t)\|^2 + \frac{3\rho_s^2 d_1^2 \ell_{g,1}^2}{2}.$$

3544 Substituting these inequalities in (201), we have

$$3546 \quad \|\tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\ 3547 \quad \leq \frac{3}{2\rho_v^2} \|\nabla_{\mathbf{x}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \\ 3548 \quad + \frac{3}{2\rho_v^2} \|\nabla_{\mathbf{x}} g_t(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t)\|^2 + \frac{3\rho_s^2 d_1^2 \ell_{g,1}^2}{2\rho_v^2}. \quad (202)$$

3552 For the second term on the right-hand side of (199), we have

$$3554 \quad \|\nabla_{\mathbf{x}} f_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\ 3555 \quad \leq 3 \|\nabla_{\mathbf{x}} f_{t+1,\rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\ 3556 \quad + 3 \|\nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\ 3557 \quad + 3 \|\nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\ 3558 \quad \leq 3 \|\nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_t, \mathbf{y}_t)\|^2 + \frac{3\rho_s^2 d_1^2 \ell_{f,1}^2}{2}, \quad (203)$$

3562 where the last inequality follows from Eq. (133).

3575 From (202), (203) and (200), we get

$$\begin{aligned}
 & \mathbb{E} \|\nabla_{\mathbf{x}} f_{t+1,\rho}(\mathbf{z}_{t+1}) + \tilde{\nabla}_{\mathbf{xy}}^2 g_{t+1}(\mathbf{z}_{t+1}) - \hat{\mathbf{d}}_{t+1}^{\mathbf{x}}\|^2 \\
 & \leq (1 - \eta_{t+1})^2 \mathbb{E} \|\nabla_{\mathbf{x}} f_{t,\rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{xy}}^2 g_t(\mathbf{z}_t) - \hat{\mathbf{d}}_t^{\mathbf{x}}\|^2 \\
 & + 36 \mathbb{E} \|\nabla_{\mathbf{x}} f_t(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{x}} f_{t+1}(\mathbf{x}_t, \mathbf{y}_t)\|^2 + 18 \rho_s^2 d_1^2 \ell_{f,1}^2 \\
 & + \frac{18}{\rho_v^2} \mathbb{E} \|\nabla_{\mathbf{x}} g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \\
 & + \frac{18}{\rho_v^2} \mathbb{E} \|\nabla_{\mathbf{x}} g_t(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t) - \nabla_{\mathbf{x}} g_{t+1}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t)\|^2 + \frac{18 \rho_s^2 d_1^2 \ell_{g,1}^2}{\rho_v^2} \\
 & + 6(12 \ell_{f,1}^2 + \frac{9 \ell_{g,1}^2}{2 \rho_v^2}) d_1 \mathbb{E} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6(12 \ell_{f,1}^2 + \frac{9 \ell_{g,1}^2}{2 \rho_v^2}) d_1 \mathbb{E} \|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\
 & + 27 d_1 \ell_{g,1}^2 \mathbb{E} \|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + 6(3 \ell_{f,1}^2 + \frac{3 \ell_{g,1}^2}{4 \rho_v^2}) d_1^2 \rho_s^2 + 3 \eta_{t+1}^2 (\frac{\hat{\sigma}_{g_x}^2}{b \rho_v^2} + \frac{\hat{\sigma}_{f_x}^2}{b}).
 \end{aligned}$$

□

3591
3592
3593
3594
3595 **Lemma D.20.** Suppose Assumptions 3.2, B2., B3., and 3.4 hold. Then, for the sequence of functions $\{f_{t,\rho}\}_{t=1}^T$ defined in
3596 Eq. (17), we have

$$\begin{aligned}
 & \sum_{t=1}^T (f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - f_{t,\rho}(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1}))) \\
 & \leq 2M + V_T + \ell_{f,1} \left(1 + 2 \frac{\ell_{g,1}}{\mu_g} \right) T (\rho_s^2 + \rho_r^2).
 \end{aligned}$$

3603 Here, V_T is defined in (10); and M is defined in Assumption 3.4.

3608 *Proof.* Note that, we have

$$\begin{aligned}
 & \sum_{t=1}^T (f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - f_{t,\rho}(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1}))) \\
 & = \sum_{t=1}^T (f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t))) \tag{204}
 \end{aligned}$$

$$\begin{aligned}
 & + \sum_{t=1}^T (f_t(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1}))) \tag{205}
 \end{aligned}$$

$$\begin{aligned}
 & + \sum_{t=1}^T (f_t(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1})) - f_{t,\rho}(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1}))). \tag{206}
 \end{aligned}$$

3621 From (127), we have

$$(204) \leq T \frac{\ell_{f,1}(\rho_s^2 + \rho_r^2)}{2}, \tag{207}$$

3625 and

$$(206) \leq T \frac{\ell_{f,1}(\rho_s^2 + \rho_r^2)}{2}. \tag{208}$$

3630 Moreover, from Lemma D.7, we have

$$\begin{aligned}
 (205) &= \sum_{t=1}^T (f_t(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \\
 &\quad + \sum_{t=1}^T (f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1}))) \\
 &\quad + \sum_{t=1}^T (f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1})) - f_t(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1}))) \\
 &\leq \ell_{f,1} \sum_{t=1}^T \|\hat{\mathbf{y}}_t^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t)\| + \ell_{f,1} \sum_{t=1}^T \|\hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1}) - \mathbf{y}_t^*(\mathbf{x}_{t+1})\| \\
 &\quad + \sum_{t=1}^T (f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1}))) \\
 &\leq 2T\ell_{f,1} \frac{\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g} + \sum_{t=1}^T (f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1}))). \tag{209}
 \end{aligned}$$

3649 For the last term of the above inequality, we have

$$\begin{aligned}
 \sum_{t=1}^T (f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_t(\mathbf{x}_{t+1}, \mathbf{y}_t^*(\mathbf{x}_{t+1}))) &= f_1(\mathbf{x}_1, \mathbf{y}_1^*(\mathbf{x}_1)) - f_T(\mathbf{x}_{T+1}, \mathbf{y}_T^*(\mathbf{x}_{T+1})) \\
 &\quad + \sum_{t=2}^T (f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - f_{t-1}(\mathbf{x}_t, \mathbf{y}_{t-1}^*(\mathbf{x}_t))) \\
 &\leq 2M + V_T,
 \end{aligned}$$

3658 which implies that

$$(205) \leq 2T\ell_{f,1} \frac{\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g} + 2M + V_T. \tag{210}$$

3663 From (207), (208), and (210), we get the desired result. \square

3664
 3665
 3666 **Lemma D.21.** Suppose that Assumptions 3.2 and 3.3 hold. Let $f_{t,\rho}$ be defined as in (17). Then, for $\hat{\mathbf{d}}_t^x$ generated by
 3667 Algorithm 2, for all $t \in [T]$, we have

$$\begin{aligned}
 \mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^x - \nabla f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) \right\|^2 \right] &\leq 4\mathbb{E} \left[\|e_t^L\|^2 \right] + 4\ell_{g,2}^2 \rho_v^2 p^4 \\
 &\quad + 2M_f^2 \left(\mathbb{E}[\hat{\theta}_t^y] + \mathbb{E}[\hat{\theta}_t^v] \right) := A_t, \tag{211}
 \end{aligned}$$

3673 where e_t^L is defined in Lemma D.15, and $\hat{\theta}_t^y$, $\hat{\theta}_t^v$ are as defined in (146). Additionally, M_f is given in Lemma D.2.

3676 *Proof.* From $\|a + b\|^2 \leq 2(\|a\|^2 + \|b\|^2)$, we get

$$\begin{aligned}
 \mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^x - \nabla f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) \right\|^2 \right] &\leq 2\mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^x - \mathbf{d}_{t,\rho}^x \right\|^2 \right] \tag{212a} \\
 &\quad + 2\mathbb{E} \left[\left\| \mathbf{d}_{t,\rho}^x - \nabla f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) \right\|^2 \right], \tag{212b}
 \end{aligned}$$

3685 where $\mathbf{d}_{t,\rho}^{\mathbf{x}}$ is defined in (120c). From Lemma D.15, we have

$$3687 \quad (212a) \leq 4\mathbb{E} \left[\|e_t^L\|^2 \right] + 4\ell_{g,2}^2 \rho_{\mathbf{v}}^2 p^4. \quad (213)$$

3689 Moreover, from Eq. (123a), we get

$$3691 \quad (212b) \leq 2M_f^2 \left(\mathbb{E}[\hat{\theta}_t^{\mathbf{y}}] + \mathbb{E}[\hat{\theta}_t^{\mathbf{v}}] \right). \quad (214)$$

3693 Substituting (213) and (214) into (212), we conclude the desired result. \square

3695 **Lemma D.22.** Suppose Assumptions 3.2, 3.3, and 3.4 hold. Let the sequence of functions $\{f_{t,\rho}\}_{t=1}^T$ be defined in (17), and
 3696 $\mathcal{P}_{\mathcal{X},\alpha_t}$ be given in Definition B.1. Then, for any positive choice of step sizes as $\alpha_t \leq 1/4L_f$, for all $t \in [T]$, Algorithm 2
 3697 guarantees the following bound:

$$\begin{aligned} 3699 & \sum_{t=1}^T (\alpha_t - L_f \alpha_t^2) \mathbb{E} \left[\|\mathcal{P}_{\mathcal{X},\alpha_t}(\mathbf{x}_t; \nabla f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right] \\ 3700 & \leq 12M + 6V_T + \sum_{t=1}^T (6\alpha_t - 3L_f \alpha_t^2) A_t \\ 3701 & + \sum_{t=1}^T \left(6\ell_{f,1}(1 + 2\frac{\ell_{g,1}}{\mu_g}) + \frac{3\ell_{f,1}\ell_{g,1}}{\mu_g}(\alpha_t - L_f \alpha_t^2) \right) (\rho_{\mathbf{s}}^2 + \rho_{\mathbf{r}}^2), \end{aligned} \quad (215)$$

3708 where V_T and A_t are respectively defined in Eq. (10) and Lemma D.21.

3712 *Proof.* Due to the L_f -smoothness of f_t function by Lemma C.1, $f_{t,\rho}$ is L_f -smooth as well. Hence,

$$\begin{aligned} 3713 & f_{t,\rho}(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1})) - f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) \\ 3714 & \leq \langle \nabla f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \frac{L_f}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ 3715 & = -\alpha_t \langle \nabla f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)), \mathcal{P}_{\mathcal{X},\alpha_t}(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}}) \rangle + \frac{L_f \alpha_t^2}{2} \left\| \mathcal{P}_{\mathcal{X},\alpha_t}(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}}) \right\|^2. \end{aligned} \quad (216)$$

3719 For the first term on the R.H.S of Eq. (216), we have that

$$\begin{aligned} 3722 & -\mathbb{E} \left\langle \nabla f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)), \mathcal{P}_{\mathcal{X},\alpha_t}(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}}) \right\rangle \\ 3723 & = -\mathbb{E} \left\langle \hat{\mathbf{d}}_t^{\mathbf{x}}, \mathcal{P}_{\mathcal{X},\alpha_t}(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}}) \right\rangle \\ 3724 & - \mathbb{E} \left\langle \nabla f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - \hat{\mathbf{d}}_t^{\mathbf{x}}, \mathcal{P}_{\mathcal{X},\alpha_t}(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}}) \right\rangle \\ 3725 & \leq -\frac{1}{2} \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X},\alpha_t}(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}}) \right\|^2 \right] + \frac{1}{2} \mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^{\mathbf{x}} - \nabla f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) \right\|^2 \right] \\ 3726 & \leq -\frac{1}{2} \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X},\alpha_t}(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}}) \right\|^2 \right] + \frac{A_t}{2}, \end{aligned} \quad (217)$$

3733 where the first inequality follows from Lemma B.8; the last inequality follows from Lemma D.21.

3734 Plugging the bound (217) into (216), we have that

$$\begin{aligned} 3736 & \mathbb{E} [f_{t,\rho}(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1})) - f_{t,\rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t))] \\ 3737 & \leq \frac{(L_f \alpha_t^2 - \alpha_t)}{2} \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X},\alpha_t}(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}}) \right\|^2 \right] + \frac{\alpha_t A_t}{2}, \end{aligned}$$

3740 which can be rearranged into

$$\begin{aligned} & (\alpha_t - L_f \alpha_t^2) \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}} \right) \right\|^2 \right] \\ & \leq 2 \mathbb{E} [f_{t, \rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - f_{t, \rho}(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1}))] + \alpha_t A_t. \end{aligned} \quad (218)$$

3745 In addition, we have

$$\begin{aligned} & \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \right) \right\|^2 \right] \\ & \leq 3 \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}} \right) - \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) \right) \right\|^2 \right] \\ & + 3 \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) \right) - \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \right) \right\|^2 \right] \\ & + 3 \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}} \right) \right\|^2 \right] \\ & \leq 3 \mathbb{E} \left[\left\| \hat{\mathbf{d}}_t^{\mathbf{x}} - \nabla f_{t, \rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) \right\|^2 \right] \\ & + 3 \mathbb{E} \left[\left\| \nabla f_{t, \rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \right\|^2 \right] \\ & + 3 \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}} \right) \right\|^2 \right], \end{aligned}$$

3742 where the second inequality follows from non-expansiveness of the projection operator.

3743 Then, from Lemma D.21 and Assumption 3.3, we have

$$\begin{aligned} & \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \right) \right\|^2 \right] \\ & \leq 3A_t + 3\ell_{f,1} \mathbb{E} \left[\left\| \hat{\mathbf{y}}_t^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t) \right\|^2 \right] + 3 \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}} \right) \right\|^2 \right] \\ & \leq 3A_t + 3\ell_{f,1} \frac{\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g} + 3 \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \hat{\mathbf{d}}_t^{\mathbf{x}} \right) \right\|^2 \right], \end{aligned} \quad (219)$$

3744 where the last inequality is by Lemma D.7.

3745 Combining (218) and (219) and summing over $t = 1$ to T , we have

$$\begin{aligned} & \sum_{t=1}^T (\alpha_t - L_f \alpha_t^2) \mathbb{E} \left[\left\| \mathcal{P}_{\mathcal{X}, \alpha_t} \left(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \right) \right\|^2 \right] \\ & \leq 6 \sum_{t=1}^T (f_{t, \rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - f_{t, \rho}(\mathbf{x}_{t+1}, \hat{\mathbf{y}}_t^*(\mathbf{x}_{t+1}))) \\ & + \frac{3\ell_{f,1}\ell_{g,1}}{\mu_g} (\rho_s^2 + \rho_r^2) \sum_{t=1}^T (\alpha_t - L_f \alpha_t^2) + 3 \sum_{t=1}^T (2\alpha_t - L_f \alpha_t^2) A_t \\ & \leq 12M + 6V_T + 6\ell_{f,1} \left(1 + 2 \frac{\ell_{g,1}}{\mu_g} \right) T (\rho_s^2 + \rho_r^2) \\ & + \frac{3\ell_{f,1}\ell_{g,1}}{\mu_g} (\rho_s^2 + \rho_r^2) \sum_{t=1}^T (\alpha_t - L_f \alpha_t^2) + 3 \sum_{t=1}^T (2\alpha_t - L_f \alpha_t^2) A_t, \end{aligned}$$

3746 where the second inequality is due to Lemma D.20. \square

3747 **Lemma D.23.** Let the sequence $\{(\mathbf{x}_t, \mathbf{y}_t, \mathbf{v}_t)\}_{t=1}^T$ be generated by Algorithm 2.

3795 (a) Then, we have

$$3796 \quad \|y_{t+1} - y_t\|^2 \leq 2\beta_t^2 \|e_t^{g,\rho}\|^2 + 2\beta_t^2 \|\nabla_y g_{t,\rho}(x_t, y_t)\|^2,$$

3797 where $e_t^{g,\rho}$ is defined in (142).

3800 (b) Suppose Assumptions 3.2, B2, and B3, hold. Then, we have

$$3804 \quad \|x_{t+1} - x_t\|^2 \leq 4\alpha_t^2 \|\mathcal{P}_{\mathcal{X},\alpha_t}(x_t; \nabla f_{t,\rho}(x_t, y_t^*(x_t)))\|^2 \\ 3805 \quad + \frac{4\ell_{f,1}\ell_{g,1}\alpha_t^2(\rho_s^2 + \rho_r^2)}{\mu_g} + 2A_t\alpha_t^2, \quad (220)$$

3808 where A_t is defined in (211).

3812 (c) Suppose Assumptions B1., B2. and B3. hold. Then, we have

$$3814 \quad \|v_{t+1} - v_t\|^2 \leq 2\delta_t^2 \|e_t^M\|^2 + 3d_2^2 \ell_{f,1}^2 \delta_t^2 \rho_r^2 \\ 3815 \quad + (12\ell_{f,0}^2 + 6\ell_{g,1}^2 p^2) \delta_t^2 + 6\ell_{g,1}^2 \frac{\delta_t^2}{\rho_v^2} \hat{\theta}_t^y,$$

3818 where e_t^M and $\hat{\theta}_t^y$ are defined in (158) and (146), respectively.

3824 **Proof. For part (a):** From Algorithm 2, we have

$$3826 \quad \|y_{t+1} - y_t\|^2 = \beta_t^2 \|\hat{d}_t^y\|^2 \\ 3827 \quad \leq 2\beta_t^2 \|\hat{d}_t^y - \nabla_y g_{t,\rho}(x_t, y_t)\|^2 + 2\beta_t^2 \|\nabla_y g_{t,\rho}(x_t, y_t)\|^2 \\ 3828 \quad = 2\beta_t^2 \|e_t^{g,\rho}\|^2 + 2\beta_t^2 \|\nabla_y g_{t,\rho}(x_t, y_t)\|^2, \quad (221)$$

3831 and from (220), we get

3832 **For part (b):**

3833 From the update rule in Algorithm 2, we obtain

$$3835 \quad \|x_t - x_{t+1}\|^2 = \alpha_t^2 \left\| \mathcal{P}_{\mathcal{X},\alpha_t}(x_t; \hat{d}_t^x) \right\|^2 \\ 3836 \quad \leq 2\alpha_t^2 \left(\|\mathcal{P}_{\mathcal{X},\alpha_t}(x_t; \nabla f_{t,\rho}(x_t, y_t^*(x_t)))\|^2 \right. \\ 3837 \quad \left. + \left\| \mathcal{P}_{\mathcal{X},\alpha_t}(x_t; \hat{d}_t^x) - \mathcal{P}_{\mathcal{X},\alpha_t}(x_t; \nabla f_{t,\rho}(x_t, y_t^*(x_t))) \right\|^2 \right) \\ 3838 \quad \leq 2\alpha_t^2 \left(\|\mathcal{P}_{\mathcal{X},\alpha_t}(x_t; \nabla f_{t,\rho}(x_t, y_t^*(x_t)))\|^2 \right. \\ 3839 \quad \left. + \left\| \hat{d}_t^x - \nabla f_{t,\rho}(x_t, y_t^*(x_t)) \right\|^2 \right) \\ 3840 \quad \leq 2\alpha_t^2 \left(\|\mathcal{P}_{\mathcal{X},\alpha_t}(x_t; \nabla f_{t,\rho}(x_t, y_t^*(x_t)))\|^2 + A_t \right), \quad (222)$$

3847 where the first inequality is by $(a+b)^2 \leq 2a^2 + 2b^2$; the second inequality follows from non-expansiveness of the projection
3848 operator; and the last inequality follows from Lemma D.21.
3849

3850 The first term in the above inequality can be bounded as
 3851

$$\begin{aligned}
 & \| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t))) \|^2 \\
 & \leq 2 \| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t))) - \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 \\
 & + 2 \| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 \\
 & \leq 2 \| \nabla f_{t, \rho}(\mathbf{x}_t, \hat{\mathbf{y}}_t^*(\mathbf{x}_t)) - \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \|^2 \\
 & + 2 \| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 \\
 & \leq 2\ell_{f,1} \| \hat{\mathbf{y}}_t^*(\mathbf{x}_t) - \mathbf{y}_t^*(\mathbf{x}_t) \|^2 + 2 \| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 \\
 & \leq 2\ell_{f,1} \frac{\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g} + 2 \| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2,
 \end{aligned} \tag{223}$$

3862 where the last inequality follows from Lemma D.7.
 3863

3864 Based on (223) and (222), we get
 3865

$$\| \mathbf{x}_t - \mathbf{x}_{t+1} \|^2 \leq 2\alpha_t^2 \left(2 \| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 + \frac{2\ell_{f,1}\ell_{g,1}(\rho_s^2 + \rho_r^2)}{\mu_g} + A_t \right).$$

3866 **For part (c):** Note that, we have
 3867

$$\begin{aligned}
 \| \mathbf{v}_{t+1} - \mathbf{v}_t \|^2 &= \delta_t^2 \| \hat{\mathbf{d}}_t^\mathbf{v} \|^2 \\
 &\leq 2\delta_t^2 \| \hat{\mathbf{d}}_t^\mathbf{v} - \nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{z}_t) - \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t) \|^2 + 2\delta_t^2 \| \nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{z}_t) + \tilde{\nabla}_{\mathbf{y}}^2 g_t(\mathbf{z}_t) \|^2 \\
 &= 2\delta_t^2 \| e_t^M \|^2 \\
 &+ 2\delta_t^2 \| \nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) + \frac{1}{2\rho_v} (\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t)) \|^2 \\
 &\leq 2\delta_t^2 \| e_t^M \|^2 + 6\delta_t^2 \| \nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) \|^2 \\
 &+ \frac{3\delta_t^2}{2\rho_v^2} \| \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) \|^2 + \frac{3\delta_t^2}{2\rho_v^2} \| \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t) \|^2.
 \end{aligned} \tag{224}$$

3880 From Assumption B3., Lemma B.3 and (8), we have
 3881

$$\begin{aligned}
 \| \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) \|^2 &\leq \ell_{g,1}^2 \| \mathbf{y}_t + \rho_v \mathbf{v}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t) \|^2 \\
 &\leq 2\ell_{g,1}^2 \| \rho_v \mathbf{v}_t \|^2 + 2\ell_{g,1}^2 \| \mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t) \|^2 \\
 &\leq 2\ell_{g,1}^2 \rho_v^2 p^2 + 2\ell_{g,1}^2 \| \mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t) \|^2.
 \end{aligned} \tag{225}$$

3887 Similarly, we get
 3888

$$\| \nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t - \rho_v \mathbf{v}_t) \|^2 \leq 2\ell_{g,1}^2 \rho_v^2 p^2 + 2\ell_{g,1}^2 \| \mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t) \|^2. \tag{226}$$

3890 Moreover, from Eq. (133) and Assumption B1., we have
 3891

$$\begin{aligned}
 \| \nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) \|^2 &\leq 2 \| \nabla_{\mathbf{y}} f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t) \|^2 + 2 \| \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t) \|^2 \\
 &\leq \frac{d_2^2 \ell_{f,1}^2 \rho_r^2}{2} + 2 \| \nabla_{\mathbf{y}} f_t(\mathbf{x}_t, \mathbf{y}_t) \|^2 \\
 &\leq \frac{d_2^2 \ell_{f,1}^2 \rho_r^2}{2} + 2\ell_{f,0}^2.
 \end{aligned} \tag{227}$$

3898 Substituting (225), (226) and (227), into (224), we get
 3899

$$\begin{aligned}
 \| \mathbf{v}_{t+1} - \mathbf{v}_t \|^2 &\leq 2\delta_t^2 \| e_t^M \|^2 + 3d_2^2 \ell_{f,1}^2 \delta_t^2 \rho_r^2 \\
 &+ (12\ell_{f,0}^2 + 6\ell_{g,1}^2 p^2) \delta_t^2 + \frac{6\ell_{g,1}^2}{\rho_v^2} \delta_t^2 \| \mathbf{y}_t - \hat{\mathbf{y}}_t^*(\mathbf{x}_t) \|^2.
 \end{aligned}$$

□

D.6. Proof of Theorem 4.2

Proof. Since $(1 - \gamma_{t+1})^2 \leq 1 - \gamma_{t+1}$ and $\gamma_{t+1} = c_\gamma \alpha_t$, from (143), we have

$$\begin{aligned} \mathbb{E}\|e_{t+1}^{g_\rho}\|^2 - \mathbb{E}\|e_t^{g_\rho}\|^2 &\leq -c_\gamma \alpha_t \mathbb{E}\|e_t^{g_\rho}\|^2 \\ &+ 12(1 - \gamma_{t+1})^2 \mathbb{E}\|\nabla_y g_{t-1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\ &+ 9d_2^2 \ell_{g,1}^2 (1 - \gamma_{t+1})^2 \rho_r^2 + 24d_2 \ell_{g,1}^2 (1 - \gamma_{t+1})^2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ &+ 24d_2 \ell_{g,1}^2 (1 - \gamma_{t+1})^2 \mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 + 2\frac{\hat{\sigma}_{g_y}^2}{b} \gamma_{t+1}^2. \end{aligned} \quad (228)$$

Since $(1 - \eta_{t+1})^2 \leq 1 - \eta_{t+1}$ and $\eta_{t+1} = c_\eta \alpha_t$, from (197), we have

$$\begin{aligned} \mathbb{E}\|e_{t+1}^L\|^2 - \mathbb{E}\|e_t^L\|^2 &\leq -c_\eta \alpha_t \mathbb{E}\|e_t^L\|^2 + 36\mathbb{E}\|\nabla_x f_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_x f_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\ &+ \left(18d_1^2 \ell_{f,1}^2 + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_v^2})d_1^2\right) \rho_s^2 + 18d_1^2 \ell_{g,1}^2 \frac{\rho_s^2}{\rho_v^2} \\ &+ \frac{36}{\rho_v^2} \mathbb{E}\|\nabla_x g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_x g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \\ &+ 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_v^2})d_1 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_v^2})d_1 \mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\ &+ 27\ell_{g,1}^2 d_1 \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + 3(\frac{\hat{\sigma}_{g_x}^2}{b\rho_v^2} + \frac{\hat{\sigma}_{f_x}^2}{b})\eta_{t+1}^2. \end{aligned} \quad (229)$$

Since $(1 - \lambda_{t+1})^2 \leq 1 - \lambda_{t+1}$ and $\lambda_{t+1} = c_\lambda \alpha_t$, from (160), we have

$$\begin{aligned} \mathbb{E}\|e_{t+1}^M\|^2 - \mathbb{E}\|e_t^M\|^2 &\leq -c_\lambda \alpha_t \mathbb{E}\|e_t^M\|^2 + 36\mathbb{E}\|\nabla_y f_{t+1}(\mathbf{x}_t, \mathbf{y}_t) - \nabla_y f_t(\mathbf{x}_t, \mathbf{y}_t)\|^2 \\ &+ \left(18d_2^2 \ell_{f,1}^2 + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_v^2})d_2^2\right) \rho_r^2 + 18d_2^2 \ell_{g,1}^2 \frac{\rho_r^2}{\rho_v^2} \\ &+ \frac{36}{\rho_v^2} \mathbb{E}\|\nabla_y g_{t+1}(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t) - \nabla_y g_t(\mathbf{x}_t, \mathbf{y}_t + \rho_v \mathbf{v}_t)\|^2 \\ &+ 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_v^2})d_2 \mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_v^2})d_2 \mathbb{E}\|\mathbf{y}_{t+1} - \mathbf{y}_t\|^2 \\ &+ 27d_2 \ell_{g,1}^2 \mathbb{E}\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + 3(\frac{\hat{\sigma}_{g_y}^2}{b\rho_v^2} + \frac{\hat{\sigma}_{f_y}^2}{b})\lambda_{t+1}^2. \end{aligned} \quad (230)$$

Combining the outcomes .

Let

$$\begin{aligned} \Lambda := \Gamma \sum_{t=1}^T (\mathbb{E}[\hat{\theta}_{t+1}^y] - \mathbb{E}[\hat{\theta}_t^y]) + \Upsilon \sum_{t=1}^T (\mathbb{E}[\hat{\theta}_{t+1}^v] - \mathbb{E}[\hat{\theta}_t^v]) + \frac{1}{\Phi} \sum_{t=1}^T (\mathbb{E}\|e_{t+1}^{g_\rho}\|^2 - \mathbb{E}\|e_t^{g_\rho}\|^2) \\ + \frac{1}{\Psi} \sum_{t=1}^T (\mathbb{E}\|e_{t+1}^M\|^2 - \mathbb{E}\|e_t^M\|^2) + \frac{1}{\Omega} \sum_{t=1}^T (\mathbb{E}\|e_{t+1}^L\|^2 - \mathbb{E}\|e_t^L\|^2). \end{aligned}$$

3960 Here, we have

$$\begin{aligned}
 \Gamma &= \frac{11M_f^2}{L_{\mu_g}c_\beta}, \quad \Upsilon = \frac{52M_f^2}{L_{\mu_g}c_\delta}, \quad \Phi = \max \left\{ 240 \frac{d_2\ell_{g,1}^2}{L_f}, \frac{12d_2\ell_{g,1}^2 L_{\mu_g}^2 c_\beta^2}{L_f M_f^2} \right\}, \\
 \Psi &= \max \left\{ 720 \frac{d_2\ell_{f,1}^2}{L_f}, 27 \frac{L_{\mu_g}\ell_{g,1}^2 d_2 c_\delta}{\Upsilon L_f}, \frac{144d_2\ell_{f,1}^2(\mu_g + \ell_{g,1})c_\beta}{L_f \Gamma}, \frac{36\ell_{f,1}^2 d_2 L_{\mu_g}^2 c_\beta^2}{L_f M_f^2} \right\}, \\
 \Omega &= \max \left\{ 720 \frac{d_1\ell_{f,1}^2}{L_f}, 27 \frac{L_{\mu_g}\ell_{g,1}^2 d_1 c_\delta}{\Upsilon L_f}, \frac{144d_1\ell_{f,1}^2(\mu_g + \ell_{g,1})c_\beta}{L_f \Gamma}, \frac{36\ell_{f,1}^2 d_1 L_{\mu_g}^2 c_\beta^2}{L_f M_f^2} \right\},
 \end{aligned} \tag{231}$$

3970 with

$$\begin{aligned}
 c_\beta &\geq \sqrt{1760 \frac{L_y^2 M_f^2}{L_{\mu_g}^2}}, \quad c_\delta \geq \sqrt{33280 \frac{\nu^2 M_f^2}{L_{\mu_g}^2 \mu_g^2} (1 + 2L_y^2)}, \\
 c &\geq \left(\max \left\{ 4L_f, c_\beta(\mu_g + \ell_{g,1}), \frac{48L_{\mu_g}^2 d_2 \ell_{g,1}^2 c_\beta^2}{M_f^2 \Phi} \right\} \right)^3 + 1, \\
 c_v &= \max \left\{ 1080\ell_{g,1}^2, \frac{324}{M_f^2} \ell_{g,1}^4 c_\delta^2, \frac{54L_{\mu_g}^2}{M_f^2} \ell_{g,1}^2 c_\beta^2, \frac{216}{\Gamma} \ell_{g,1}^2 c_\beta (\mu_g + \ell_{g,1}) \right\} \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right), \\
 c_\gamma &= \frac{26M_f^2 \Phi}{L_{\mu_g}^2}, \quad c_\eta = 26\Omega, \quad c_\lambda = \frac{10\Upsilon}{L_{\mu_g}} c_\delta \Psi.
 \end{aligned} \tag{232}$$

3984
3985
3986
3987
3988
3989
3990
3991
3992
3993
3994
3995
3996
3997
3998
3999
4000
4001
4002
4003
4004
4005
4006
4007
4008
4009
4010
4011
4012
4013
4014

4015 By adding (229), (228), (230), (147), and (184), along with (215) and considering the fact that α_t decreases with respect to
 4016 t , and by applying Lemma D.23, we obtain:

$$\begin{aligned} 4018 & \sum_{t=1}^T A(\alpha_t, \beta_t, \delta_t) \mathbb{E} \left[\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right] + \Lambda \\ 4019 & \leq 12M + 6V_T + \sum_{t=1}^T B(\alpha_t, \beta_t, \delta_t) \hat{\theta}_t^{\mathbf{y}} + \sum_{t=1}^T C(\alpha_t, \beta_t, \delta_t) \hat{\theta}_t^{\mathbf{y}} \\ 4020 & \end{aligned} \quad (233a)$$

$$4021 + \frac{4\ell_{f,1}\ell_{g,1}}{\mu_g} \sum_{t=1}^T E(\beta_t, \delta_t, \rho_{\mathbf{v}}) \alpha_t^2 (\rho_{\mathbf{s}}^2 + \rho_{\mathbf{r}}^2) + \sum_{t=1}^T L(\alpha_t, \beta_t, \delta_t) \mathbb{E} \|e_t^L\|^2 \quad (233b)$$

$$4022 + \frac{8\ell_{g,2}^2 p^4 \Upsilon}{L_{\mu_g}} \sum_{t=1}^T \delta_t \rho_{\mathbf{v}}^2 + 4\ell_{g,2}^2 p^4 \sum_{t=1}^T (6\alpha_t - 3L_f \alpha_t^2 + 2\alpha_t^2 E(\beta_t, \delta_t, \rho_{\mathbf{v}})) \rho_{\mathbf{v}}^2 \quad (233c)$$

$$4023 + \left(\frac{12}{L_{\mu_g}} \frac{\Gamma}{\beta_T} + \frac{48\nu^2}{L_{\mu_g} \mu_g^2} \frac{\Upsilon}{\delta_T} \right) H_{2,T} + \sum_{t=1}^T M(\alpha_t, \beta_t, \delta_t) \mathbb{E} \|e_t^M\|^2 \quad (233d)$$

$$4024 + \sum_{t=1}^T Q(\alpha_t, \beta_t, \delta_t) \mathbb{E} \|e_t^{g_{\rho}}\|^2 + \sum_{t=1}^T S(\alpha_t, \beta_t, \delta_t) \mathbb{E} [\|\nabla_{\mathbf{y}} g_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2] \quad (233e)$$

$$4025 + \sum_{t=1}^T Z(\alpha_t) (3d_2^2 \ell_{f,1}^2 \delta_t^2 \rho_{\mathbf{r}}^2 + (12\ell_{f,0}^2 + 6\ell_{g,1}^2 p^2) \delta_t^2) \quad (233f)$$

$$4026 + \frac{36}{\Psi} D_{\mathbf{y}, T} + \frac{36}{\Omega} D_{\mathbf{x}, T} + \frac{12}{\Phi} G_{\mathbf{y}, T} + \frac{18}{\Psi \rho_{\mathbf{v}}^2} G_{\mathbf{v}, T} + \frac{18}{\Omega \rho_{\mathbf{v}}^2} G_{\mathbf{x}, T} \quad (233g)$$

$$4027 + 2 \sum_{t=1}^T \frac{\gamma_{t+1}^2}{\Phi} \frac{\hat{\sigma}_{g_{\mathbf{y}}}^2}{b} + 3 \left(\frac{\hat{\sigma}_{g_{\mathbf{y}}}^2}{b \rho_{\mathbf{v}}^2} + \frac{\hat{\sigma}_{f_{\mathbf{y}}}^2}{b} \right) \sum_{t=1}^{T+1} \frac{\lambda_{t+1}^2}{\Psi} + 3 \left(\frac{\hat{\sigma}_{g_{\mathbf{x}}}^2}{b \rho_{\mathbf{v}}^2} + \frac{\hat{\sigma}_{f_{\mathbf{x}}}^2}{b} \right) \sum_{t=1}^T \frac{\eta_{t+1}^2}{\Omega} \quad (233h)$$

$$4028 + R(\rho_{\mathbf{v}}) T \rho_{\mathbf{r}}^2 + \tilde{R}(\rho_{\mathbf{v}}) T \rho_{\mathbf{s}}^2 + 18d_1^2 \ell_{g,1}^2 \frac{T \rho_{\mathbf{s}}^2}{\Omega \rho_{\mathbf{v}}^2} + 18d_2^2 \ell_{g,1}^2 \frac{T \rho_{\mathbf{r}}^2}{\Psi \rho_{\mathbf{v}}^2} + \sum_{t=1}^T D(\alpha_t, \beta_t, \delta_t) (\rho_{\mathbf{s}}^2 + \rho_{\mathbf{r}}^2). \quad (233i)$$

4046 Here

$$\begin{aligned} 4047 & E(\beta_t, \delta_t, \rho_{\mathbf{v}}) := \frac{4L_{\mathbf{y}}^2}{L_{\mu_g}} \frac{\Gamma}{\beta_t} + \frac{16\nu^2}{L_{\mu_g} \mu_g^2} (2L_{\mathbf{y}}^2 + 1) \frac{\Upsilon}{\delta_t} \\ 4048 & + 24d_2 \frac{\ell_{g,1}^2}{\Phi} + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_{\mathbf{v}}^2}) (\frac{d_2}{\Psi} + \frac{d_1}{\Omega}), \\ 4049 & A(\alpha_t, \beta_t, \delta_t) := \alpha_t - L_f \alpha_t^2 - 4E(\beta_t, \delta_t, \rho_{\mathbf{v}}) \alpha_t^2, \\ 4050 & B(\alpha_t, \beta_t, \delta_t) := -\frac{L_{\mu_g}}{4} \Upsilon \delta_t + 2M_f^2 (6\alpha_t - 3L_f \alpha_t^2 + 2\alpha_t^2 E(\beta_t, \delta_t, \rho_{\mathbf{v}})), \\ 4051 & Z(\alpha_t) := 27\ell_{g,1}^2 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right), \\ 4052 & C(\alpha_t, \beta_t, \delta_t) := -\frac{L_{\mu_g}}{2} \Gamma \beta_t + Z(\alpha_t) 6\ell_{g,1}^2 \frac{\delta_t^2}{\rho_{\mathbf{v}}^2} + 2M_f^2 (6\alpha_t - 3L_f \alpha_t^2 + 2\alpha_t^2 E(\beta_t, \delta_t, \rho_{\mathbf{v}})). \end{aligned} \quad (234)$$

4070 Moreover,

$$\begin{aligned}
 M(\alpha_t, \beta_t, \delta_t) &:= -\frac{\lambda_{t+1}}{\Psi} + Z(\alpha_t)2\delta_t^2 + \frac{8\Upsilon}{L_{\mu_g}}\delta_t, \\
 D(\alpha_t, \beta_t, \delta_t) &:= 6\ell_{f,1}(1 + 2\frac{\ell_{g,1}}{\mu_g}) + \frac{3\ell_{f,1}\ell_{g,1}}{\mu_g}(\alpha_t - L_f\alpha_t^2) \\
 &\quad + \frac{24\ell_{g,1}}{L_{\mu_g}\mu_g}\frac{\Gamma}{\beta_t} + \frac{96\ell_{g,1}\nu^2}{L_{\mu_g}\mu_g^3}\frac{\Upsilon}{\delta_t}, \\
 F(\alpha_t) &:= 24d_2\frac{\ell_{g,1}^2}{\Phi} + (72\ell_{f,1}^2 + \frac{27\ell_{g,1}^2}{\rho_v^2})(\frac{d_2}{\Psi} + \frac{d_1}{\Omega}), \\
 S(\alpha_t, \beta_t, \delta_t) &:= -\frac{2\beta_t\Gamma}{\mu_g + \ell_{g,1}} + \beta_t^2\Gamma + 2F(\alpha_t)\beta_t^2, \\
 Q(\alpha_t, \beta_t, \delta_t) &:= \frac{2}{L_{\mu_g}}\Gamma\beta_t - \frac{\gamma_{t+1}}{\Phi} + 2F(\alpha_t)\beta_t^2, \\
 R(\rho_v) &:= 9d_2\frac{\ell_{g,1}^2}{\Phi} + 18d_2^2\frac{\ell_{f,1}^2}{\Psi} + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_v^2})\frac{d_2^2}{\Psi}, \\
 \acute{R}(\rho_v) &:= 18d_1^2\frac{\ell_{f,1}^2}{\Omega} + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_v^2})\frac{d_1^2}{\Omega}, \\
 L(\alpha_t, \beta_t, \delta_t) &:= -\frac{\eta_{t+1}}{\Omega} + 4(6\alpha_t - 3L_f\alpha_t^2 + 2\alpha_t^2E(\beta_t, \delta_t, \rho_v)).
 \end{aligned} \tag{235}$$

4093 We then provide bounds for the terms in (233a)-(233i).

4094 Note that, we have

$$\begin{aligned}
 E(\beta_t, \delta_t, \rho_v) &:= \frac{4L_y^2}{L_{\mu_g}}\frac{\Gamma}{\beta_t} + \frac{16\nu^2}{L_{\mu_g}\mu_g^2}(2L_y^2 + 1)\frac{\Upsilon}{\delta_t} \\
 &\quad + 24d_2\frac{\ell_{g,1}^2}{\Phi} + 6(12\ell_{f,1}^2 + \frac{9\ell_{g,1}^2}{2\rho_v^2})(\frac{d_2}{\Psi} + \frac{d_1}{\Omega}),
 \end{aligned}$$

4100 which together with $\beta_t = c_\beta\alpha_t$, $\delta_t = c_\delta\alpha_t$, we have

$$\begin{aligned}
 \alpha_t^2E(\beta_t, \delta_t, \rho_v) &= \frac{4L_y^2}{L_{\mu_g}}\frac{\Gamma\alpha_t^2}{\beta_t} + \frac{16\nu^2}{L_{\mu_g}\mu_g^2}(2L_y^2 + 1)\frac{\Upsilon\alpha_t^2}{\delta_t} \\
 &\quad + 24d_2\frac{\ell_{g,1}^2}{\Phi}\alpha_t^2 + (72\ell_{f,1}^2\alpha_t^2 + \frac{27\ell_{g,1}^2}{\rho_v^2}\alpha_t^2)(\frac{d_2}{\Psi} + \frac{d_1}{\Omega}) \\
 &\leq \frac{44L_y^2}{L_{\mu_g}^2}M_f^2\frac{\alpha_t}{c_\beta^2} + \frac{832\nu^2}{L_{\mu_g}^2\mu_g^2}(1 + 2L_y^2)M_f^2\frac{\alpha_t}{c_\delta^2} \\
 &\quad + 6\frac{d_2\ell_{g,1}^2}{L_f\Phi}\alpha_t + (\frac{18\ell_{f,1}^2}{L_f}\alpha_t + \frac{27\ell_{g,1}^2}{c_v}\alpha_t)(\frac{d_2}{\Psi} + \frac{d_1}{\Omega}) \\
 &\leq \frac{\alpha_t}{8},
 \end{aligned} \tag{236}$$

4115 where the first inequality is by $\Gamma = \frac{11M_f^2}{L_{\mu_g}c_\beta}$, $\Upsilon = \frac{52M_f^2}{L_{\mu_g}c_\delta}$ in (231), $\rho_v^2 = c_v\alpha_t$ and $\alpha_t \leq 1/4L_f$; the second inequality follows
 4116 from $c_\beta \geq \sqrt{1760\frac{L_y^2M_f^2}{L_{\mu_g}^2}}$, $c_\delta \geq \sqrt{33280\frac{\nu^2M_f^2}{L_{\mu_g}^2\mu_g^2}(1 + 2L_y^2)}$, in (232); and $\Phi = 240\frac{d_2\ell_{g,1}^2}{L_f}$, $\Psi = 720\frac{d_2\ell_{f,1}^2}{L_f}$, $\Omega = 720\frac{d_1\ell_{f,1}^2}{L_f}$
 4117 and $c_v \geq 1080\ell_{g,1}^2(\frac{d_2}{\Psi} + \frac{d_1}{\Omega})$ in (231).

4125 Moreover, we have

$$\begin{aligned}
 A(\alpha_t, \beta_t, \delta_t) &= \alpha_t - L_f \alpha_t^2 - 4E(\beta_t, \delta_t, \rho_v) \alpha_t^2 \\
 &\geq \alpha_t - L_f \alpha_t^2 - \frac{\alpha_t}{2} \\
 &\geq \frac{\alpha_t}{4},
 \end{aligned} \tag{237}$$

4132 where the last inequality is by $\alpha_t \leq 1/4L_f$ in (232).

4133 **Bounding (233a)**.

4134 From $\delta_t = c_\delta \alpha_t$, we have

$$\begin{aligned}
 B(\alpha_t, \beta_t, \delta_t) &= -\frac{L_{\mu_g}}{4} \Upsilon \delta_t + 2M_f^2 (6\alpha_t - 3L_f \alpha_t^2 + 2\alpha_t^2 E(\beta_t, \delta_t, \rho_v)) \\
 &\leq -\frac{L_{\mu_g}}{4} \Upsilon c_\delta \alpha_t + 12M_f^2 \alpha_t - 6M_f^2 L_f \alpha_t^2 + \frac{M_f^2}{2} \alpha_t \\
 &\leq \left(-\frac{L_{\mu_g}}{4} \Upsilon c_\delta + \frac{25}{2} M_f^2 \right) \alpha_t \\
 &\leq -\frac{1}{2} M_f^2 \alpha_t,
 \end{aligned} \tag{238}$$

4144 where the first inequality follows from (236); the last inequality is by $\Upsilon = \frac{52M_f^2}{L_{\mu_g} c_\delta}$ in (231).

4145 From (234), we obtain

$$Z(\alpha_t) = 27\ell_{g,1}^2 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right).$$

4150 Thus, from $\beta_t = c_\beta \alpha_t$, $\delta_t = c_\delta \alpha_t$ and $\rho_v^2 = c_v \alpha_t$, we have

$$\begin{aligned}
 C(\alpha_t, \beta_t, \delta_t) &= -\frac{L_{\mu_g}}{2} \Gamma \beta_t + 162 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) \ell_{g,1}^4 \frac{\delta_t^2}{\rho_v^2} \\
 &\quad + 2M_f^2 (6\alpha_t - 3L_f \alpha_t^2 + 2\alpha_t^2 E(\beta_t, \delta_t, \rho_v)) \\
 &\leq -\frac{L_{\mu_g}}{2} \Gamma c_\beta \alpha_t + 162 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) \ell_{g,1}^4 \frac{c_\delta^2}{c_v} \alpha_t + \frac{9}{2} M_f^2 \alpha_t \\
 &= -\frac{11}{2} M_f^2 \alpha_t + 162 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) \ell_{g,1}^4 \frac{c_\delta^2}{c_v} \alpha_t + \frac{9}{2} M_f^2 \alpha_t \\
 &\leq -\frac{1}{2} M_f^2 \alpha_t,
 \end{aligned} \tag{239}$$

4163 where the first inequality follows from (236); the second equality follows from $\Gamma = \frac{11M_f^2}{L_{\mu_g} c_\beta}$ in (231); the last inequality is by

4164 $c_v \geq \frac{324}{M_f^2} \ell_{g,1}^4 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) c_\delta^2$.

4165 Thus, from (238) and (239), we get

$$(233a) \leq \mathcal{O}(V_T). \tag{240}$$

4170 **Bounding (233b)**.

4180 From (236), we also obtain

$$\begin{aligned}
 & \frac{4\ell_{f,1}\ell_{g,1}}{\mu_g} \sum_{t=1}^T E(\beta_t, \delta_t, \rho_{\mathbf{v}}) \alpha_t^2 (\rho_{\mathbf{s}}^2 + \rho_{\mathbf{r}}^2) \\
 & \leq \frac{4\ell_{f,1}\ell_{g,1}}{\mu_g} \sum_{t=1}^T \frac{\alpha_t}{8} (\rho_{\mathbf{s}}^2 + \rho_{\mathbf{r}}^2) \\
 & = \mathcal{O} \left(\sum_{t=1}^T \alpha_t (\rho_{\mathbf{s}}^2 + \rho_{\mathbf{r}}^2) \right). \tag{241}
 \end{aligned}$$

4191 From (235) and $\eta_{t+1} = c_\eta \alpha_t$, we have

$$\begin{aligned}
 L(\alpha_t, \beta_t, \delta_t) &= -\frac{\eta_{t+1}}{\Omega} + 4(6\alpha_t - 3L_f \alpha_t^2 + 2\alpha_t^2 E(\beta_t, \delta_t, \rho_{\mathbf{v}})) \\
 &\leq -\frac{c_\eta}{\Omega} \alpha_t + 25\alpha_t \\
 &\leq -\alpha_t,
 \end{aligned}$$

4198 where the last inequality is by $c_\eta \geq 26\Omega$ and (236).

4199 Thus, we get

$$\sum_{t=1}^T L(\alpha_t, \beta_t, \delta_t) \mathbb{E} \|e_t^L\|^2 \leq 0. \tag{242}$$

4204 From (242) and (241), we have

$$(233b) \leq \mathcal{O} \left(\sum_{t=1}^T \alpha_t (\rho_{\mathbf{s}}^2 + \rho_{\mathbf{r}}^2) \right). \tag{243}$$

4209 **Bounding (233c).**

4210 From $\delta_t = c_\delta \alpha_t$ and (236), we have

$$\begin{aligned}
 & \frac{8\ell_{g,2}^2 p^4 \Upsilon}{L_{\mu_g}} \sum_{t=1}^T \delta_t \rho_{\mathbf{v}}^2 + 4\ell_{g,2}^2 p^4 \sum_{t=1}^T (6\alpha_t - 3L_f \alpha_t^2 + 2\alpha_t^2 E(\beta_t, \delta_t, \rho_{\mathbf{v}})) \rho_{\mathbf{v}}^2 \\
 & \leq \frac{8\ell_{g,2}^2 p^4 \Upsilon}{L_{\mu_g}} \sum_{t=1}^T c_\delta \alpha_t \rho_{\mathbf{v}}^2 + 4\ell_{g,2}^2 p^4 \sum_{t=1}^T \frac{25}{4} \alpha_t \rho_{\mathbf{v}}^2.
 \end{aligned}$$

4217 Thus, from $\rho_{\mathbf{v}}^2 = c_{\mathbf{v}} \alpha_t$, we have

$$(233c) \leq \mathcal{O} \left(\sum_{t=1}^T \alpha_t^2 \right). \tag{244}$$

4223 **Bounding (233d).**

4224 From (234), we obtain

$$Z(\alpha_t) = 27\ell_{g,1}^2 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right). \tag{245}$$

4235 From (235), $\lambda_{t+1} = c_\lambda \alpha_t$ and $\delta_t = c_\delta \alpha_t$, we have

$$\begin{aligned} M(\alpha_t, \beta_t, \delta_t) &= -\frac{\lambda_{t+1}}{\Psi} + Z(\alpha_t) 2\delta_t^2 + \frac{8\Upsilon}{L_{\mu_g}} \delta_t \\ &= -\frac{c_\lambda \alpha_t}{\Psi} + 27\ell_{g,1}^2 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) 2c_\delta^2 \alpha_t^2 + \frac{8\Upsilon}{L_{\mu_g}} c_\delta \alpha_t \\ &\leq -\frac{2\Upsilon}{L_{\mu_g}} c_\delta \alpha_t + \frac{27}{4L_f} \ell_{g,1}^2 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) 2c_\delta^2 \alpha_t \\ &\leq -\frac{\Upsilon}{L_{\mu_g}} c_\delta \alpha_t, \end{aligned}$$

4246 where the first inequality is by $c_\lambda \geq \frac{10\Upsilon}{L_{\mu_g}} c_\delta \Psi$ and $\alpha_t \leq 1/4L_f$; the last inequality follows from $\Psi \geq 27 \frac{L_{\mu_g}}{\Upsilon L_f} \ell_{g,1}^2 d_2 c_\delta$ and
 4247 $\Omega \geq 27 \frac{L_{\mu_g}}{\Upsilon L_f} \ell_{g,1}^2 d_1 c_\delta$.

4248 Since $\beta_t = c_\beta \alpha_t$ and $\delta_t = c_\delta \alpha_t$, we get

$$\begin{aligned} (233d) &= \left(\frac{12}{L_{\mu_g}} \frac{\Gamma}{\beta_T} + \frac{48\nu^2}{L_{\mu_g} \mu_g^2} \frac{\Upsilon}{\delta_T} \right) H_{2,T} + \sum_{t=1}^T M(\alpha_t, \beta_t, \delta_t) \mathbb{E} \|e_t^M\|^2 \\ &\leq \mathcal{O}\left(\frac{H_{2,T}}{\alpha_T}\right). \end{aligned} \tag{246}$$

Bounding (233e).

4258 From (235), we have

$$F(\alpha_t) = 24d_2 \frac{\ell_{g,1}^2}{\Phi} + (72\ell_{f,1}^2 + \frac{27\ell_{g,1}^2}{\rho_v^2}) \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) \tag{247}$$

4263 From (235), $\gamma_{t+1} = c_\gamma \alpha_t$, $\beta_t = c_\beta \alpha_t$, we have

$$\begin{aligned} Q(\alpha_t, \beta_t, \delta_t) &= -\frac{\gamma_{t+1}}{\Phi} + \frac{2}{L_{\mu_g}} \Gamma \beta_t + 2F(\alpha_t) \beta_t^2 \\ &= -\frac{c_\gamma \alpha_t}{\Phi} + \frac{22M_f^2}{L_{\mu_g}^2} \alpha_t + \left(24d_2 \frac{\ell_{g,1}^2}{\Phi} + (72\ell_{f,1}^2 + \frac{27\ell_{g,1}^2}{c_v \alpha_t}) \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) \right) 2c_\beta^2 \alpha_t^2 \\ &\leq -\frac{4M_f^2}{L_{\mu_g}^2} \alpha_t + \left(24d_2 \frac{\ell_{g,1}^2}{\Phi} \alpha_t^2 + (72\ell_{f,1}^2 \alpha_t^2 + \frac{27\ell_{g,1}^2 \alpha_t}{c_v}) \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) \right) 2c_\beta^2 \\ &\leq -\frac{4M_f^2}{L_{\mu_g}^2} \alpha_t + \left(\frac{6d_2}{L_f} \frac{\ell_{g,1}^2}{\Phi} \alpha_t + (\frac{18}{L_f} \ell_{f,1}^2 \alpha_t + \frac{27\ell_{g,1}^2 \alpha_t}{c_v}) \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) \right) 2c_\beta^2 \\ &\leq -\frac{M_f^2}{L_{\mu_g}^2} \alpha_t, \end{aligned} \tag{248}$$

4279 where the first equality is by $\Gamma = \frac{11M_f^2}{L_{\mu_g} c_\beta}$ and $\rho_v^2 = c_v \alpha_t$; the first inequality follows from $c_\gamma \geq \frac{26M_f^2 \Phi}{L_{\mu_g}^2}$; the second inequality
 4280 is by $\alpha_t \leq 1/4L_f$; the last inequality follows from $c_v \geq \frac{54L_f^2}{M_f^2} \ell_{g,1}^2 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) c_\beta^2$, $\Phi \geq \frac{12d_2 \ell_{g,1}^2 L_{\mu_g}^2 c_\beta^2}{L_f M_f^2}$, and $\Psi \geq \frac{36\ell_{f,1}^2 d_2 L_{\mu_g}^2 c_\beta^2}{L_f M_f^2}$,
 4281 and $\Omega \geq \frac{36\ell_{f,1}^2 d_1 L_{\mu_g}^2 c_\beta^2}{L_f M_f^2}$.

4290 From (235), $\beta_t = c_\beta \alpha_t$, $\rho_{\mathbf{v}}^2 = c_{\mathbf{v}} \alpha_t$ and (247), we have

$$\begin{aligned}
 4292 \quad S(\alpha_t, \beta_t, \delta_t) &= -\frac{2\beta_t \Gamma}{\mu_g + \ell_{g,1}} + \beta_t^2 \Gamma + 2F(\alpha_t) \beta_t^2 \\
 4293 \quad &= -\frac{2c_\beta \alpha_t \Gamma}{\mu_g + \ell_{g,1}} + c_\beta^2 \alpha_t^2 \Gamma + \left(24d_2 \frac{\ell_{g,1}^2}{\Phi} + (72\ell_{f,1}^2 + \frac{27\ell_{g,1}^2}{c_{\mathbf{v}} \alpha_t})(\frac{d_2}{\Psi} + \frac{d_1}{\Omega}) \right) 2c_\beta^2 \alpha_t^2 \\
 4294 \quad &\leq -\frac{c_\beta \alpha_t \Gamma}{\mu_g + \ell_{g,1}} + \left(24d_2 \frac{\ell_{g,1}^2}{\Phi} \alpha_t^2 + (72\ell_{f,1}^2 \alpha_t^2 + \frac{27\ell_{g,1}^2 \alpha_t}{c_{\mathbf{v}}})(\frac{d_2}{\Psi} + \frac{d_1}{\Omega}) \right) 2c_\beta^2 \\
 4295 \quad &\leq -\frac{c_\beta \alpha_t \Gamma}{\mu_g + \ell_{g,1}} + \left(\frac{6d_2}{L_f} \frac{\ell_{g,1}^2}{\Phi} \alpha_t + (\frac{18}{L_f} \ell_{f,1}^2 \alpha_t + \frac{27\ell_{g,1}^2 \alpha_t}{c_{\mathbf{v}}})(\frac{d_2}{\Psi} + \frac{d_1}{\Omega}) \right) 2c_\beta^2 \\
 4296 \quad &\leq -\frac{c_\beta \alpha_t \Gamma}{4(\mu_g + \ell_{g,1})}, \tag{249}
 \end{aligned}$$

4305 where the first inequality follows from $\alpha_t \leq 1/c_\beta(\mu_g + \ell_{g,1})$; the second inequality is by $\alpha \leq 1/4L_f$; the last inequality is by
 4306 $c_{\mathbf{v}} \geq \frac{216}{\Gamma} \ell_{g,1}^2 (\frac{d_2}{\Psi} + \frac{d_1}{\Omega}) c_\beta (\mu_g + \ell_{g,1})$ and $\Phi \geq \frac{24d_2 \ell_{g,1}^2 (\mu_g + \ell_{g,1})}{L_f c_\beta \Gamma}$, and $\Psi \geq \frac{144d_2 \ell_{f,1}^2 (\mu_g + \ell_{g,1}) c_\beta}{L_f \Gamma}$, and $\Omega \geq \frac{144d_1 \ell_{f,1}^2 (\mu_g + \ell_{g,1}) c_\beta}{L_f \Gamma}$.
 4307 Thus, we get

$$4309 \quad (233e) = \sum_{t=1}^T Q(\alpha_t, \beta_t, \delta_t) \mathbb{E} \|e_t^{g_{\rho}}\|^2 + \sum_{t=1}^T S(\alpha_t, \beta_t, \delta_t) \mathbb{E} [\|\nabla_{\mathbf{y}} g_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t)\|^2] \leq 0. \tag{250}$$

4313 Bounding (233f).

4314 From (234), we obtain

$$4315 \quad Z(\alpha_t) = 27\ell_{g,1}^2 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right).$$

4318 Thus, from $\delta_t = c_\delta \alpha_t$, we have

$$\begin{aligned}
 4320 \quad (233f) &= \sum_{t=1}^T Z(\alpha_t) (3d_2^2 \ell_{f,1}^2 \delta_t^2 \rho_{\mathbf{r}}^2 + (12\ell_{f,0}^2 + 6\ell_{g,1}^2 p^2) \delta_t^2) \\
 4321 \quad &= \sum_{t=1}^T 27\ell_{g,1}^2 \left(\frac{d_2}{\Psi} + \frac{d_1}{\Omega} \right) (3d_2^2 \ell_{f,1}^2 \rho_{\mathbf{r}}^2 + (12\ell_{f,0}^2 + 6\ell_{g,1}^2 p^2)) c_\delta^2 \alpha_t^2 \\
 4322 \quad &= \mathcal{O} \left(\sum_{t=1}^T (d_1 + d_2)(\alpha_t^2 \rho_{\mathbf{r}}^2 + \alpha_t^2) \right). \tag{251}
 \end{aligned}$$

4329 **Bounding (233g)**. From $\rho_{\mathbf{v}}^2 = c_{\mathbf{v}} \alpha_t$, we have

$$\begin{aligned}
 4331 \quad (233g) &= \frac{36}{\Psi} D_{\mathbf{y},T} + \frac{36}{\Omega} D_{\mathbf{x},T} + \frac{12}{\Phi} G_{\mathbf{y},T} + \frac{36}{\Psi \rho_{\mathbf{v}}^2} G_{\mathbf{v},T} + \frac{36}{\Omega \rho_{\mathbf{v}}^2} G_{\mathbf{x},T} \\
 4332 \quad &= \mathcal{O} \left(D_{\mathbf{y},T} + D_{\mathbf{x},T} + G_{\mathbf{y},T} + \frac{1}{\alpha_T} (G_{\mathbf{v},T} + G_{\mathbf{x},T}) \right). \tag{252}
 \end{aligned}$$

4335
 4336
 4337
 4338
 4339
 4340
 4341
 4342
 4343
 4344

4345 **Bounding (233h)**. From $\gamma_{t+1} = c_\gamma \alpha_t$, $\eta_{t+1} = c_\eta \alpha_t$, $\lambda_{t+1} = c_\lambda \alpha_t$ and $\rho_v^2 = c_v \alpha_t$, we have

$$\begin{aligned}
 (233h) &= 2 \sum_{t=1}^T \frac{\gamma_{t+1}^2}{\Phi} \frac{\hat{\sigma}_{g_y}^2}{b} + 3 \left(\frac{\hat{\sigma}_{g_y}^2}{b \rho_v^2} + \frac{\hat{\sigma}_{f_y}^2}{b} \right) \sum_{t=1}^{T+1} \frac{\lambda_{t+1}^2}{\Psi} + 3 \left(\frac{\hat{\sigma}_{g_x}^2}{b \rho_v^2} + \frac{\hat{\sigma}_{f_x}^2}{b} \right) \sum_{t=1}^T \frac{\eta_{t+1}^2}{\Omega} \\
 &= 2 \sum_{t=1}^T \frac{c_\gamma^2 \alpha_t^2}{\Phi} \frac{\hat{\sigma}_{g_y}^2}{b} + 3 \left(\frac{\hat{\sigma}_{g_y}^2}{b \rho_v^2} + \frac{\hat{\sigma}_{f_y}^2}{b} \right) \sum_{t=1}^{T+1} \frac{c_\lambda^2 \alpha_t^2}{\Psi} + 3 \left(\frac{\hat{\sigma}_{g_x}^2}{b \rho_v^2} + \frac{\hat{\sigma}_{f_x}^2}{b} \right) \sum_{t=1}^T \frac{c_\eta^2 \alpha_t^2}{\Omega} \\
 &= \mathcal{O} \left(\left(\frac{\hat{\sigma}_{g_y}^2}{b} + \frac{\hat{\sigma}_{g_y}^2}{b \alpha_t} + \frac{\hat{\sigma}_{f_y}^2}{b} + \frac{\hat{\sigma}_{g_x}^2}{b \alpha_t} + \frac{\hat{\sigma}_{f_x}^2}{b} \right) \sum_{t=1}^T \alpha_t^2 \right). \tag{253}
 \end{aligned}$$

4356 **Bounding (233i)**. From $\beta_t = c_\beta \alpha_t$, $\delta_t = c_\delta \alpha_t$, we have

$$\begin{aligned}
 D(\alpha_t, \beta_t, \delta_t) &= 6\ell_{f,1}(1 + 2\frac{\ell_{g,1}}{\mu_g}) + \frac{3\ell_{f,1}\ell_{g,1}}{\mu_g}(\alpha_t - L_f \alpha_t^2) + \frac{24\ell_{g,1}}{L_{\mu_g} \mu_g} \frac{\Gamma}{\beta_t} + \frac{96\ell_{g,1}\nu^2}{L_{\mu_g} \mu_g^3} \frac{\Upsilon}{\delta_t} \\
 &= 6\ell_{f,1}(1 + 2\frac{\ell_{g,1}}{\mu_g}) + \frac{3\ell_{f,1}\ell_{g,1}}{\mu_g}(\alpha_t - L_f \alpha_t^2) + \frac{24\ell_{g,1}}{L_{\mu_g} \mu_g} \frac{\Gamma}{c_\beta \alpha_t} + \frac{96\ell_{g,1}\nu^2}{L_{\mu_g} \mu_g^3} \frac{\Upsilon}{c_\delta \alpha_t} \\
 &= \mathcal{O} \left(\alpha_t + \frac{1}{\alpha_t} \right),
 \end{aligned}$$

4365 and

$$\sum_{t=1}^T D(\alpha_t, \beta_t, \delta_t) (\rho_s^2 + \rho_r^2) := \mathcal{O} \left(\sum_{t=1}^T (\alpha_t + \frac{1}{\alpha_t}) (\rho_s^2 + \rho_r^2) \right). \tag{254}$$

4370 Moreover, we have

$$\begin{aligned}
 R(\rho_v) &= 9d_2^2 \frac{\ell_{g,1}^2}{\Phi} + 18d_2^2 \frac{\ell_{f,1}^2}{\Psi} + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_v^2}) \frac{d_2^2}{\Psi} = \mathcal{O} \left((1 + \frac{1}{\rho_v^2}) d_2^2 \right), \\
 \acute{R}(\rho_v) &= 18d_1^2 \frac{\ell_{f,1}^2}{\Omega} + 6(3\ell_{f,1}^2 + \frac{3\ell_{g,1}^2}{4\rho_v^2}) \frac{d_1^2}{\Omega} = \mathcal{O} \left((1 + \frac{1}{\rho_v^2}) d_1^2 \right),
 \end{aligned}$$

4377 which, implies that

$$\begin{aligned}
 R(\rho_v)T\rho_r^2 + \acute{R}(\rho_v)T\rho_s^2 + 18d_1^2\ell_{g,1}^2 \frac{T\rho_s^2}{\Omega\rho_v^2} + 18d_2^2\ell_{g,1}^2 \frac{T\rho_r^2}{\Psi\rho_v^2} \\
 = \mathcal{O} \left((1 + \frac{1}{\rho_v^2}) T(d_1^2\rho_s^2 + d_2^2\rho_r^2) + \frac{T}{\rho_v^2} (d_2^2\rho_r^2 + d_1^2\rho_s^2) \right). \tag{255}
 \end{aligned}$$

4384 From (254), (255) and $\rho_v^2 = c_v \alpha_t$, we get

$$(233i) \leq \mathcal{O} \left(\sum_{t=1}^T (\alpha_t + \frac{1}{\alpha_t}) (\rho_s^2 + \rho_r^2) + (1 + \frac{1}{\alpha_T}) T(d_2^2\rho_r^2 + d_1^2\rho_s^2) \right). \tag{256}$$

4390 **Combining the outcomes (233i)**. Combining inequalities (240), (243), (244), (246), (250), (251), (252), (253), and (256)

leads to

$$\begin{aligned}
 & \frac{\alpha_T}{2} \sum_{t=1}^T \mathbb{E} \left[\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right] + \Lambda \\
 & \leq \mathcal{O} \left(V_T + \sum_{t=1}^T \alpha_t (\rho_s^2 + \rho_r^2) + \sum_{t=1}^T \alpha_t^2 + \frac{H_{2,T}}{\alpha_T} + \sum_{t=1}^T (d_1 + d_2)(\alpha_t^2 \rho_r^2 + \alpha_t^2) \right) \\
 & + \mathcal{O} \left(D_{\mathbf{y}, T} + D_{\mathbf{x}, T} + G_{\mathbf{y}, T} + \frac{1}{\alpha_T} (G_{\mathbf{v}, T} + G_{\mathbf{x}, T}) \right) \\
 & + \mathcal{O} \left(\sum_{t=1}^T \left(\frac{\hat{\sigma}_{g_y}^2 \alpha_t^2}{b} + \frac{\hat{\sigma}_{f_y}^2 \alpha_t^2}{b} + \frac{\hat{\sigma}_{g_x}^2 \alpha_t^2}{b} + \frac{\hat{\sigma}_{f_x}^2 \alpha_t^2}{b} + \frac{\hat{\sigma}_{g_{xy}}^2 \alpha_t^2}{b} \right) \right) \\
 & + \mathcal{O} \left(\sum_{t=1}^T \left(\alpha_t + \frac{1}{\alpha_t} \right) (\rho_s^2 + \rho_r^2) + (1 + \frac{1}{\alpha_T}) T (d_2^2 \rho_r^2 + d_1^2 \rho_s^2) \right).
 \end{aligned}$$

From the definition of Λ in (102), we have

$$\begin{aligned}
 -\Lambda &= \Gamma \sum_{t=1}^T (\mathbb{E}[\theta_t^y] - \mathbb{E}[\theta_{t+1}^y]) + \Upsilon \sum_{t=1}^T (\mathbb{E}[\theta_t^v] - \mathbb{E}[\theta_{t+1}^v]) + \frac{1}{\Phi} \sum_{t=1}^T (\mathbb{E}\|e_t^g\|^2 - \mathbb{E}\|e_{t+1}^g\|^2) \\
 &+ \frac{1}{\Psi} \sum_{t=1}^T (\mathbb{E}\|e_t^v\|^2 - \mathbb{E}\|e_{t+1}^v\|^2) + \frac{1}{\Omega} \sum_{t=1}^T (\mathbb{E}\|e_t^f\|^2 - \mathbb{E}\|e_{t+1}^f\|^2) \\
 &\leq \Gamma \theta_1^y + \Upsilon \theta_1^v + \frac{\hat{\sigma}_{g_y}^2}{\Phi} + \frac{\hat{\sigma}_{g_{yy}}^2 + \hat{\sigma}_{f_y}^2}{\Psi} + \frac{\hat{\sigma}_{g_{xy}}^2 + \hat{\sigma}_{f_x}^2}{\Omega}.
 \end{aligned} \tag{257}$$

From (23), we have

$$\hat{\sigma}^2 := \hat{\sigma}_{g_y}^2 + \hat{\sigma}_{g_{yy}}^2 + \hat{\sigma}_{f_y}^2 + \hat{\sigma}_{g_{xy}}^2 + \hat{\sigma}_{f_x}^2.$$

From (26), we have

$$\begin{aligned}
 \alpha_t &= \frac{1}{(d_1 + d_2)^{3/4}(c + t)^{1/3}}, \quad \beta_t = c_\beta \alpha_t, \quad \delta_t = c_\delta \alpha_t, \quad \rho_v^2 = c_v \alpha_t, \\
 \gamma_{t+1} &= c_\gamma \alpha_t, \quad \eta_{t+1} = c_\eta \alpha_t, \quad \lambda_{t+1} = c_\lambda \alpha_t, \quad \rho_r^2 = \frac{1}{d_2^2 T}, \quad \rho_s^2 = \frac{1}{d_1^2 T}, \\
 b &= \frac{T^{1/3}}{(d_1 + d_2)^{3/2}}, \quad \bar{b} = \frac{T^{2/3}}{(d_1 + d_2)^{3/4}}.
 \end{aligned} \tag{258}$$

4455 Thus, using (257), (258), and rearranging the terms, we get
 4456

$$\begin{aligned}
 & \sum_{t=1}^T \mathbb{E} \left[\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right] \\
 & \leq \frac{2}{\alpha_T} \mathcal{O} \left(V_T + \sum_{t=1}^T \alpha_t (\rho_s^2 + \rho_r^2) + \sum_{t=1}^T \alpha_t^2 + \frac{H_{2,T}}{\alpha_T} + \sum_{t=1}^T (d_1 + d_2)(\alpha_t^2 \rho_r^2 + \alpha_t^2) \right) \\
 & + \frac{2}{\alpha_T} \mathcal{O} \left(D_{\mathbf{y}, T} + D_{\mathbf{x}, T} + G_{\mathbf{y}, T} + \frac{1}{\alpha_T} (G_{\mathbf{v}, T} + G_{\mathbf{x}, T}) \right) \\
 & + \frac{2}{\alpha_T} \mathcal{O} \left(\sum_{t=1}^T \left(\frac{\hat{\sigma}_{g_y}^2 \alpha_t^2}{\bar{b}} + \frac{\hat{\sigma}_{g_y}^2 \alpha_t}{\bar{b}} + \frac{\hat{\sigma}_{f_y}^2 \alpha_t^2}{b} + \frac{\hat{\sigma}_{g_x}^2 \alpha_t}{\bar{b}} + \frac{\hat{\sigma}_{f_x}^2 \alpha_t^2}{b} \right) \right) \\
 & + \frac{2}{\alpha_T} \mathcal{O} \left(\sum_{t=1}^T \left(\alpha_t + \frac{1}{\alpha_t} \right) (\rho_s^2 + \rho_r^2) + (1 + \frac{1}{\alpha_T}) T (d_2^2 \rho_r^2 + d_1^2 \rho_s^2) \right) \\
 & + \frac{2}{\alpha_T} \mathcal{O} (\theta_1^{\mathbf{y}} + \theta_1^{\mathbf{v}} + \hat{\sigma}^2) \\
 & \leq \mathcal{O} \left((d_1 + d_2)^{3/4} T^{1/3} (V_T + D_{\mathbf{y}, T} + D_{\mathbf{x}, T} + G_{\mathbf{y}, T} + \Delta_1 + \hat{\sigma}^2) \right. \\
 & \quad \left. + (d_1 + d_2)^{3/2} T^{2/3} (H_{2,T} + G_{\mathbf{v}, T} + G_{\mathbf{x}, T}) \right), \tag{259}
 \end{aligned}$$

4477 where second inequality holds because we have
 4478

$$\begin{aligned}
 \sum_{t=1}^T \alpha_t^3 &= \sum_{t=1}^T \frac{1}{(d_1 + d_2)^{9/4}(c+t)} \leq \sum_{t=1}^T \frac{1}{(d_1 + d_2)^{9/4}(1+t)} \leq \frac{\log(T+1)}{(d_1 + d_2)^{9/4}}, \\
 \sum_{t=1}^T \alpha_t^2 &= \sum_{t=1}^T \frac{1}{(d_1 + d_2)^{3/2}(c+t)^{2/3}} \leq \sum_{t=1}^T \frac{1}{(d_1 + d_2)^{3/2}(1+t)^{2/3}} \leq \frac{T^{1/3}}{(d_1 + d_2)^{3/2}}, \\
 \sum_{t=1}^T \alpha_t &= \sum_{t=0}^T \frac{1}{(d_1 + d_2)^{3/4}(c+t)^{1/3}} \leq \sum_{t=1}^T \frac{1}{(d_1 + d_2)^{3/4}(1+t)^{1/3}} \leq \frac{3T^{2/3}}{2(d_1 + d_2)^{3/4}}, \\
 \sum_{t=1}^T \frac{1}{\alpha_t} &= \sum_{t=0}^T (d_1 + d_2)^{3/4}(c+t)^{1/3} \leq \frac{3}{2}(d_1 + d_2)^{3/4}T^{4/3}.
 \end{aligned}$$

4490 Then, note that, we have
 4491

$$\begin{aligned}
 & \frac{1}{2} \sum_{t=1}^T \mathbb{E} \left[\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right] \\
 & \leq \sum_{t=1}^T \mathbb{E} \left[\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right] \\
 & + \sum_{t=1}^T \mathbb{E} \left[\|\mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) - \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t, \rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)))\|^2 \right].
 \end{aligned}$$

4510 From non-expansiveness of the projection operator and Lemma D.4, we have

$$\begin{aligned}
 & \| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) - \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 \\
 & \leq \| \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) - \nabla f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t)) \|^2 \\
 & \leq \frac{(\rho_s d_1 + \rho_r d_2)^2 \ell_{f,1}^2}{4} \\
 & \leq \frac{(\rho_s^2 d_1^2 + \rho_r^2 d_2^2) \ell_{f,1}^2}{2}.
 \end{aligned}$$

4519 This implies

$$\begin{aligned}
 & \frac{1}{2} \sum_{t=1}^T \mathbb{E} \left[\| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 \right] \\
 & \leq \sum_{t=1}^T \mathbb{E} \left[\| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_{t,\rho}(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 \right] + \frac{T(\rho_s^2 d_1^2 + \rho_r^2 d_2^2) \ell_{f,1}^2}{2}.
 \end{aligned}$$

4527 Applying the upper bound in (259) yields

$$\begin{aligned}
 & \frac{1}{2} \sum_{t=1}^T \mathbb{E} \left[\| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 \right] \\
 & \leq \mathcal{O} \left((d_1 + d_2)^{3/4} T^{1/3} (V_T + D_{\mathbf{y},T} + D_{\mathbf{x},T} + G_{\mathbf{y},T} + \Delta_1 + \hat{\sigma}^2) \right. \\
 & \quad \left. + (d_1 + d_2)^{3/2} T^{2/3} (H_{2,T} + G_{\mathbf{v},T} + G_{\mathbf{x},T}) \right) \\
 & \quad + \frac{T(\rho_s^2 d_1^2 + \rho_r^2 d_2^2) \ell_{f,1}^2}{2}.
 \end{aligned}$$

4538 Thus, from $\rho_r^2 = \frac{1}{d_2^2 T}$ and $\rho_s^2 = \frac{1}{d_1^2 T}$, we get

$$\begin{aligned}
 & \sum_{t=1}^T \mathbb{E} \left[\| \mathcal{P}_{\mathcal{X}, \alpha_t}(\mathbf{x}_t; \nabla f_t(\mathbf{x}_t, \mathbf{y}_t^*(\mathbf{x}_t))) \|^2 \right] \\
 & \leq \mathcal{O} \left((d_1 + d_2)^{3/4} T^{1/3} (V_T + D_{\mathbf{y},T} + D_{\mathbf{x},T} + G_{\mathbf{y},T} + \Delta_1 + \hat{\sigma}^2) \right. \\
 & \quad \left. + (d_1 + d_2)^{3/2} T^{2/3} (H_{2,T} + G_{\mathbf{v},T} + G_{\mathbf{x},T}) \right).
 \end{aligned}$$

4547 This completes the proof. □

4549
4550
4551
4552
4553
4554
4555
4556
4557
4558
4559
4560
4561
4562
4563
4564