# DOUBLY SMOOTHED DECENTRALIZED STOCHASTIC MINIMAX OPTIMIZATION ALGORITHM

## Anonymous authors

Paper under double-blind review

#### **ABSTRACT**

Decentralized stochastic minimax optimization has recently attracted significant attention due to its applications in machine learning. However, existing state-of-the-art methods use learning rates of different scales for the primal and dual variables, making them difficult to tune in practice. To address this problem, this paper proposes a novel doubly smoothed decentralized stochastic minimax algorithm. Specifically, in terms of algorithm design, we update both the primal and dual variables using smoothed gradients and introduce novel approaches to handle the computation and communication of the auxiliary variables introduced by the smoothing technique. On the theoretical side, for nonconvex-PL problems, our convergence analysis reveals that the learning rates for the primal and dual variables are of the same scale. Moreover, the order of the condition number in our convergence rate is improved to  $O(\kappa^{3/2})$ . To the best of our knowledge, this is the first time it has been improved to such a favorable order. Finally, extensive experimental results validate the effectiveness of our algorithm.

## 1 Introduction

In this paper, we focus on the following decentralized stochastic minimax optimization problem:

$$\min_{x \in \mathbb{R}^{d_1}} \max_{y \in \mathbb{R}^{d_2}} f(x, y) \triangleq \frac{1}{K} \sum_{k=1}^{K} f^{(k)}(x, y) , \qquad (1)$$

where  $x \in \mathbb{R}^{d_1}$  is the primal variable,  $y \in \mathbb{R}^{d_2}$  is the dual variable,  $f^{(k)}(x,y) = \mathbb{E}[f^{(k)}(x,y;\xi^{(k)})]$  is the loss function on the k-th (where  $k \in \{1,\cdots,K\}$ ) worker, and  $\xi^{(k)}$  denotes the random sample on the k-th worker. Throughout this paper, it is assumed that f(x,y) is nonconvex in x and satisfies the Polyak-Lojasiewicz (PL) condition in y.

Stochastic minimax optimization has attracted increasing attention in the machine learning community recently because it finds numerous applications, such as generative adversarial networks (Goodfellow et al., 2014), adversarially robust learning (Madry et al., 2017), distributionally robust learning (Duchi et al., 2021), imbalanced data classification (Ying et al., 2016), policy evaluation (Zhang et al., 2021), etc. Moreover, in real-world machine learning applications, the training data is typically distributed on different devices. To take advantage of the distributed data to train the aforementioned machine learning models, decentralized minimax optimization has been actively studied in recent years. For example, Xian et al. (2021); Huang & Chen (2023) proposed decentralized stochastic variance-reduced gradient descent ascent algorithm based on the STORM gradient estimator (Cutkosky & Orabona, 2019), while Zhang et al. (2021; 2024) proposed to use the SPIDER gradient estimator (Fang et al., 2018; Nguyen et al., 2017). Recently, Huang et al. (2024) developed a decentralized adaptive minimax algorithm and established its convergence rate for nonconvex-strongly-concave problems.

However, most existing decentralized minimax optimization algorithms suffer from a significant limitation. Specifically, to ensure convergence, the learning rate for the primal variable is set on a different scale than that for the dual variable. For example, Xian et al. (2021); Zhang et al. (2024); Chen et al. (2024); Huang & Chen (2023) prove that the ratio between the learning rate of the primal variable and that of the dual variable has to be  $O(1/\kappa^2)$  for nonconvex-strongly-concave (or nonconvex-PL) problems, where  $\kappa > 1$  is the condition number. Since  $\kappa$  is an unknown parameter, it is difficult to tune their learning rates to ensure convergence in practice. To address this issue, a recent

Table 1: The communication complexity (i.e., iteration complexity) of different decentralized stochastic minimax algorithms that using variance-reduced gradients. **N-PL**: denotes nonconvex-PL problems. **N-SCV**: denotes nonconvex-strongly-concave problems. **LR Ratio**: the ratio between the learning rate of the primal variable and that of the dual variable.  $\kappa$ : denotes condition number.  $1 - \lambda$ : denotes spectral gap.  $\epsilon$ : denotes solution accuracy. Note that Smoothed-SAGDA is a *single-machine* algorithm *without using variance-reduced gradients*. DGDA-VR and DREAM depend on the condition number, scaling as  $\kappa^2$ , in the cost of a large batch size  $O\left(\frac{\kappa}{\epsilon}\right)$ . DREAM achieves a better dependence on the spectral gap in the cost of performing multi-round communication in each iteration.

Algorithms	Communication Complexity	Batch Size	Problem Clas	s LR Ratio
Smoothed-SAGDA (Yang et al., 2022	$O\left(\frac{\kappa^2}{\epsilon^4}\right)$	O(1)	N-PL	O(1)
DM-HSGD (Xian et al., 2021)	$O\left(\frac{\kappa^3}{(1-\lambda)^2\epsilon^3}\right)$ $O\left(\frac{\kappa^2}{(1-\lambda)^2\epsilon^2}\right)$	O(1)	N-SCV	$O(1/\kappa^2)$
DGDA-VR (Zhang et al., 2024)	$O\left(\frac{\kappa^2}{(1-\lambda)^2\epsilon^2}\right)$	$O\left(\frac{\kappa}{\epsilon}\right)$	N-SCV	$O(1/\kappa^2)$
DREAM (Chen et al., 2024)	$O\left(\frac{\kappa^2}{\sqrt{1-\lambda}\epsilon^2}\right)$	$O\left(\frac{\kappa}{\epsilon}\right)$	N-SCV	$O(1/\kappa^2)$
DM-GDA (Huang & Chen, 2023)	$O\left(\frac{\kappa^3}{(1-\lambda)^2\epsilon^3}\right)$	O(1)	N-PL	$O(1/\kappa^2)$
Ours (Corollary 4.3)	$O\left(\frac{\kappa^{3/2}}{(1-\lambda)^2\epsilon^3}\right)$	O(1)	N-PL	O(1)

work (Yang et al., 2022) in the single-machine setting demonstrates that the smoothing technique proposed by Zhang et al. (2020) allows primal and dual variables to use learning rates of the same scale, that is, with a ratio of the order of O(1). However, the convergence rate  $^1$   $O(1/\epsilon^4)$  of Yang et al. (2022) is inferior to  $O(1/\epsilon^3)$  of Xian et al. (2021); Huang & Chen (2023) because it just uses the standard stochastic gradient. Then, a natural question arises:

Can we develop a decentralized smoothed minimax optimization algorithm that achieves a better convergence rate while using same-scale learning rates for the primal and dual variables?

Actually, there are unique challenges when applying the smoothing technique to decentralized minimax optimization in order to improve the convergence rate, as outlined below.

Challenge-1: How to incorporate the variance reduction technique into the smoothing technique to achieve a faster convergence rate? Existing minimax optimization algorithms with the smoothing technique in a single machine setting are based on the *deterministic gradient* (Zhang et al., 2020) or the *unbiased stochastic gradient* (Yang et al., 2022). Directly extending their smoothing technique to decentralized *stochastic* minimax optimization will lead to a slow convergence rate. For example, (Yang et al., 2022) can only achieve a  $O(1/\epsilon^4)$  convergence rate to achieve the  $\epsilon$ -accuracy solution for a nonconvex-PL problem, while the existing decentralized minimax optimization algorithm (Huang & Chen, 2023) can achieve a  $O(1/\epsilon^3)$  convergence rate for the same problem class by using the variance reduction technique. However, due to the existence of the auxiliary variable in the smoothing technique, it is unclear how to leverage the variance reduction technique to accelerate its convergence rate. For example, it is unclear which component in the smoothed gradient should use variance reduction and how to control the gradient bias to guarantee the fast convergence rate.

Challenge-2: How to compute and communicate the auxiliary variable in the smoothing technique and how does it affect the communication complexity? The standard smoothing technique introduces an auxiliary variable to smooth the loss landscape with respect to the primal variable to improve the convergence rate. However, in a decentralized setting, it is unclear how to update and communicate the auxiliary variable. In particular, due to the strong dependence between the original variable and the auxiliary variable, it remains unclear whether the communication of the auxiliary variable, especially given that our algorithm introduces auxiliary variables for both the primal and dual variables, will improve or degrade the communication complexity, for example, by affecting the dependence on the spectral gap or condition number in the convergence rate.

<sup>&</sup>lt;sup>1</sup>In the introduction, we omit other factors in the convergence rate for clarity.

To answer the aforementioned questions, we develop a novel decentralized algorithm based on the smoothing technique: the doubly smoothed decentralized stochastic gradient descent ascent with momentum (Smoothed<sup>2</sup>-DSGDAM) algorithm, which brings the following contributions:

- In terms of algorithm design, we apply the smoothing technique to both the primal and dual variables. Importantly, we propose a novel and feasible approach to incorporate the variance reduction technique into the smoothed gradient regarding both the primal and dual variables. More importantly, our algorithm demonstrates how to update and communicate the auxiliary variable introduced by the smoothing technique in the decentralized setting. As far as we know, this is the first time to show how to handle the auxiliary variable and reduced the gradient variance for the decentralized smoothed minimax algorithm.
- In terms of convergence analysis, we establish the convergence rate of our algorithm for nonconvex-PL minimax problems. In particular, on the one hand, for a nonconvex-PL minimax problem, the smoothing technique with a variance-reduced gradient can make the convergence rate enjoy a better dependence on the condition number  $\kappa$ , i.e., in the order of  $O(\kappa^{3/2})$ , which is better than the dependence  $O(\kappa^3)$  in existing decentralized non-smoothed minimax algorithms (Xian et al., 2021; Huang & Chen, 2023) and the dependence  $O(\kappa^2)$  in smoothed minimax algorithms (Yang et al., 2022) in the single-machine setting  $^2$ . To the best of our knowledge, this is the first time the dependence on the condition number is improved to  $O(\kappa^{3/2})$ . On the other hand, our convergence analysis shows that the ratio between the learning rate of the primal variable and that of the dual variable can be improved from  $O(1/\kappa^2)$  of Xian et al. (2021); Zhang et al. (2024); Chen et al. (2024); Huang & Chen (2023) to O(1), and the convergence rate can be improved from  $O(1/\epsilon^4)$  of Yang et al. (2022) to  $O(1/\epsilon^3)$ . To the best of our knowledge, this is the first time that a decentralized stochastic minimax optimization algorithm can achieve such a fast convergence rate with the same-scale learning rate. A detailed comparison between our algorithm and existing algorithms can be found in Table 1.

Finally, the extensive experimental results validate the performance of our proposed algorithm.

### 2 RELATED WORKS

#### 2.1 STOCHASTIC MINIMAX OPTIMIZATION

Due to the widespread application of stochastic minimax optimization in machine learning, numerous stochastic optimization algorithms (Lin et al., 2020; Luo et al., 2020; Huang et al., 2022; Qiu et al., 2020; Guo et al., 2021; Yang et al., 2020; 2022; Chen et al., 2022) have been developed recently. In particular, the nonconvex-strongly-concave and nonconvex-PL problems have been extensively studied. For the former, Lin et al. (2020) established the convergence rate of the stochastic gradient descent ascent (SGDA) algorithm for nonconvex-strongly-concave problems. Following that, a couple of variance-reduced gradient methods (Luo et al., 2020; Huang et al., 2022; Qiu et al., 2020; Guo et al., 2021) have been developed to accelerate its convergence rate. Specifically, Huang et al. (2022); Qiu et al. (2020) combined the STORM gradient estimator (Cutkosky & Orabona, 2019) with SGDA and established its convergence rate. Luo et al. (2020) investigated the convergence rate when incorporating the SPIDER gradient estimator (Fang et al., 2018) into SGDA. For the latter, Yang et al. (2020) investigated the convergence rate for the alternating stochastic gradient descent ascent (ASGDA) algorithm. Chen et al. (2022) studied the convergence rate for the finite-sum minimax problem when combining the SPIDER gradient estimator with ASGDA.

The smoothing technique for the minimax problem was first studied for nonconvex-concave problems in Zhang et al. (2020). Specifically, it established the convergence rate of the full alternating gradient (AGDA) descent ascent algorithm when incorporating the smoothing technique. Later, Yang et al. (2022) applied this technique to nonconvex-PL problems and established its convergence rate for SGDA. In fact, due to the efficacy of the smoothing technique, it has been applied to various settings, such as nonconvex-nonconcave problems with the one-sided KŁ condition (Zheng et al., 2023), constrained optimization problems (Pu et al., 2024), etc, which are beyond the scope of this paper.

#### 2.2 DECENTRALIZED STOCHASTIC MINIMAX OPTIMIZATION

To facilitate decentralized optimization for minimax problems, a great amount of effort (Tsaknakis et al., 2020; Zhang et al., 2021; Xian et al., 2021; Gao, 2022; Zhang et al., 2024; Chen et al., 2024;

<sup>&</sup>lt;sup>2</sup>Here, to make a fair comparison, the existing methods considered use a batch size of O(1), rather than large batch sizes.

Xu, 2024) has recently been made. For example, Tsaknakis et al. (2020) developed a decentralized gradient descent ascent algorithm by using the full gradient for local computation and the gradient tracking technique for communication. Xian et al. (2021) proposed a decentralized minimax algorithm based on the STORM gradient estimator and established its convergence rate for the stochastic setting. Zhang et al. (2021) developed a decentralized minimax algorithm based on the SPIDER gradient estimator and established its convergence rate for the finite-sum setting. Later, its convergence rate for the stochastic setting was established in Zhang et al. (2024). Moreover, Gao (2022) incorporated the ZeroSARAH gradient estimator into the decentralized minimax algorithm and provided convergence analysis for the finite-sum setting. Recently, Chen et al. (2024) studied the convergence rate of the decentralized minimax algorithm when using the PAGE gradient estimator (Li et al., 2021). More recently, Huang et al. (2024) introduced the adaptive learning rate to decentralized minimax optimization and established the corresponding convergence rate. Note that all these existing methods restrict their focus to the nonconvex-strongly-concave problem.

Recently, Huang & Chen (2023) developed a decentralized minimax algorithm for nonconvex-PL problems, where the STORM gradient estimator is used for local updates and the gradient tracking technique is used for communication between workers. To our knowledge, in the distributed setting, the smoothing technique has only been studied for federated centralized learning in Shen et al. (2024). Specifically, each worker uses the standard unbiased stochastic gradient to do local update and the central server uses the smoothing technique to assist the update of the dual variable. As a result, the additional variable introduced by the smoothing technique behaves as a single-machine setting. Thus, it is easy to handle this variable in convergence analysis. All in all, the smoothing technique has not been studied for decentralized minimax optimization and it is unclear how to apply it from the algorithm design perspective and how to handle it from the convergence analysis perspective.

## 3 Method

#### 3.1 PROBLEM SETUP

We introduce the following assumptions with respect to the loss function and communication topology, which have been widely used in the existing literature (Yang et al., 2022; Xian et al., 2021; Huang & Chen, 2023; Zhang et al., 2021; 2024; Chen et al., 2024).

**Assumption 3.1.** (Smoothness) For any  $k \in \{1, 2, \dots, K\}$ , the loss function on the k-th worker satisfies the mean-squared Lipschitz smoothness, i.e., for any  $(x_1, y_1) \in \mathbb{R}^{d_1} \times \mathbb{R}^{d_2}$  and  $(x_2, y_2) \in \mathbb{R}^{d_1} \times \mathbb{R}^{d_2}$ , there exists a constant value L > 0 such that  $\mathbb{E}[\|\nabla_x f^{(k)}(x_1, y_1; \xi^{(k)}) - \nabla_x f^{(k)}(x_2, y_2; \xi^{(k)})\|^2] \le L^2(\|x_1 - x_2\|^2 + \|y_1 - y_2\|^2)$  and  $\mathbb{E}[\|\nabla_y f^{(k)}(x_1, y_1; \xi^{(k)}) - \nabla_y f^{(k)}(x_2, y_2; \xi^{(k)})\|^2] \le L^2(\|x_1 - x_2\|^2 + \|y_1 - y_2\|^2)$ .

**Assumption 3.2.** (PL condition) For any fixed  $x \in \mathbb{R}^{d_1}$ , the set of solutions of the optimization problem with respect to y,  $\max_{y \in \mathbb{R}^{d_2}} f(x,y)$ , is not empty and the optimal value is finite. Furthermore, for any  $x \in \mathbb{R}^{d_1}$ , there exists a constant value  $\mu > 0$  such that  $\|\nabla_y f(x,y)\|^2 \ge 2\mu(f(x,y^*) - f(x,y))$ , where  $y^* = \arg\max_{y \in \mathbb{R}^{d_2}} f(x,y)$ .

**Assumption 3.3.** (Variance) For any  $k \in \{1, 2, \cdots, K\}$ , the stochastic gradients with respect to x and y of the loss function on the k-th worker are unbiased estimators and their variances are upper bounded as:  $\mathbb{E}[\|\nabla_x f^{(k)}(x, y; \xi^{(k)}) - \nabla_x f^{(k)}(x, y)\|^2] \le \sigma^2$  and  $\mathbb{E}[\|\nabla_y f^{(k)}(x, y; \xi^{(k)}) - \nabla_y f^{(k)}(x, y)\|^2] \le \sigma^2$ , where  $\sigma > 0$  is a constant value.

**Assumption 3.4.** (Communication graph) The element  $w_{ij}$  of the adjacency matrix  $W \in \mathbb{R}^{K \times K}$  of the communication graph is non-negative, with a positive value indicating that worker-i is connected to worker-j, and a value of zero indicating they are disconnected. Moreover, W is doubly stochastic and symmetric, and its eigenvalues satisfy  $|\lambda_K| \leq |\lambda_{K-1}| \leq \cdots \leq |\lambda_2| < |\lambda_1| = 1$ .

By denoting  $\lambda=|\lambda_2|$ , the spectral gap of the adjacency matrix can by represented by  $1-\lambda$ . Moreover, we use  $\mathcal{N}_k$  to denote the neighboring worker of the k-th worker, and use  $\kappa=L/\mu$  to represent the condition number. In addition, we use  $\bar{a}_t=\frac{1}{K}\sum_{k=1}^K a_t^{(k)}$  to denote the average value of any  $\{a_t^{(k)}\}_{k=1}^K$  in the t-th iteration.

## 3.2 SMOOTHED<sup>2</sup>-DSGDAM

The essential idea of the smoothing technique is to introduce a regularization term such that the original nonconvex function becomes strongly convex. As a result, the update of the primal and dual

216

237 238

239

240

241

242

243

244 245

246

247

248

249

250

251

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

Algorithm 1 Doubly Smoothed Decentralized Stochastic Gradient Descent Ascent with Momentum (Smoothed<sup>2</sup>-DSGDAM)

```
217
218
                                                                                                           Input: \eta > 0, \rho_x > 0, \rho_y > 0, \beta_x > 0, \beta_y > 0, \hat{\beta}_x > 0, \hat{\beta}_y > 0, \rho_x \eta^2 < 1, \rho_y \eta^2 < 1, \hat{\beta}_x \eta < 1, \hat{\beta}_y \eta < 1.
219
                                                                                                                                                 Initialization on worker k: x_0^{(k)} = x_0, y_0^{(k)} = y_0, \hat{x}_0^{(k)} = \hat{x}_0, \quad x_0^{(k)} = \hat{x}_0,
220
221
222
223
                                                                                                                      1: for t = 0, \dots, T - 1, worker k do
                                                                                                             1: for t = 0, \dots, T-1, worker k do
2: Update x: \tilde{x}_{t+1}^{(k)} = \sum_{j \in \mathcal{N}_k} w_{kj} x_t^{(j)} - \beta_x p_t^{(k)}, x_{t+1}^{(k)} = x_t^{(k)} + \eta(\tilde{x}_{t+1}^{(k)} - x_t^{(k)}),
3: Update y: \tilde{y}_{t+1}^{(k)} = \sum_{j \in \mathcal{N}_k} w_{kj} y_t^{(j)} + \beta_y q_t^{(k)}, y_{t+1}^{(k)} = y_t^{(k)} + \eta(\tilde{y}_{t+1}^{(k)} - y_t^{(k)}),
4: Update \hat{x}: \tilde{x}_{t+1}^{(k)} = \sum_{j \in \mathcal{N}_k} w_{kj} \hat{x}_t^{(j)} + \hat{\beta}_x (x_{t+1}^{(k)} - \hat{x}_t^{(k)}), \hat{x}_{t+1}^{(k)} = \hat{x}_t^{(k)} + \eta(\tilde{x}_{t+1}^{(k)} - \hat{x}_t^{(k)}),
5: Update \hat{y}: \tilde{y}_{t+1}^{(k)} = \sum_{j \in \mathcal{N}_k} w_{kj} \hat{y}_t^{(j)} + \hat{\beta}_y (y_{t+1}^{(k)} - \hat{y}_t^{(k)}), \hat{y}_{t+1}^{(k)} = \hat{y}_t^{(k)} + \eta(\tilde{y}_{t+1}^{(k)} - \hat{y}_t^{(k)}),
224
225
226
227
228
                                                                                                                                                                          \begin{array}{lll} & & & & \\ u_{t+1}^{(k)} & = & (1 & - & \rho_x \eta^2)(u_t^{(k)} & - & & \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \xi_{t+1}^{(k)})) \\ \nabla_x F^{(k)}(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \hat{x}_{t+1}^{(k)}, \hat{y}_{t+1}^{(k)}; \xi_{t+1}^{(k)}) & , \\ v_{t+1}^{(k)} & = & (1 & - & \rho_y \eta^2)(v_t^{(k)} & - & \nabla_y F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \xi_{t+1}^{(k)})) \\ \nabla_y F^{(k)}(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \hat{x}_{t+1}^{(k)}, \hat{y}_{t+1}^{(k)}; \xi_{t+1}^{(k)}) & , \\ \text{Gradient tracking:} & & \\ n^{(k)} & = & \ddots & \\ n^{(k)} & =
229
230
231
232
                                                                                                                                                                                 p_{t+1}^{(k)} = \sum_{j \in \mathcal{N}_k} w_{kj} p_t^{(k)} + u_{t+1}^{(k)} - u_t^{(k)} , \quad q_{t+1}^{(k)} = \sum_{j \in \mathcal{N}_k} w_{kj} q_t^{(k)} + v_{t+1}^{(k)} - v_t^{(k)} ,
235
236
```

variables can be well coordinated to avoid divergence. Inspired by this, we introduce the doubly smoothed loss function, which adds the regularization term to both the primal and dual variables such that the nonconvex-PL loss function becomes strongly-convex-strongly-concave. Specifically, the doubly smoothed loss function is defined as follows:

$$F(x,y;\hat{x},\hat{y}) = \frac{1}{K} \sum_{k=1}^{K} \underbrace{f^{(k)}(x,y) + \frac{\gamma_1}{2} \|x - \hat{x}\|^2 - \frac{\gamma_2}{2} \|y - \hat{y}\|^2}_{F^{(k)}(x,y;\hat{x},\hat{y})},$$
 (2)

where  $\gamma_1 > 0$  and  $\gamma_2 > 0$  are hyperparameters, and  $\hat{x} \in \mathbb{R}^{d_1}$ ,  $\hat{y} \in \mathbb{R}^{d_2}$  are the auxiliary variables for the primal and dual variables, respectively. Here,  $\gamma_1$  and  $\gamma_2$  are set such that  $F(x,y;\hat{x},\hat{y})$  is strongly convex with respect to x and strongly concave with respect to y. For example, we can set  $\gamma_1 = 2L$  and  $\gamma_2 = 2L$ . Note that most existing works in the single-machine setting, such as (Zhang et al., 2020; Yang et al., 2022) apply the smoothing technique to a single variable. Only a recent work (Zheng et al., 2023) uses it for both variables for nonconvex-nonconcave problems. However, it focuses on the deterministic setting, failing to handle the biased stochastic gradient estimator and the decentralized communication. Hence, a new algorithm design and convergence analysis are required to address the challenges caused by them.

Based on the smoothed loss function in Eq. (2), the k-th worker can compute the stochastic gradient with respect to the primal and dual variables in the t-th iteration as follows:

$$\nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \xi_t^{(k)}) = \nabla_x f^{(k)}(x_t^{(k)}, y_t^{(k)}; \xi_t^{(k)}) + \gamma_1(x_t^{(k)} - \hat{x}_t^{(k)}),$$

$$\nabla_y F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \xi_t^{(k)}) = \nabla_y f^{(k)}(x_t^{(k)}, y_t^{(k)}; \xi_t^{(k)}) - \gamma_2(y_t^{(k)} - \hat{y}_t^{(k)}).$$
(3)

In terms of the smoothed loss function in Eq. (2) and the stochastic gradients in Eq. (3), we develop a novel doubly smoothed decentralized stochastic gradient descent ascent with momentum (Smoothed<sup>2</sup>-DSGDAM) algorithm in Algorithm 1. Generally speaking, we apply the variance reduction technique, STORM (Cutkosky & Orabona, 2019), to the stochastic gradient on each worker to update the primal and dual variables, and use the gradient tracking technique to conduct communication between different workers. However, there are two unique challenges when designing our Smoothed<sup>2</sup>-DSGDAM algorithm: 1) How to apply the variance reduction technique in the presence of the smoothing term? 2) How to update and communicate the auxiliary variables  $\hat{x}$ and  $\hat{y}$  to guarantee convergence in the decentralized setting?

As for the first challenge regarding variance reduction, there are actually two ways to apply variance reduction. Specifically, we can apply it to the original stochastic gradient  $\nabla_x f^{(k)}(x_t^{(k)}, y_t^{(k)}; \xi_t^{(k)})$  or to the smoothed gradient  $\nabla_x F^{(k)}(x_t^{(k)},y_t^{(k)};\hat{x}_t^{(k)},\hat{y}_t^{(k)};\xi_t^{(k)})$ . However, computing the variance-reduced gradient  $u_t^{(k)}$  for the original stochastic gradient  $\nabla_x f^{(k)}(x_{t_-}^{(k)},y_t^{(k)};\xi_t^{(k)})$  will complicate the convergence analysis, when bounding a critical term  $\langle \nabla_x F(\bar{x}_t,\bar{y}_t;\hat{x}_t,\hat{y}_t),\bar{x}_{t+1}-\bar{x}_t \rangle$ , where  $\bar{x}_t,\bar{y}_t,\bar{x}_t,\bar{y}_t$ , and  $\bar{y}_t$  denote the averaged variable across workers.

Specifically, when computing the STORM gradient estimator  $u_t^{(k)}$  for the smoothed gradient  $\nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \xi_t^{(k)})$ , we can bound it as follows:

$$\langle \nabla_x F(\bar{x}_t, \bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t), \bar{x}_{t+1} - \bar{x}_t \rangle = -\eta \beta_x \langle \nabla_x F(\bar{x}_t, \bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t), \bar{u}_t \rangle$$

$$= -\frac{\eta \beta_x}{2} \|\bar{u}_t\|^2 - \frac{\eta \beta_x}{2} \|\nabla_x F(\bar{x}_t, \bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t)\|^2 + \frac{\eta \beta_x}{2} \|\nabla_x F(\bar{x}_t, \bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t) - \bar{u}_t\|^2. \tag{4}$$

All three terms in the last step are straightforward to handle.

However, when computing the STORM gradient estimator  $u_t^{(k)}$  for the original stochastic gradient  $\nabla_x f^{(k)}(x_t^{(k)}, y_t^{(k)}; \xi_t^{(k)})$ , we have

$$\langle \nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \hat{\bar{x}}_{t}, \hat{\bar{y}}_{t}), \bar{x}_{t+1} - \bar{x}_{t} \rangle$$

$$= -\eta \beta_{x} \langle \nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \hat{\bar{x}}_{t}, \hat{\bar{y}}_{t}), \gamma_{1}(\bar{x}_{t} - \hat{\bar{x}}_{t}) \rangle - \eta \beta_{x} \langle \nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \hat{\bar{x}}_{t}, \hat{\bar{y}}_{t}), \bar{u}_{t} \rangle$$

$$\leq \frac{\eta \beta_{x}}{4} \|\nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \hat{\bar{x}}_{t}, \hat{\bar{y}}_{t})\|^{2} + 4\eta \beta_{x} \gamma_{1}^{2} \|\bar{x}_{t} - \hat{\bar{x}}_{t}\|^{2}$$

$$- \frac{\eta \beta_{x}}{2} \|\bar{u}_{t}\|^{2} - \frac{\eta \beta_{x}}{2} \|\nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \hat{\bar{x}}_{t}, \hat{\bar{y}}_{t})\|^{2} + \frac{\eta \beta_{x}}{2} \|\nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \hat{\bar{x}}_{t}, \hat{\bar{y}}_{t}) - \bar{u}_{t}\|^{2}. \tag{5}$$

Here, the last term should be further handled by  $\|\nabla_x F(\bar{x}_t,\bar{y}_t;\hat{\bar{x}}_t,\bar{\bar{y}}_t) - \bar{u}_t\|^2 \leq 2\|\nabla_x f(\bar{x}_t,\bar{y}_t) - \bar{u}_t\|^2 + 2\gamma_1^2\|\bar{x}_t - \bar{\bar{x}}_t\|^2$ . Then, it can be seen that this approach introduces an addition term (marked in blue), which can make it more challenging to select hyperparameters to handle  $\|\bar{x}_t - \bar{\hat{x}}_t\|^2$ .

In addition to this problem, if STORM is applied to the original stochastic gradient, whenever  $\|\bar{x}_{t+1} - \bar{x}_t\|^2$  appears, it should be decomposed into two terms:  $\|\bar{u}_t\|^2$  and  $\|\bar{x}_t - \hat{\bar{x}}_t\|^2$ . In contrast, the smoothed one only needs to replace  $\|\bar{x}_{t+1} - \bar{x}_t\|^2$  with  $\eta^2 \beta_x^2 \|\bar{u}_t\|^2$ , which is much easier for the downstream proof. All in all, applying the variance reduction technique to the original stochastic gradient could significantly complicate the proof. Therefore, we apply it to the smoothed gradient  $\nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \xi_t^{(k)})$ , which is shown in Step 6 of Algorithm 1.

Regarding the update and communication of the variable  $\hat{x}$  and  $\hat{y}$ , it has not been studied in the existing decentralized optimization literature. A straightforward approach is to update  $\hat{x}$  (and  $\hat{y}$ ) locally without communication. However, in convergence analysis, we have to handle the negative term,  $-\|\bar{x}_{t+1} - \bar{x}_t\|^2$  (See Lemma C.1), and positive term,  $\|\hat{x}_{t+1}^{(k)} - \hat{x}_t^{(k)}\|^2$  (See Lemma D.1), simultaneously. Specifically, we need to convert  $\|\hat{x}_{t+1}^{(k)} - \hat{x}_t^{(k)}\|^2$  to  $\|\bar{x}_{t+1} - \bar{x}_t\|^2$  based on the consensus error  $\|\hat{x}_t^{(k)} - \bar{x}_t\|^2$ . If there is no communication operation for  $\hat{x}$ , it is difficult to control the consensus error. In fact, it may be exploding. To address this challenge, we propose the following approach for the update and communication of  $\hat{x}$  (and  $\hat{y}$ ):

$$\tilde{\hat{x}}_{t+1}^{(k)} = \sum_{j \in \mathcal{N}_t} w_{kj} \hat{x}_t^{(j)} + \hat{\beta}_x (x_{t+1}^{(k)} - \hat{x}_t^{(k)}) , \quad \hat{x}_{t+1}^{(k)} = \hat{x}_t^{(k)} + \eta (\tilde{\hat{x}}_{t+1}^{(k)} - \hat{x}_t^{(k)}) , \tag{6}$$

where  $\hat{\beta}_x > 0$  and  $\eta > 0$  are hyperparameters. The first step in Eq. (6) can be viewed as the update of the communicated variable  $\sum_{j \in \mathcal{N}_k} w_{kj} \hat{x}_t^{(j)}$  with the local gradient  $x_{t+1}^{(k)} - \hat{x}_t^{(k)}$ , and the second step is a convex combination between the intermediate variable  $\tilde{x}_{t+1}^{(k)}$  and the local variable  $\hat{x}_t^{(k)}$ . With such an update and communication strategy, we can bound the consensus error regarding the auxiliary variable as shown in our Lemma D.3, where the coefficient  $1 - \frac{\eta(1-\lambda^2)}{4}$  is important to shrink the consensus error.

In summary, the smoothing technique brings new challenges for algorithm design in the decentralized setting. In our algorithm, we develop novel strategies to handle variance reduction and the update and communication of the auxiliary variables. Therefore, our algorithm design is novel.

## 4 Convergence Analysis

Before presenting the convergence rate of our algorithm, we introduce the following stationary measures, which were introduced in (Yang et al., 2022).

**Definition 4.1.** A solution (x,y) is termed the  $(\epsilon_1,\epsilon_2)$ -stationary solution if  $\|\nabla_x f(x,y)\| \le \epsilon_1$  and  $\|\nabla_y f(x,y)\| \le \epsilon_2$ . A solution x is termed the  $\epsilon$ -stationary solution if  $\|\nabla \Phi(x)\| \le \epsilon$ , where  $\Phi(x) = f(x,y^*(x))$  and  $y^*(x) = \arg\max_{y \in \mathbb{R}^{d_2}} f(x,y)$ .

Based on the assumptions in Section 3, we establish the convergence rate of our Algorithm 1 in the following theorem.

**Theorem 4.2.** Given Assumptions 3.1-3.4, when  $\rho_x > 0$ ,  $\rho_y > 0$ ,  $\gamma = O(L)$ , the condition about  $\eta$  and  $\beta_x$  in Eq. (164), and those about  $\beta_y$ ,  $\hat{\beta}_x$ ,  $\hat{\beta}_y$  in Eq. (54) hold, Algorithm 1 has the following convergence upper bound:

$$\frac{1}{T} \sum_{t=0}^{T-1} \left( \mathbb{E}[\|\nabla_x f(\bar{x}_t, \bar{y}_t)\|^2] + \kappa \mathbb{E}[\|\nabla_y f(\bar{x}_t, \bar{y}_t)\|^2] \right) \le O(\kappa \rho_x^2 \eta^4 \sigma^2) + O(\kappa \rho_y^2 \eta^4 \sigma^2) 
+ O\left(\frac{\kappa \mathcal{P}_0}{\beta_x \eta T}\right) + O\left(\frac{\kappa}{T} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\nabla_x f^{(k)}(x_0, y_0)\|^2]\right) + O\left(\frac{\kappa}{T} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\nabla_y f^{(k)}(x_0, y_0)\|^2]\right) 
+ O\left(\frac{\kappa \sigma^2}{\rho_x \eta^2 T B}\right) + O\left(\frac{\kappa \sigma^2}{\rho_y \eta^2 T B}\right) + O\left(\frac{\kappa \sigma^2}{T}\right) + O\left(\frac{\kappa \rho_x \eta^2 \sigma^2}{K}\right) + O\left(\frac{\kappa \rho_y \eta^2 \sigma^2}{K}\right). \tag{7}$$

where  $\mathcal{P}_0 = F(x_0, y_0; \hat{x}_0, \hat{y}_0) - 2F_d(y_0; \hat{x}_0, \hat{y}_0) + 2q(\hat{x}_0)$ , whose definitions can be found in Eq. (11). Corollary 4.3. Given Assumptions 3.1-3.4, by setting  $\beta_x = O((1-\lambda)^2)$ ,  $\beta_y = O((1-\lambda)^2)$ ,  $\hat{\beta}_x = O\left(\frac{(1-\lambda)^2}{\kappa}\right)$ ,  $\hat{\beta}_y = O((1-\lambda)^2)$ ,  $\eta = O\left(\frac{K\epsilon}{\kappa^{1/2}}\right)$ ,  $\rho_x = O\left(\frac{1}{K}\right)$ ,  $\rho_y = O\left(\frac{1}{K}\right)$ ,  $\theta_y = O\left(\frac{1}{K}\right)$ ,  $\theta_y = O\left(\frac{\kappa^{1/2}}{\kappa}\right)$ ,  $\theta_y = O\left(\frac{\kappa^{1/2}}{\kappa}\right)$ , Algorithm 1 can achieve the  $\theta_y = O(\kappa, \epsilon/\sqrt{\kappa})$ -stationary solution, where  $\theta_y = O(\kappa, \epsilon/\sqrt{\kappa})$ -stationary solution accuracy, and  $\theta_y = O(\kappa, \epsilon/\sqrt{\kappa})$ -stationary solution.

**Remark 4.4.** The actual learning rate of the primal variable is  $\beta_x \eta = O\left(\frac{K(1-\lambda)^2 \epsilon}{\kappa^{1/2}}\right)$ , and that of the dual variable is  $\beta_y \eta = O\left(\frac{K(1-\lambda)^2 \epsilon}{\kappa^{1/2}}\right)$ . Obviously, they are on the same scale in terms of condition number  $\kappa$ , solution accuracy  $\epsilon$ , and spectral gap  $1-\lambda$ . In addition, the constant batch size based methods, including DM-HSGD (Xian et al., 2021) and DM-GDA (Huang & Chen, 2023), use the learning rate for the primal variable in the order of  $O\left(\frac{K(1-\lambda)^2 \epsilon}{\kappa^3}\right)$  and that for the dual variable is  $O\left(\frac{K(1-\lambda)^2 \epsilon}{\kappa^1}\right)$ . Apparently, our algorithm can allow a larger learning rate. Moreover, when the number of workers K=1, the spectral gap becomes  $1-\lambda=1$ . Our learning rates  $O\left(\frac{\epsilon}{\kappa^{1/2}}\right)$  are larger than  $O\left(\frac{\epsilon^2}{\kappa^1}\right)$  of the single-machine method, Smoothed-SAGDA (Yang et al., 2022).

The primal—dual learning rate ratio is important because the loss function often exhibits distinct properties for the two variables. When the loss function is nonconvex in the primal variable but satisfies the PL condition in the dual variable, optimizing the primal variable becomes significantly more challenging, and a smaller learning rate is commonly used (See Table 1) to maintain stability and prevent divergence. In contrast, with the smoothing technique, the loss function becomes strongly convex in the primal variable and strongly concave in the dual variable, resulting a well-behaved loss landscape that permits a larger primal learning rate.

Remark 4.5. To compare the convergence rate of our algorithm in Corollary 4.3 with existing algorithms in Table 1, we need to translate the  $O(\epsilon,\epsilon/\sqrt{\kappa})$ -stationary solution to the  $O(\epsilon)$ -stationary solution. In particular, (Yang et al., 2022) shows that we can apply stochastic gradient descent ascent algorithm to the optimization problem:  $\min_{x\in\mathbb{R}^{d_1}}\max_{y\in\mathbb{R}^{d_2}}f(x,y)+L\|x-x'\|^2$ , where x' is the output of our Algorithm 1. Since this problem satisfies the PL condition in both x and y, the iteration complexity for the translation is in the order of  $\tilde{O}(\frac{1}{\epsilon^2})$ , which is apparently dominated by  $T=O\left(\frac{\kappa^{3/2}}{K(1-\lambda)^2\epsilon^3}\right)$ . Therefore, the iteration complexity to find the  $O(\epsilon)$ -stationary solution is still  $T=O\left(\frac{\kappa^{3/2}}{K(1-\lambda)^2\epsilon^3}\right)$ .

The proof structure and all technical details is provided in Appendix B-E.

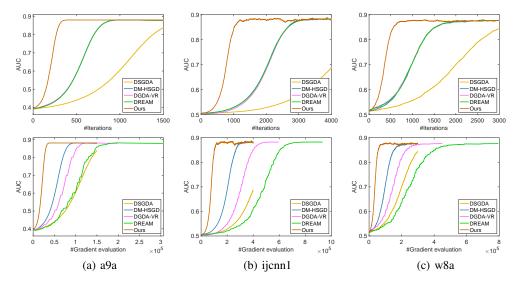


Figure 1: Test AUC vs. Iterations and Gradient Evaluations (Random Graph).

## 5 EXPERIMENTS

In this section, we conduct extensive experiments on AUC maximization, which is defined in Appendix A, to verify the performance of our Algorithm 1.

#### 5.1 EXPERIMENTAL SETTINGS

We employ three benchmark datasets: a9a, w8a, and ijcnn1, which can be found from LIBSVM Data website  $^3$ . In our experiments, 80% of samples are used as the training set, while the remaining 20% are used for testing. The training samples are randomly distributed across ten workers, where K=10 in our experiment. To evaluate the performance of our algorithm, we compare it with the state-of-the-art decentralized optimization algorithms: DSGDA (Tsaknakis et al., 2020), DM-HSGD  $^4$  (Xian et al., 2021), DGDA-VR (Zhang et al., 2024), and DREAM (Chen et al., 2024). Notably, for DSGDA, we use the stochastic gradient descent ascent instead of the full gradient as described in their paper. For DM-HSGD, the STORM gradient estimator is employed. DGDA-VR leverages the SPIDER gradient estimator in the stochastic setting, while DREAM utilizes the PAGE estimator.

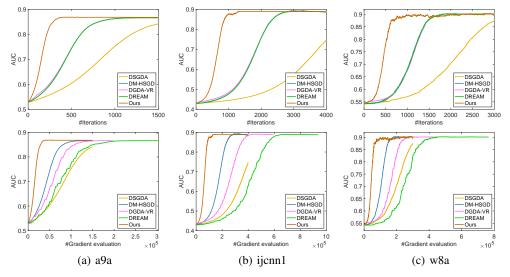


Figure 2: Test AUC vs Iterations and Gradient Evaluations (Line Graph).

<sup>3</sup>https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/

<sup>&</sup>lt;sup>4</sup>Note that DM-GDA is the same as DM-HSGD; they differ only in their convergence analysis under different assumptions.

Specifically, we consider two types of communication networks: 1) an Erdos-Renyi random graph with an edge probability of 0.5, and 2) a line communication network where each worker is connected to only two neighboring workers. Throughout all experiments, we fix the solution accuracy  $\epsilon$  at 0.01 and use a batch size b of 100. For the a9a and ijcnn1 datasets, the step size of all methods is set to 0.01. Specifically, in our method,  $\beta_x$ ,  $\beta_y$ ,  $\hat{\beta}_x$ , and  $\hat{\beta}_y$  are each set to 0.1, while  $\eta$  is set to 0.1, ensuring that their product equals 0.01. For the w8a dataset, the step size of all methods is set to 0.05. In this case,  $\beta_x$ ,  $\beta_y$ ,  $\hat{\beta}_x$ , and  $\hat{\beta}_y$  are each set to 0.5, while  $\eta$  remains 0.1, ensuring that their product equals 0.05. Moreover, according to the theoretical results of the baseline methods, the learning rate of the dual variable in DSGDA, DM-HSGD, and DGDA-VR is scaled by  $1/\kappa$ , while the learning rate of the primal variable is scaled by  $1/\kappa^3$ . For DREAM, scaling is 1 for the dual variable and  $1/\kappa^2$  for the primal variable. Both learning rates in our method are scaled by  $1/\kappa^{1/2}$ . In our experiments, we assume  $\kappa=1.5$ . Additionally, in our method,  $\gamma_1$  and  $\gamma_2$  are assigned a value of 0.01. For DM-HSGD, the coefficient of the STORM estimator is set to 0.01. Additionally, DGDA-VR computes the full gradient every 100 iterations, while for DREAM, the probability of the PAGE estimator is set to  $\frac{\sqrt{b}}{b\sqrt{K}}$ .

#### 5.2 EXPERIMENTAL RESULTS

For the random communication graph, we present test AUC versus the number of iterations and gradient evaluations in Figure 1. As shown in Figure 1, our algorithm achieves significantly faster convergence than all baseline methods in terms of the number of iterations, demonstrating its superior efficiency. Furthermore, Figure 1 also indicates that our method also converges more quickly when measured by the number of gradient evaluations, highlighting its lower sample complexity. Notably, DGDA-VR and DREAM incur significantly higher computational cost due to periodic full-gradient computation. These results underscore the efficacy of our algorithm in optimizing performance while maintaining computational efficiency. For the line communication graph, we also present test AUC versus the number of iterations and gradient evaluations in Figure 2. Our method continues to exhibit faster convergence compared to the baseline methods, further validating its effectiveness.

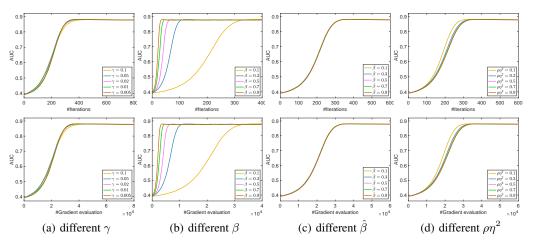


Figure 3: Test AUC under hyperparameters (Random Graph, a9a).

Finally, we evaluate the performance of our method under different values of  $\gamma$ ,  $\beta$ ,  $\hat{\beta}$ , and  $\rho\eta$  in Figure 3, where we set  $\gamma_x = \gamma_y = \gamma$ ,  $\beta_x = \beta_y = \beta$ ,  $\hat{\beta}_x = \hat{\beta}_y = \hat{\beta}$ , and  $\rho_x = \rho_y = \rho$ . Our method is robust to all hyperparameters except  $\beta$ , so they do not require fine-tuning. Since  $\beta$  only scales the learning rate, we fix its value, leaving the learning rate  $\eta$  as the only hyperparameter to tune.

#### 6 CONCLUSION

In this paper, we developed a novel decentralized minimax optimization algorithm based on the smoothing technique. In particular, our algorithm demonstrates how to incorporate the variance-reduced gradient in the presence of the auxiliary variable and how to perform communication for the auxiliary variable. Moreover, our algorithm can achieve a better dependence on the condition number than all existing methods, which confirms the significance of our algorithm. Finally, experimental results confirm the effectiveness of our algorithm.

**Ethics statement** This research complies with the ICLR Code of Ethics. The study is purely theoretical and methodological, and it does not involve human participants or personally identifiable information. The datasets used in this paper are publicly accessible sources. Our algorithm is designed for advancing machine learning research and does not raise foreseeable risks regarding safety, fairness, privacy, or security.

**Reproducibility statement** To facilitate reproducibility, we provide a comprehensive description of the algorithm and its underlying assumptions in Section 3, with complete theoretical analyses and proofs in Section 4 and Appendix B- E. Experimental details, including datasets, hyperparameter settings, and the communication graph used in decentralized network, are presented in Section 5. Upon acceptance, we will release the full source code to ensure that all reported results can be reliably reproduced.

**The Use of Large Language Models (LLMs)** LLMs were used only to aid in polishing the writing and improving readability of the manuscript. No part of the research ideation, algorithm design, analysis, or experimental results relied on LLMs.

## REFERENCES

- Lesi Chen, Boyuan Yao, and Luo Luo. Faster stochastic algorithms for minimax optimization under polyak-{\L} ojasiewicz condition. In *Advances in Neural Information Processing Systems*, 2022.
- Lesi Chen, Haishan Ye, and Luo Luo. An efficient stochastic algorithm for decentralized nonconvex-strongly-concave minimax optimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 1990–1998. PMLR, 2024.
- Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex sgd. *Advances in neural information processing systems*, 32, 2019.
- John C Duchi, Peter W Glynn, and Hongseok Namkoong. Statistics of robust optimization: A generalized empirical likelihood approach. *Mathematics of Operations Research*, 46(3):946–969, 2021.
- Cong Fang, Chris Junchi Li, Zhouchen Lin, and Tong Zhang. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. *Advances in neural information processing systems*, 31, 2018.
- Hongchang Gao. Decentralized stochastic gradient descent ascent for finite-sum minimax problems. *arXiv preprint arXiv:2212.02724*, 2022.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- Zhishuai Guo, Yi Xu, Wotao Yin, Rong Jin, and Tianbao Yang. On stochastic moving-average estimators for non-convex optimization. *arXiv preprint arXiv:2104.14840*, 2021.
- Feihu Huang and Songcan Chen. Near-optimal decentralized momentum method for nonconvex-pl minimax problems. *arXiv preprint arXiv:2304.10902*, 2023.
- Feihu Huang, Shangqian Gao, Jian Pei, and Heng Huang. Accelerated zeroth-order and first-order momentum methods from mini to minimax optimization. *Journal of Machine Learning Research*, 23(36):1–70, 2022.
- Yan Huang, Xiang Li, Yipeng Shen, Niao He, and Jinming Xu. Achieving near-optimal convergence for distributed minimax optimization with adaptive stepsizes. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=07IN4nsaIO.
- Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtárik. Page: A simple and optimal probabilistic gradient estimator for nonconvex optimization. In *International conference on machine learning*, pp. 6286–6295. PMLR, 2021.

- Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pp. 6083–6093. PMLR, 2020.
  - Luo Luo, Haishan Ye, Zhichao Huang, and Tong Zhang. Stochastic recursive gradient descent ascent for stochastic nonconvex-strongly-concave minimax problems. *arXiv* preprint arXiv:2001.03724, 2020.
    - Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
    - Lam M Nguyen, Jie Liu, Katya Scheinberg, and Martin Takáč. Sarah: A novel method for machine learning problems using stochastic recursive gradient. In *International Conference on Machine Learning*, pp. 2613–2621. PMLR, 2017.
    - Wenqiang Pu, Kaizhao Sun, and Jiawei Zhang. Smoothed proximal lagrangian method for nonlinear constrained programs. *arXiv preprint arXiv:2408.15047*, 2024.
    - Shuang Qiu, Zhuoran Yang, Xiaohan Wei, Jieping Ye, and Zhaoran Wang. Single-timescale stochastic nonconvex-concave optimization for smooth nonlinear td learning. *arXiv* preprint arXiv:2008.10103, 2020.
    - Wei Shen, Minhui Huang, Jiawei Zhang, and Cong Shen. Stochastic smoothed gradient descent ascent for federated minimax optimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 3988–3996. PMLR, 2024.
    - Ioannis Tsaknakis, Mingyi Hong, and Sijia Liu. Decentralized min-max optimization: Formulations, algorithms and applications in network poisoning attack. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5755–5759. IEEE, 2020.
    - Wenhan Xian, Feihu Huang, Yanfu Zhang, and Heng Huang. A faster decentralized algorithm for nonconvex minimax problems. *Advances in Neural Information Processing Systems*, 34, 2021.
    - Yangyang Xu. Decentralized gradient descent maximization method for composite nonconvex strongly-concave minimax problems. *SIAM Journal on Optimization*, 34(1):1006–1044, 2024.
    - Junchi Yang, Negar Kiyavash, and Niao He. Global convergence and variance-reduced optimization for a class of nonconvex-nonconcave minimax problems. *arXiv preprint arXiv:2002.09621*, 2020.
    - Junchi Yang, Antonio Orvieto, Aurelien Lucchi, and Niao He. Faster single-loop algorithms for minimax optimization without strong concavity. In *International Conference on Artificial Intelligence and Statistics*, pp. 5485–5517. PMLR, 2022.
    - Yiming Ying, Longyin Wen, and Siwei Lyu. Stochastic online auc maximization. *Advances in neural information processing systems*, 29:451–459, 2016.
    - Peiran Yu, Guoyin Li, and Ting Kei Pong. Kurdyka-łojasiewicz exponent via inf-projection. *Foundations of Computational Mathematics*, 22(4):1171–1217, 2022.
    - Jiawei Zhang, Peijun Xiao, Ruoyu Sun, and Zhiquan Luo. A single-loop smoothed gradient descent-ascent algorithm for nonconvex-concave min-max problems. *Advances in neural information processing systems*, 33:7377–7389, 2020.
    - Xin Zhang, Zhuqing Liu, Jia Liu, Zhengyuan Zhu, and Songtao Lu. Taming communication and sample complexities in decentralized policy evaluation for cooperative multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 34:18825–18838, 2021.
    - Xuan Zhang, Gabriel Mancino-Ball, Necdet Serhat Aybat, and Yangyang Xu. Jointly improving the sample and communication complexities in decentralized stochastic minimax optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 20865–20873, 2024.
    - Taoli Zheng, Linglingzhi Zhu, Anthony Man-Cho So, José Blanchet, and Jiajin Li. Universal gradient descent ascent method for nonconvex-nonconcave minimax optimization. *Advances in Neural Information Processing Systems*, 36:54075–54110, 2023.

## A AUC MAXIMIZATION

Specifically, we focus on the AUC maximization problem (Ying et al., 2016) for the binary classification task, which is formulated as the following minimax optimization problem (Note that we have included the smoothed term  $\gamma_1/2||x-\hat{x}||^2$ ,  $\gamma_2/2||y-\hat{y}||^2$ ):

$$\min_{x,\tilde{x}_{1},\tilde{x}_{2}} \max_{y} \frac{1}{K} \sum_{k=1}^{K} \frac{1}{n} \sum_{i=1}^{n} \left( (1-p)(x^{T} a_{i}^{(k)} - \tilde{x}_{1})^{2} \mathbb{I}_{[b_{i}^{(k)} = 1]} + 2(1+y) \left( px^{T} a_{i}^{(k)} \mathbb{I}_{[b_{i}^{(k)} = -1]} - (1-p)x^{T} a_{i}^{(k)} \mathbb{I}_{[b_{i}^{(k)} = 1]} \right) + p(x^{T} a_{i}^{(k)} - \tilde{x}_{2})^{2} \mathbb{I}_{[b_{i}^{(k)} = -1]} - p(1-p)y^{2} + \rho \sum_{i=1}^{d} \frac{x_{j}^{2}}{1+x_{j}^{2}} + \frac{\gamma_{1}}{2} \|x - \hat{x}\|^{2} - \frac{\gamma_{2}}{2} \|y - \hat{y}\|^{2} \right), \tag{8}$$

where  $x \in \mathbb{R}^d$  is the classifier's parameter,  $\tilde{x}_1 \in \mathbb{R}$ ,  $\tilde{x}_2 \in \mathbb{R}$ ,  $y \in \mathbb{R}$  are the parameters to compute the AUC loss,  $\hat{x}$  and  $\hat{y}$  are the auxiliary variables.  $(a_i^{(k)}, b_i^{(k)})$  is the i-th sample's feature and label on the k-th worker, p is the prior probability of positive class,  $\mathbb{I}$  is an indicator function,  $\rho$  is a hyperparameter for the regularization term, and  $\gamma_1 > 0$ ,  $\gamma_2 > 0$  are hyperparameters for the auxiliary variable. In our experiments, we set  $\rho$  to 0.001. Notably, this optimization problem satisfies the nonconvex-PL optimization problem, which can be efficiently solved using our proposed Algorithm 1.

## B THE STRUCTURE OF THE PROOF FOR THEOREM 4.2

To make our proof easy to follow, we provide an overview diagram in Figure 4.

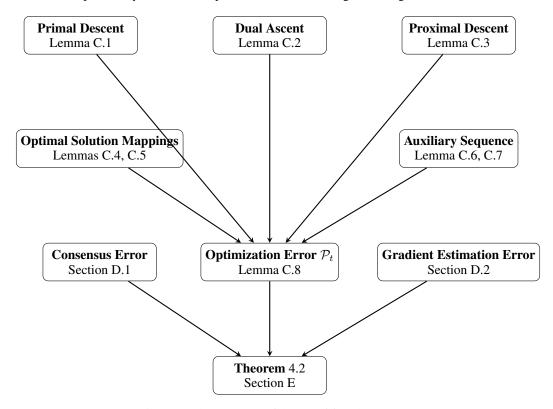


Figure 4: The structure of the proof for Theorem 4.2

It is worth noting that the STORM gradient estimator is a biased gradient estimator, so existing convergence analyzes based on the deterministic gradient (Zhang et al., 2020; Zheng et al., 2023)

and the unbiased gradient estimator (Yang et al., 2022) cannot be applied directly to our algorithm. Moreover, most existing stochastic smoothing methods typically apply smoothing only to the primal variable, which makes their analysis insufficient for our algorithm.

In Figure 4, there are actually two key components in our proof: 1) the optimization error related to doubly smoothing, 2) the consensus error and the gradient estimation error related to the decentralized setting. In Section C, we provide the lemmas for bounding the optimization error. This includes:

- descent-ascent update lemmas (Lemma C.1, Lemma C.2, Lemma C.3);
- optimal solution mappings (Lemma C.4, Lemma C.5);
- auxiliary sequences (Lemma C.6, Lemma C.7).

These results are used in a potential function as Eq.(51):

$$\mathcal{P}_t = \mathbb{E}[F(\bar{x}_t, \bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t)] - 2\mathbb{E}[F_d(\bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t)] + 2\mathbb{E}[q(\bar{\hat{x}}_t)],$$

to establish the overall optimization error bound  $\mathcal{P}_{t+1} - \mathcal{P}_t$  in Lemma C.8. It is worth noted that Lemma C.8 demonstrates that optimization error is affected by the consensus error caused by the decentralized setting and gradient estimation errors. Therefore, in Section D, we address two types of error in the decentralized setting:

- the consensus error, including that of auxiliary variables introduced by smoothing (Section D.1);
- the gradient estimation error from the STORM update (Section D.2).

After establishing all supporting lemmas, we proceed to derive the convergence rate through a novel potential function  $\mathcal{L}_t$ , which intergrates the optimization error in Lemma C.8 and the consensus error and gradient estimation error together as follows:

$$\begin{split} \mathcal{L}_{t} &= \underbrace{\mathcal{P}_{t}}_{\text{optimization error}} + c_{1} \underbrace{\mathbb{E}[\|\frac{1}{K} \sum_{k=1}^{K} u_{t}^{(k)} - \frac{1}{K} \sum_{k=1}^{K} \nabla_{x} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]}_{\text{gradient estimation error}} \\ &+ c_{2} \underbrace{\mathbb{E}[\|\frac{1}{K} \sum_{k=1}^{K} v_{t}^{(k)} - \frac{1}{K} \sum_{k=1}^{K} \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]}_{\text{gradient estimation error}} \\ &+ c_{3} \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{x}_{t} - x_{t}^{(k)}\|^{2}] + c_{4} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{y}_{t} - y_{t}^{(k)}\|^{2}] + c_{5} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{x}_{t} - \hat{x}_{t}^{(k)}\|^{2}]} \\ &+ c_{10} \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{y}_{t} - \hat{y}_{t}^{(k)}\|^{2}] + c_{6} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{p}_{t} - p_{t}^{(k)}\|^{2}] + c_{7} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{q}_{t} - q_{t}^{(k)}\|^{2}]} \\ &+ c_{8} \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|u_{t}^{(k)} - \nabla_{x} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]}_{\text{gradient estimation error}} \\ &+ c_{9} \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]}_{\text{gradient estimation error}} \\ &+ c_{9} \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]}_{\text{gradient estimation error}} \\ &+ c_{9} \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]}_{\text{gradient estimation error}} \\ &+ \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]}_{\text{gradient estimation error}} \\ &+ \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}}_{\text{gradient estimation error}} \\ &+ \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}}_{\text{gradient estimation error}} \\ &+ \underbrace{\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{$$

By selecting appropriate hyperparameters, as detailed in Section E, we establish the convergence guarantee stated in Theorem 4.2. The construction of this proof framework is both technically intricate and conceptually non-trivial, underscoring the novelty and difficulty of our analysis.

#### **B.1** TERMINOLOGIES

To establish the convergence rate of Algorithm 1, we introduce the following symbols:

$$X_{t} = [x_{t}^{(1)}, x_{t}^{(2)}, \cdots, x_{t}^{(K)}] \in \mathbb{R}^{d_{1} \times K}, \ \tilde{X}_{t} = [\tilde{x}_{t}^{(1)}, \tilde{x}_{t}^{(2)}, \cdots, \tilde{x}_{t}^{(K)}] \in \mathbb{R}^{d_{1} \times K},$$

$$Y_{t} = [y_{t}^{(1)}, y_{t}^{(2)}, \cdots, y_{t}^{(K)}] \in \mathbb{R}^{d_{2} \times K}, \ \tilde{Y}_{t} = [\tilde{y}_{t}^{(1)}, \tilde{y}_{t}^{(2)}, \cdots, \tilde{y}_{t}^{(K)}] \in \mathbb{R}^{d_{2} \times K},$$

$$\hat{X}_{t} = [\hat{x}_{t}^{(1)}, \hat{x}_{t}^{(2)}, \cdots, \hat{x}_{t}^{(K)}] \in \mathbb{R}^{d_{1} \times K}, \ \tilde{X}_{t} = [\tilde{x}_{t}^{(1)}, \tilde{x}_{t}^{(2)}, \cdots, \tilde{x}_{t}^{(K)}] \in \mathbb{R}^{d_{1} \times K},$$

$$\hat{Y}_{t} = [\hat{y}_{t}^{(1)}, \hat{y}_{t}^{(2)}, \cdots, \hat{y}_{t}^{(K)}] \in \mathbb{R}^{d_{2} \times K}, \ \tilde{Y}_{t} = [\tilde{y}_{t}^{(1)}, \tilde{y}_{t}^{(2)}, \cdots, \tilde{y}_{t}^{(K)}] \in \mathbb{R}^{d_{2} \times K},$$

$$U_{t} = [u_{t}^{(1)}, u_{t}^{(2)}, \cdots, u_{t}^{(K)}] \in \mathbb{R}^{d_{1} \times K}, \ V_{t} = [v_{t}^{(1)}, v_{t}^{(2)}, \cdots, v_{t}^{(K)}] \in \mathbb{R}^{d_{2} \times K},$$

$$P_{t} = [p_{t}^{(1)}, p_{t}^{(2)}, \cdots, p_{t}^{(K)}] \in \mathbb{R}^{d_{1} \times K}, \ Q_{t} = [q_{t}^{(1)}, q_{t}^{(2)}, \cdots, q_{t}^{(K)}] \in \mathbb{R}^{d_{2} \times K},$$

$$\bar{X}_{t} = \frac{1}{K} X_{t} \mathbf{1} \mathbf{1}^{T}, \ \bar{Y}_{t} = \frac{1}{K} Y_{t} \mathbf{1} \mathbf{1}^{T}, \ \bar{X}_{t} = \frac{1}{K} \hat{X}_{t} \mathbf{1} \mathbf{1}^{T}, \ \bar{Y}_{t} = \frac{1}{K} \hat{Y}_{t} \mathbf{1} \mathbf{1}^{T},$$

$$\bar{U}_{t} = \frac{1}{K} U_{t} \mathbf{1} \mathbf{1}^{T}, \ \bar{V}_{t} = \frac{1}{K} V_{t} \mathbf{1} \mathbf{1}^{T}, \ \bar{P}_{t} = \frac{1}{K} P_{t} \mathbf{1} \mathbf{1}^{T}, \ \bar{Q}_{t} = \frac{1}{K} Q_{t} \mathbf{1} \mathbf{1}^{T},$$

$$(9)$$

where  $\mathbf{1} = [1, 1, \dots, 1]^T \in \mathbb{R}^K$ . Based on these terminologies, the update of  $x, y, \hat{x}, \hat{y}, p$ , and q in Algorithm 1 is represented as follows:

$$\tilde{X}_{t+1} = X_t W - \beta_x P_t , X_{t+1} = X_t + \eta (\tilde{X}_{t+1} - X_t) , 
\tilde{Y}_{t+1} = Y_t W + \beta_y Q_t , Y_{t+1} = Y_t + \eta (\tilde{Y}_{t+1} - Y_t) , 
\tilde{X}_{t+1} = \hat{X}_t W + \hat{\beta}_x (X_{t+1} - \hat{X}_t) , \hat{X}_{t+1} = \hat{X}_t + \eta (\tilde{X}_{t+1} - \hat{X}_t) , 
\tilde{Y}_{t+1} = \hat{Y}_t W + \hat{\beta}_y (Y_{t+1} - \hat{Y}_t) , \hat{Y}_{t+1} = \hat{Y}_t + \eta (\tilde{Y}_{t+1} - \hat{Y}_t) , 
P_{t+1} = P_t W + U_{t+1} - U_t , Q_{t+1} = Q_t W + V_{t+1} - V_t , 
\bar{X}_{t+1} = \bar{X}_t - \beta_x \eta \bar{U}_t , \bar{Y}_{t+1} = \bar{Y}_t + \beta_y \eta \bar{V}_t , 
\tilde{X}_{t+1} = \tilde{X}_t + \hat{\beta}_x \eta (\bar{X}_{t+1} - \tilde{X}_t) , \bar{Y}_{t+1} = \tilde{Y}_t + \hat{\beta}_y \eta (\bar{Y}_{t+1} - \bar{Y}_t) .$$
(10)

Note that  $\bar{P}_t = \bar{U}_t$  and  $\bar{Q}_t = \bar{V}_t$ .

Moreover, following (Yang et al., 2022; Zheng et al., 2023), we introduce the following auxiliary functions and variables for convergence analysis:

$$F_{d}(y; \hat{x}, \hat{y}) = \min_{x \in \mathbb{R}^{d_1}} F(x, y; \hat{x}, \hat{y}) , \quad \text{dual function}$$

$$F_{p}(x; \hat{x}, \hat{y}) = \max_{y \in \mathbb{R}^{d_2}} F(x, y; \hat{x}, \hat{y}) , \quad \text{primal function}$$

$$g(\hat{x}, \hat{y}) = \min_{x \in \mathbb{R}^{d_1}} \max_{y \in \mathbb{R}^{d_2}} F(x, y; \hat{x}, \hat{y}) ,$$

$$p(\hat{y}) = \min_{\hat{x} \in \mathbb{R}^{d_1}} \sup_{y \in \mathbb{R}^{d_2}} F(x, y; \hat{x}, \hat{y}) ,$$

$$x^*(y; \hat{x}, \hat{y}) = \arg\min_{x \in \mathbb{R}^{d_1}} F(x, y; \hat{x}, \hat{y}) ,$$

$$y^*(x; \hat{x}, \hat{y}) = \arg\max_{y \in \mathbb{R}^{d_2}} F(x, y; \hat{x}, \hat{y}) ,$$

$$x^*(\hat{x}, \hat{y}) \triangleq x^*(y^*(\hat{x}, \hat{y}); \hat{x}, \hat{y}) = \arg\min_{x \in \mathbb{R}^{d_1}} F_{p}(x; \hat{x}, \hat{y}) ,$$

$$y^*(\hat{x}, \hat{y}) \triangleq y^*(x^*(\hat{x}, \hat{y}); \hat{x}, \hat{y}) = \arg\max_{y \in \mathbb{R}^{d_2}} F_{d}(y; \hat{x}, \hat{y}) ,$$

$$\hat{x}^*(\hat{y}) = \arg\min_{\hat{x} \in \mathbb{R}^{d_1}} g(\hat{x}, \hat{y}) , \quad \hat{y}^*(\hat{x}) = \arg\max_{\hat{y} \in \mathbb{R}^{d_2}} g(\hat{x}, \hat{y}) ,$$

$$y^+(\hat{x}_t, \hat{y}_t) = y_t + \beta_y \eta \nabla_y F_{d}(y_t; \hat{x}_t, \hat{y}_t) ,$$

$$\hat{y}^+(\hat{x}_{t+1}) = \hat{y}_t + \hat{\beta}_y \eta (y^*(\hat{x}_t, \hat{y}_t) - \hat{y}_t) . \quad (11)$$

#### **B.2** Function Properties

**Lemma B.1.** (Zheng et al., 2023) Given Assumptions 3.1-3.4, then  $F(x,y;\hat{x},\hat{y})$  is  $(\gamma_1 + L)$ -smooth and  $(\gamma_1 - L)$ -strongly convex with respect to x.  $F(x,y;\hat{x},\hat{y})$  is  $(\gamma_2 + L)$ -smooth and  $(\gamma_2 - L)$ -strongly concave with respect to y.

**Lemma B.2.** (Zheng et al., 2023) Given Assumptions 3.1-3.4, the following inequality holds:

$$\begin{split} &\|x^*(y_1;\hat{x},\hat{y}) - x^*(y_2;\hat{x},\hat{y})\| \leq C_{x_{y\hat{x}\hat{y}}^1} \|y_1 - y_2\| \;, \\ &\|x^*(y;\hat{x}_1,\hat{y}) - x^*(y;\hat{x}_2,\hat{y})\| \leq C_{x_{y\hat{x}\hat{y}}^2} \|\hat{x}_1 - \hat{x}_2\| \;, \\ &\|x^*(\hat{x}_1,\hat{y}) - x^*(\hat{x}_2,\hat{y})\| \leq C_{x_{\hat{x}\hat{y}}^1} \|\hat{x}_1 - \hat{x}_2\| \;, \\ &\|y^*(x_1;\hat{x},\hat{y}) - y^*(x_2;\hat{x},\hat{y})\| \leq C_{y_{x\hat{x}\hat{y}}^1} \|x_1 - x_2\| \;, \\ &\|y^*(x;\hat{x},\hat{y}_1) - y^*(x;\hat{x},\hat{y}_2)\| \leq C_{y_{x\hat{x}\hat{y}}^3} \|\hat{y}_1 - \hat{y}_2\| \;, \\ &\|y^*(\hat{x}_1,\hat{y}) - y^*(\hat{x}_2,\hat{y})\| \leq C_{y_{x\hat{y}}^1} \|\hat{x}_1 - \hat{x}_2\| \;, \\ &\|y^*(\hat{x},\hat{y}_1) - y^*(\hat{x},\hat{y}_2)\| \leq C_{y_{x\hat{y}}^2} \|\hat{y}_1 - \hat{y}_2\| \;, \end{split}$$

where

$$\begin{split} &C_{x_{y\hat{x}\hat{y}}^{1}} = \frac{\gamma_{1}}{\gamma_{1} - L} \;, \quad C_{x_{y\hat{x}\hat{y}}^{2}} = \frac{\gamma_{1}}{\gamma_{1} - L} \;, \quad C_{x_{\hat{x}\hat{y}}^{1}} = \frac{\gamma_{1}}{\gamma_{1} - L} \;, \\ &C_{y_{x\hat{x}\hat{y}}^{1}} = \frac{\gamma_{2}}{\gamma_{2} - L} \;, \quad C_{y_{x\hat{x}\hat{y}}^{3}} = \frac{\gamma_{2}}{\gamma_{2} - L} \;, \quad C_{y_{\hat{x}\hat{y}}^{1}} = \frac{\gamma_{1}}{\gamma_{2} - L} C_{x_{y\hat{x}\hat{y}}^{1}} + 1 \;, C_{y_{\hat{x}\hat{y}}^{2}} = \frac{\gamma_{2}}{\gamma_{2} - L} \;. \end{split} \tag{13}$$

**Lemma B.3.** (Zheng et al., 2023) Given Assumptions 3.1-3.4, then  $F_d(y; \hat{x}, \hat{y})$  is  $L_d$ -smooth, where  $L_d = LC_{x_{n\hat{x}\hat{y}}} + L + \gamma_2$ .

**Lemma B.4.** Given Assumptions 3.1-3.4, by defining  $y^+(\hat{x}_t, \hat{y}_t) = y_t + \beta_y \eta \nabla_y F_d(y_t; \hat{x}_t, \hat{y}_t)$ , the following inequality holds:

$$||y_t - y^*(\hat{x}_t, \hat{y}_t)|| \le \frac{1}{\beta_y \eta(\gamma_2 - L)} ||y_t - y^+(\hat{x}_t, \hat{y}_t)||.$$
(14)

*Proof.* Due to  $y^*(\hat{x}_t, \hat{y}_t) = \arg\max_{y \in \mathbb{R}^{d_2}} F_d(y; \hat{x}_t, \hat{y}_t)$ , for any  $y \in \mathbb{R}^{d_2}$ , we have

$$\langle y - y^*(\hat{x}_t, \hat{y}_t), \nabla_y F_d(y^*(\hat{x}_t, \hat{y}_t); \hat{x}_t, \hat{y}_t) \rangle \le 0.$$
 (15)

By taking  $y = y_t$ , we have

$$\langle y_t - y^*(\hat{x}_t, \hat{y}_t), \nabla_y F_d(y^*(\hat{x}_t, \hat{y}_t); \hat{x}_t, \hat{y}_t) \rangle \le 0.$$
 (16)

In addition, because  $F_d(y; \hat{x}, \hat{y})$  is  $(\gamma_2 - L)$ -strongly concave with respect to y, we have

$$\langle y_t - y^*(\hat{x}_t, \hat{y}_t), \nabla_y F_d(y_t; \hat{x}_t, \hat{y}_t) - \nabla_y F_d(y^*(\hat{x}_t, \hat{y}_t); \hat{x}_t, \hat{y}_t) \rangle + (\gamma_2 - L) \|y_t - y^*(\hat{x}_t, \hat{y}_t)\|^2 \le 0.$$
(17)

By combining the above two inequalities, we have

$$\langle y_t - y^*(\hat{x}_t, \hat{y}_t), \nabla_y F_d(y_t; \hat{x}_t, \hat{y}_t) \rangle + (\gamma_2 - L) \|y_t - y^*(\hat{x}_t, \hat{y}_t)\|^2 \le 0.$$
 (18)

Then, we can obtain

$$(\gamma_{2} - L)\|y_{t} - y^{*}(\hat{x}_{t}, \hat{y}_{t})\|^{2} \leq \langle y^{*}(\hat{x}_{t}, \hat{y}_{t}) - y_{t}, \nabla_{y} F_{d}(y_{t}; \hat{x}_{t}, \hat{y}_{t}) \rangle$$

$$\leq \|y_{t} - y^{*}(\hat{x}_{t}, \hat{y}_{t})\|\|\nabla_{y} F_{d}(y_{t}; \hat{x}_{t}, \hat{y}_{t})\| = \|y_{t} - y^{*}(\hat{x}_{t}, \hat{y}_{t})\|\|\frac{y^{+}(\hat{x}_{t}, \hat{y}_{t}) - y_{t}}{\beta_{y}\eta}\|.$$
(19)

As a result, we have

$$||y_t - y^*(\hat{x}_t, \hat{y}_t)|| \le \frac{1}{\beta_u \eta(\gamma_2 - L)} ||y^+(\hat{x}_t, \hat{y}_t) - y_t||.$$
 (20)

**Lemma B.5.** Given Assumptions 3.1-3.4, then

$$||x_t - x^*(y_t; \hat{x}_t, \hat{y}_t)|| \le \frac{1}{\gamma_1 - L} ||\nabla_x F(x_t, y_t; \hat{x}_t, \hat{y}_t)||.$$
(21)

*Proof.* Due to  $x^*(y_t; \hat{x}_t, \hat{y}_t) = \arg\min_{x \in \mathbb{R}^{d_1}} F(x, y_t; \hat{x}_t, \hat{y}_t)$ , for any  $x \in \mathbb{R}^{d_1}$ , we have

$$\langle x - x^*(y_t; \hat{x}_t, \hat{y}_t), -\nabla_x F(x^*(y_t; \hat{x}_t, \hat{y}_t), y_t; \hat{x}_t, \hat{y}_t) \rangle \le 0.$$
 (22)

By taking  $x = x_t$ , we have

$$\langle x_t - x^*(y_t; \hat{x}_t, \hat{y}_t), \nabla_x F(x^*(y_t; \hat{x}_t, \hat{y}_t), y_t; \hat{x}_t, \hat{y}_t) \rangle \ge 0.$$
 (23)

In addition, because  $F(x, y; \hat{x}, \hat{y})$  is  $(\gamma_1 - L)$ -strongly convex with respect to x, we have

$$\langle x_t - x^*(y_t; \hat{x}_t, \hat{y}_t), \nabla_x F(x_t, y_t; \hat{x}_t, \hat{y}_t) - \nabla_x F(x^*(y_t; \hat{x}_t, \hat{y}_t), y_t; \hat{x}_t, \hat{y}_t) \rangle$$

$$\geq (\gamma_1 - L) \|x_t - x^*(y_t; \hat{x}_t, \hat{y}_t)\|^2 . \tag{24}$$

By combing the above two inequalities, we have

$$(\gamma_1 - L) \|x_t - x^*(y_t; \hat{x}_t, \hat{y}_t)\|^2 \le \langle x_t - x^*(y_t; \hat{x}_t, \hat{y}_t), \nabla_x F(x_t, y_t; \hat{x}_t, \hat{y}_t) \rangle$$

$$\le \|x_t - x^*(y_t; \hat{x}_t, \hat{y}_t)\| \|\nabla_x F(x_t, y_t; \hat{x}_t, \hat{y}_t)\|.$$
(25)

As a result, we have

$$||x_t - x^*(y_t; \hat{x}_t, \hat{y}_t)|| \le \frac{1}{\gamma_1 - L} ||\nabla_x F(x_t, y_t; \hat{x}_t, \hat{y}_t)||.$$
 (26)

## C OPTIMIZATION ERRORS

**Lemma C.1.** Given Assumptions 3.1-3.4 and  $\eta \leq \frac{1}{2\beta_x(\gamma_1 + L)}$ , the following inequality holds:

$$\mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t+1})] - \mathbb{E}[F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] \\
\leq -\frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] + \frac{\beta_{y}\eta}{2} \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] + (4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2} - \frac{\beta_{x}\eta}{4})\mathbb{E}[\|\bar{u}_{t}\|^{2}] \\
+ (\frac{3\beta_{y}\eta}{4} + \frac{\beta_{y}^{2}\eta^{2}(\gamma_{2} + L)}{2})\mathbb{E}[\|\bar{v}_{t}\|^{2}] + \frac{\beta_{x}\eta}{2}\mathbb{E}[\|\nabla_{x}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{u}_{t}\|^{2}] \\
- \frac{\gamma_{1}(2 - \hat{\beta}_{x}\eta)}{2\hat{\beta}_{x}\eta}\mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] - \frac{\gamma_{2}(\hat{\beta}_{y}\eta - 2)}{2\hat{\beta}_{y}\eta}\mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}].$$
(27)

*Proof.* Because  $F(x, y; \hat{x}, \hat{y})$  is  $(L + \gamma_1)$ -smooth with respect to x, we have

$$\mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] \\
\leq \mathbb{E}[F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] + \mathbb{E}[\langle \nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{x}_{t+1} - \bar{x}_{t} \rangle] + \frac{L + \gamma_{1}}{2} \mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] \\
= \mathbb{E}[F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] - \beta_{x} \eta \mathbb{E}[\langle \nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{u}_{t} \rangle] + \frac{\beta_{x}^{2} \eta^{2} (L + \gamma_{1})}{2} \mathbb{E}[\|\bar{u}_{t}\|^{2}] \\
= \mathbb{E}[F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] - \frac{\beta_{x} \eta}{2} \mathbb{E}[\|\nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] - \frac{\beta_{x} \eta}{2} \mathbb{E}[\|\bar{u}_{t}\|^{2}] \\
+ \frac{\beta_{x} \eta}{2} \mathbb{E}[\|\nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{u}_{t}\|^{2}] + \frac{\beta_{x}^{2} \eta^{2} (L + \gamma_{1})}{2} \mathbb{E}[\|\bar{u}_{t}\|^{2}] \\
\leq \mathbb{E}[F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] - \frac{\beta_{x} \eta}{2} \mathbb{E}[\|\nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] - \frac{\beta_{x} \eta}{4} \mathbb{E}[\|\bar{u}_{t}\|^{2}] \\
+ \frac{\beta_{x} \eta}{2} \mathbb{E}[\|\nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{u}_{t}\|^{2}], \tag{28}$$

where the last step holds due to  $\eta \leq \frac{1}{2\beta_x(\gamma_1 + L)}$ .

 In addition, because  $F(x, y; \hat{x}, \hat{y})$  is  $(L + \gamma_1)$ -smooth with respect to y, we have

$$\mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t})] \\
\leq \mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] + \mathbb{E}[\langle \nabla_{y} F(\bar{x}_{t+1}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{y}_{t+1} - \bar{y}_{t} \rangle] + \frac{\gamma_{2} + L}{2} \mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}] \\
= \mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] + \beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F(\bar{x}_{t+1}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \\
+ \beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] + \frac{\beta_{y}^{2} \eta^{2} (\gamma_{2} + L)}{2} \mathbb{E}[\|\bar{v}_{t}\|^{2}] \\
\leq \mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] + 4\beta_{y} \eta \mathbb{E}[\|\nabla_{y} F(\bar{x}_{t+1}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] + \frac{\beta_{y} \eta}{4} \mathbb{E}[\|\bar{v}_{t}\|^{2}] \\
+ \frac{\beta_{y} \eta}{2} \mathbb{E}[\|\nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] + \frac{\beta_{y} \eta}{2} \mathbb{E}[\|\bar{v}_{t}\|^{2}] + \frac{\beta_{y}^{2} \eta^{2} (\gamma_{2} + L)}{2} \mathbb{E}[\|\bar{v}_{t}\|^{2}] \\
\leq \mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] + \frac{\beta_{y} \eta}{2} \mathbb{E}[\|\nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
\leq \mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] + \frac{\beta_{y} \eta}{2} \mathbb{E}[\|\nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
+ 4\beta_{y} \eta \beta_{x}^{2} \eta^{2} L^{2} \mathbb{E}[\|\bar{u}_{t}\|^{2}] + \left(\frac{3\beta_{y} \eta}{4} + \frac{\beta_{y}^{2} \eta^{2} (\gamma_{2} + L)}{2}\right) \mathbb{E}[\|\bar{v}_{t}\|^{2}], \tag{29}$$

where the last step holds due to the following inequality.

$$\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t+1},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t}) - \nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t})\|^{2}] 
= \mathbb{E}[\|\nabla_{y}f(\bar{x}_{t+1},\bar{y}_{t}) + \gamma_{2}(\bar{y}_{t} - \bar{\hat{y}}_{t}) - \nabla_{y}f(\bar{x}_{t},\bar{y}_{t}) - \gamma_{2}(\bar{y}_{t} - \bar{\hat{y}}_{t})\|^{2}] 
\leq L^{2}\mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] \leq \beta_{x}^{2}\eta^{2}L^{2}\mathbb{E}[\|\bar{u}_{t}\|^{2}].$$
(30)

By combining Eq. (28) and Eq. (29), we have

$$\mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t})] \\
\leq \mathbb{E}[F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] - \frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] + \frac{\beta_{y}\eta}{2} \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
+ \frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{u}_{t}\|^{2}] \\
+ (4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2} - \frac{\beta_{x}\eta}{4})\mathbb{E}[\|\bar{u}_{t}\|^{2}] + (\frac{3\beta_{y}\eta}{4} + \frac{\beta_{y}^{2}\eta^{2}(\gamma_{2} + L)}{2})\mathbb{E}[\|\bar{v}_{t}\|^{2}] . \tag{31}$$

Moreover, according to the definition of  $F(x, y; \hat{x}, \hat{y})$ , we have

$$F(\bar{x}_{t+1}, \bar{y}_{t+1}; \hat{\bar{x}}_t, \bar{y}_t) - F(\bar{x}_{t+1}, \bar{y}_{t+1}; \hat{\bar{x}}_{t+1}, \bar{y}_t)$$

$$= f(\bar{x}_{t+1}, \bar{y}_{t+1}) + \frac{\gamma_1}{2} \|\bar{x}_{t+1} - \hat{\bar{x}}_t\|^2 - \frac{\gamma_2}{2} \|\bar{y}_{t+1} - \hat{\bar{y}}_t\|^2$$

$$- f(\bar{x}_{t+1}, \bar{y}_{t+1}) - \frac{\gamma_1}{2} \|\bar{x}_{t+1} - \hat{\bar{x}}_{t+1}\|^2 + \frac{\gamma_2}{2} \|\bar{y}_{t+1} - \bar{\hat{y}}_t\|^2$$

$$= \frac{\gamma_1}{2} (\|\bar{x}_{t+1} - \hat{\bar{x}}_t\|^2 - \|\bar{x}_{t+1} - \hat{\bar{x}}_{t+1}\|^2)$$

$$= \frac{\gamma_1}{2} (\|\bar{x}_{t+1} - \hat{\bar{x}}_t\|^2 - \|(1 - \hat{\beta}_x \eta)(\bar{x}_{t+1} - \hat{\bar{x}}_t)\|^2)$$

$$= \frac{\gamma_1(1 - (1 - \hat{\beta}_x \eta)^2)}{2} \|\bar{x}_{t+1} - \hat{\bar{x}}_t\|^2$$

$$= \frac{\gamma_1(1 - (1 - \hat{\beta}_x \eta)^2)}{2\hat{\beta}_x^2 \eta^2} \|\bar{x}_{t+1} - \hat{\bar{x}}_t\|^2$$

$$= \frac{\gamma_1(2 - \hat{\beta}_x \eta)}{2\hat{\beta}_x \eta} \|\bar{x}_{t+1} - \hat{\bar{x}}_t\|^2, \qquad (32)$$

where the third and fifth steps hold due to  $\bar{\hat{x}}_{t+1} = \bar{\hat{x}}_t + \hat{\beta}_x \eta(\bar{x}_{t+1} - \bar{\hat{x}}_t)$ .

Similarly, we have

$$F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_t) - F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t+1})$$

$$= f(\bar{x}_{t+1}, \bar{y}_{t+1}) + \frac{\gamma_1}{2} \|\bar{x}_{t+1} - \bar{\hat{x}}_{t+1}\|^2 - \frac{\gamma_2}{2} \|\bar{y}_{t+1} - \bar{\hat{y}}_t\|^2$$

$$- f(\bar{x}_{t+1}, \bar{y}_{t+1}) - \frac{\gamma_1}{2} \|\bar{x}_{t+1} - \bar{\hat{x}}_{t+1}\|^2 + \frac{\gamma_2}{2} \|\bar{y}_{t+1} - \bar{\hat{y}}_{t+1}\|^2$$

$$= \frac{\gamma_2}{2} \|\bar{y}_{t+1} - \bar{\hat{y}}_{t+1}\|^2 - \frac{\gamma_2}{2} \|\bar{y}_{t+1} - \bar{\hat{y}}_t\|^2$$

$$= \frac{\gamma_2(\hat{\beta}_y \eta - 2)}{2\hat{\beta}_y \eta} \|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_t\|^2.$$
(33)

By combining the above three inequalities, we have

$$\mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t+1})] - \mathbb{E}[F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] \\
= \mathbb{E}[F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t+1}) - F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) \\
+ \left(F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t})\right) + \left(F(\bar{x}_{t+1}, \bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t}) - F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\right) \\
\leq -\frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] + \frac{\beta_{y}\eta}{2} \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
+ \frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{u}_{t}\|^{2}] \\
+ (4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2} - \frac{\beta_{x}\eta}{4})\mathbb{E}[\|\bar{u}_{t}\|^{2}] + (\frac{3\beta_{y}\eta}{4} + \frac{\beta_{y}^{2}\eta^{2}(\gamma_{2} + L)}{2})\mathbb{E}[\|\bar{v}_{t}\|^{2}] \\
- \frac{\gamma_{1}(2 - \hat{\beta}_{x}\eta)}{2\hat{\beta}_{x}\eta} \mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] - \frac{\gamma_{2}(\hat{\beta}_{y}\eta - 2)}{2\hat{\beta}_{y}\eta} \mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}] . \tag{34}$$

**Lemma C.2.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\mathbb{E}[F_{d}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t+1})] - \mathbb{E}[F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] \\
\geq \beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] - \frac{\beta_{y}^{2} \eta^{2} L_{d}}{2} \mathbb{E}[\|\bar{v}_{t}\|^{2}] + \frac{\gamma_{2} (2 - \hat{\beta}_{y} \eta)}{2 \hat{\beta}_{y} \eta} \mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}] \\
+ \frac{\gamma_{1}}{2} \mathbb{E}[\langle \bar{x}_{t+1} - \bar{x}_{t}, \bar{x}_{t+1} + \bar{x}_{t} - 2x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) \rangle] .$$
(35)

*Proof.* According to the definition of  $F_d(y; \hat{x}, \hat{y})$ , we have

$$F_{d}(\bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{y}_{t+1}) - F_{d}(\bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t})$$

$$= F(x^{*}(\bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t+1}), \bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t+1}) - F(x^{*}(\bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t}), \bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t})$$

$$\geq F(x^{*}(\bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t+1}), \bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t+1}) - F(x^{*}(\bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t+1}), \bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_{t})$$

$$= \frac{\gamma_{2}}{2} (\|\bar{y}_{t+1} - \bar{\hat{y}}_{t}\|^{2} - \|\bar{y}_{t+1} - \bar{\hat{y}}_{t+1}\|^{2})$$

$$= \frac{\gamma_{2}(2 - \hat{\beta}_{y}\eta)}{2\hat{\beta}_{y}\eta} \|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_{t}\|^{2}, \qquad (36)$$

where the second step holds due to  $x^*(\bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_t) = \arg\min_{x \in \mathbb{R}^{d_1}} F(x, \bar{y}_{t+1}; \bar{\hat{x}}_{t+1}, \bar{\hat{y}}_t)$ , the last step holds as Eq. (33).

In addition, according to the definition of  $F_d(y; \hat{x}, \hat{y})$ , we have

$$F_{d}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - F_{d}(\bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t})$$

$$= F(x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}), \bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - F(x^{*}(\bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t}), \bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t})$$

$$\geq F(x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}), \bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - F(x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}), \bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t})$$

$$= \frac{\gamma_{1}}{2} (\|x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - \bar{x}_{t+1}\|^{2} - \|x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - \bar{x}_{t}\|^{2})$$

$$= \frac{\gamma_{1}}{2} \langle x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - \bar{x}_{t+1} - (x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - \bar{x}_{t}),$$

$$x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - \bar{x}_{t+1} + (x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - \bar{x}_{t}) \rangle$$

$$= \frac{\gamma_{1}}{2} \langle \bar{x}_{t+1} - \bar{x}_{t}, \bar{x}_{t+1} + \bar{x}_{t} - 2x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) \rangle,$$
(37)

where the second step holds due to  $x^*(\bar{y}_{t+1}; \bar{\hat{x}}_t, \bar{\hat{y}}_t) = \arg\min_{x \in \mathbb{R}^{d_1}} F(x, \bar{y}_{t+1}; \bar{\hat{x}}_t, \bar{\hat{y}}_t)$ , the fourth step holds due to the fact  $a^2 - b^2 = (a - b)(a + b)$ .

Moreover, because  $F_d(y; \hat{x}, \hat{y})$  is  $L_d$ -smooth, we have

$$F_{d}(\bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t}) \geq F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) + \langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{y}_{t+1} - \bar{y}_{t} \rangle - \frac{L_{d}}{2} \|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}$$

$$= F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) + \beta_{y} \eta \langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle - \frac{\beta_{y}^{2} \eta^{2} L_{d}}{2} \|\bar{v}_{t}\|^{2}. \tag{38}$$

By combining the above three inequalities, we have

$$\mathbb{E}[F_{d}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t+1})] - \mathbb{E}[F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})] \\
= \mathbb{E}[F_{d}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t+1}) - F_{d}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) \\
+ \left(F_{d}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - F_{d}(\bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t})\right) + \left(F_{d}(\bar{y}_{t+1}; \bar{x}_{t}, \bar{y}_{t}) - F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\right) \\
\geq \beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] - \frac{\beta_{y}^{2} \eta^{2} L_{d}}{2} \mathbb{E}[\|\bar{v}_{t}\|^{2}] + \frac{\gamma_{2}(2 - \hat{\beta}_{y} \eta)}{2\hat{\beta}_{y} \eta} \mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}] \\
+ \frac{\gamma_{1}}{2} \mathbb{E}[\langle \bar{x}_{t+1} - \bar{x}_{t}, \bar{x}_{t+1} + \bar{x}_{t} - 2x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) \rangle] . \tag{39}$$

**Lemma C.3.** Given Assumptions 3.1-3.4, the following inequality holds:

$$q(\bar{\hat{x}}_{t+1}) - q(\bar{\hat{x}}_t) \le \frac{\gamma_1}{2} \langle \bar{\hat{x}}_{t+1} - \bar{\hat{x}}_t, \bar{\hat{x}}_{t+1} + \bar{\hat{x}}_t - 2x^*(\bar{\hat{x}}_t, \hat{y}^*(\bar{\hat{x}}_{t+1})) \rangle . \tag{40}$$

Proof.

$$q(\bar{x}_{t+1}) - q(\bar{x}_t)$$

$$= g(\bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})) - g(\bar{x}_t, \hat{y}^*(\bar{x}_t))$$

$$= F_p(x^*(\bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})) - F_p(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_t)); \bar{x}_t, \hat{y}^*(\bar{x}_t))$$

$$\leq F_p(x^*(\bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})) - F_p(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_t, \hat{y}^*(\bar{x}_{t+1}))$$

$$\leq F_p(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})) - F_p(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_t, \hat{y}^*(\bar{x}_{t+1}))$$

$$= F(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1}))$$

$$= F(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})), y^*(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_t, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_t, \hat{y}^*(\bar{x}_{t+1}))$$

$$- F(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})), y^*(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_t, \hat{y}^*(\bar{x}_{t+1}))$$

$$- F(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})), y^*(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_t, \hat{y}^*(\bar{x}_{t+1}))$$

$$- F(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})), y^*(x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^*(\bar{x}_{t+1})); \bar{x}_t, \hat{y}^*(\bar{x}_{t+1}))$$

$$= \frac{\gamma_1}{2} (\|x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})) - \bar{x}_{t+1}\|^2 - \|x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})) - \bar{x}_t\|^2)$$

$$= \frac{\gamma_1}{2} \langle x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})) - \bar{x}_{t+1} + (x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})) - \bar{x}_t) \rangle$$

$$= \frac{\gamma_1}{2} \langle \bar{x}_{t+1} - \bar{x}_t, \bar{x}_{t+1} + \bar{x}_t - 2x^*(\bar{x}_t, \hat{y}^*(\bar{x}_{t+1})) \rangle,$$
(41)

where the second step holds due to  $g(\hat{x}, \hat{y}) = \min_{x \in \mathbb{R}^{d_1}} F_p(x; \hat{x}, \hat{y})$ , the three inequalities hold due to  $y^*(x; \hat{x}, \hat{y}) = \arg\max_y F(x, y; \hat{x}, \hat{y})$  and  $F_p(x; \hat{x}, \hat{y}) = F(x, y^*(x; \hat{x}, \hat{y}); \hat{x}, \hat{y})$ , the second to last step holds due to the fact  $a^2 - b^2 = (a - b)(a + b)$ .

**Lemma C.4.** Given Assumptions 3.1-3.4, the following inequality holds:

$$||x^*(\bar{\hat{x}}_{t+1}, \hat{y}^+(\bar{\hat{x}}_{t+1})) - x^*(\bar{\hat{x}}_{t+1}, \hat{y}^*(\bar{\hat{x}}_{t+1}))||^2$$

$$\leq \frac{2}{\gamma_1 - L} \frac{2\gamma_2^2 C_{y_{\hat{x}\hat{y}}^1}^2}{\mu} ||\bar{\hat{x}}_{t+1} - \bar{\hat{x}}_t||^2 + \frac{2}{\gamma_1 - L} \frac{2\gamma_2^2}{\mu} \left( C_{y_{\hat{x}\hat{y}}^2}^2 + \frac{(1 - \hat{\beta}_y \eta)^2}{\hat{\beta}_y^2 \eta^2} \right) ||\bar{\hat{y}}_t - \hat{y}^+(\bar{\hat{x}}_{t+1})||^2.$$
(42)

Proof.

Proof.

$$\frac{\gamma_{1} - L}{2} \|x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})) - x^{*}(\bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1}))\|^{2} \\
\leq F_{p}(x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1})) - F_{p}(x^{*}(\bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1})) \\
\leq \max_{\hat{y} \in \mathbb{R}^{d_{2}}} F_{p}(x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}) - F_{p}(x^{*}(\bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1})) \\
\leq \max_{\hat{y} \in \mathbb{R}^{d_{2}}} F_{p}(x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}) - F_{p}(x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})) \\
\leq \max_{\hat{y} \in \mathbb{R}^{d_{2}}} F_{p}(x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}) - F_{p}(x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})) \\
\leq \frac{1}{2\mu} \|\nabla_{\hat{y}} F_{p}(x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1}))\|^{2} \\
= \frac{\gamma_{2}^{2}}{2\mu} \|y^{*}(x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})) - \hat{y}^{+}(\bar{x}_{t+1})\|^{2} \\
= \frac{\gamma_{2}^{2}}{2\mu} \|y^{*}(x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})); \bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})) - \hat{y}^{*}(\bar{x}_{t+1})\|^{2} \\
\leq \frac{\gamma_{2}^{2}}{2\mu} \|y^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})) - y^{*}(\bar{x}_{t}, \hat{y}_{t})\|^{2} + (1 - \hat{\beta}_{y}\eta)^{2} \frac{\gamma_{2}^{2}}{\mu} \|\hat{y}_{t} - y^{*}(\bar{x}_{t}, \hat{y}_{t})\|^{2} \\
\leq \frac{2\gamma_{2}^{2}}{\mu} \|y^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})) - y^{*}(\bar{x}_{t}, \hat{y}^{+}(\bar{x}_{t+1}))\|^{2} \\
\leq \frac{2\gamma_{2}^{2}}{\mu} \|y^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})) - y^{*}(\bar{x}_{t}, \hat{y}_{t})\|^{2} \\
\leq \frac{2\gamma_{2}^{2}}{\mu} \|y^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1})) - y^{*}(\bar{x}_{t}, \hat{y}_{t})\|^{2} \\
\leq \frac{2\gamma_{2}^{2}}{\mu} \|y^{*}(\bar{x}_{t+1} - \bar{x}_{t})\|^{2} + \frac{2\gamma_{2}^{2}}{\mu} (C_{y_{x_{\hat{y}_{$$

where the first step holds because  $F_p(x; \hat{x}, \hat{y})$  is  $(\gamma_1 - L)$ -strongly convex with respect to x, the fourth step holds due to Theorem 5.2 of (Yu et al., 2022) with PL property being a special KL property, the fifth step holds due to the definition of  $F_p$ , the sixth step and the last step hold due to the definition  $\hat{y}^+(\bar{\hat{x}}_{t+1}) = \bar{\hat{y}}_t + \hat{\beta}_y \eta(y^*(\bar{\hat{x}}_t, \bar{\hat{y}}_t) - \bar{\hat{y}}_t).$ 

**Lemma C.5.** Given Assumptions 3.1-3.4, the following inequality holds: 

$$\mathbb{E}[\|x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - x^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1}))\|^{2}] \\
\leq 10\beta_{y}^{2}\eta^{2}C_{x_{y\hat{x}\hat{y}}}^{2}\mathbb{E}[\|\bar{v}_{t} - \nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] + 10\beta_{y}^{2}\eta^{2}L^{2}C_{x_{y\hat{x}\hat{y}}}^{2}\mathbb{E}[\|\bar{x}_{t} - x^{*}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
+ 5C_{x_{y\hat{x}\hat{y}}}^{2}\left(1 + \frac{1}{\beta_{y}^{2}\eta^{2}(\gamma_{2} - L)^{2}}\right)\mathbb{E}[\|y^{+}(\bar{x}_{t}, \bar{y}_{t}) - \bar{y}_{t}\|^{2}] \\
+ 5C_{x_{y\hat{x}\hat{y}}}^{2}C_{y_{\hat{x}\hat{y}}}^{2}\mathbb{E}[\|\bar{x}_{t} - \bar{x}_{t+1}\|^{2}] + 5C_{x_{y\hat{x}\hat{y}}}^{2}C_{y_{\hat{x}\hat{y}}}^{2}\mathbb{E}[\|\bar{y}_{t} - \hat{y}^{+}(\bar{x}_{t+1})\|^{2}]. \tag{44}$$

Proof.

$$\mathbb{E}[\|x^*(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_t) - x^*(\bar{x}_{t+1}, \hat{y}^+(\bar{x}_{t+1}))\|^2]$$

$$\leq C_{x_{1}, \hat{y}, \hat{y}}^{1} \mathbb{E}[\|\bar{y}_{t+1} - y^*(\bar{x}_{t+1}, \hat{y}^+(\bar{x}_{t+1}))\|^2]$$

$$\begin{aligned} & = C_{x_{y\hat{x}\hat{y}}}^{2} \mathbb{E}[\|\bar{y}_{t+1} - y^{+}(\bar{x}_{t}, \bar{y}_{t}) + y^{+}(\bar{x}_{t}, \bar{y}_{t}) - \bar{y}_{t} + \bar{y}_{t} - y^{*}(\bar{x}_{t}, \bar{y}_{t}) \\ & + y^{*}(\bar{x}_{t}, \bar{y}_{t}) - y^{*}(\bar{x}_{t+1}, \bar{y}_{t}) + y^{*}(\bar{x}_{t+1}, \bar{y}_{t}) - y^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1}))\|^{2} \\ & + y^{*}(\bar{x}_{t}, \bar{y}_{t}) - y^{*}(\bar{x}_{t+1}, \bar{y}_{t}) + y^{*}(\bar{x}_{t+1}, \bar{y}_{t}) - y^{*}(\bar{x}_{t+1}, \hat{y}^{+}(\bar{x}_{t+1}))\|^{2} \\ & \leq 5C_{x_{y\hat{x}\hat{y}}}^{2} \mathbb{E}[\|\bar{y}_{t+1} - y^{+}(\bar{x}_{t}, \bar{y}_{t})\|^{2}] + 5C_{x_{y\hat{x}\hat{y}}}^{2} \mathbb{E}[\|y^{+}(\bar{x}_{t}, \bar{y}_{t}) - \bar{y}_{t}\|^{2}] + 5C_{x_{y\hat{x}\hat{y}}}^{2} \mathbb{E}[\|\bar{y}_{t} - y^{*}(\bar{x}_{t}, \bar{y}_{t})\|^{2}] \\ & + 5C_{x_{y\hat{x}\hat{y}}}^{2} \mathbb{E}[\|y^{*}(\bar{x}_{t}, \bar{y}_{t}) - y^{*}(\bar{x}_{t+1}, \bar{y}_{t})\|^{2}] + 5C_{x_{y\hat{x}\hat{y}}}^{2} \mathbb{E}[\|y^{*}(\bar{x}_{t+1}, \bar{y}_{t}) - y^{*}(\bar{x}_{t+1}, \hat{y}_{t})\|^{2}] \\ & \leq 10\beta_{y}^{2} \gamma^{2} C_{x_{y\hat{x}\hat{y}}}^{2} \mathbb{E}[\|\bar{v}_{t} - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] + 10\beta_{y}^{2} \gamma^{2} L^{2} C_{x_{y\hat{x}\hat{y}}}^{2} \mathbb{E}[\|\bar{x}_{t} - x^{*}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\ & + 5C_{x_{y\hat{x}\hat{y}}}^{2} \left(1 + \frac{1}{\beta_{y}^{2} \gamma^{2} (\gamma_{2} - L)^{2}}\right) \mathbb{E}[\|y^{+}(\bar{x}_{t}, \bar{y}_{t}) - \bar{y}_{t}\|^{2}] \\ & + 5C_{x_{y\hat{x}\hat{y}}}^{2} C_{y_{\hat{x}\hat{y}}}^{2} \mathbb{E}[\|\bar{x}_{t} - \bar{x}_{t+1}\|^{2}] + 5C_{x_{y\hat{x}\hat{y}}}^{2} C_{y_{\hat{x}\hat{y}}}^{2} \mathbb{E}[\|\bar{y}_{t} - \hat{y}^{+}(\bar{x}_{t+1})\|^{2}], \end{aligned} \tag{45}$$

where the last step holds due to the following inequality:

$$\mathbb{E}[\|\bar{y}_{t+1} - y^{+}(\bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
= \mathbb{E}[\|\bar{y}_{t} + \beta_{y}\eta\bar{v}_{t} - \bar{y}_{t} - \beta_{y}\eta\nabla_{y}F(x^{*}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
= \beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t} - \nabla_{y}F(x^{*}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
\leq 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t} - \nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
+ 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y}F(x^{*}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] \\
\leq 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t} - \nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] + 2\beta_{y}^{2}\eta^{2}L^{2}\mathbb{E}[\|\bar{x}_{t} - x^{*}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] . \tag{46}$$

**Lemma C.6.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\mathbb{E}[\|\bar{\hat{y}}_{t} - \hat{y}^{+}(\bar{\hat{x}}_{t+1})\|^{2}] \leq 2\mathbb{E}[\|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_{t}\|^{2}] + 4\hat{\beta}_{y}^{2}\eta^{2}\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t}\|^{2}] + \frac{4\hat{\beta}_{y}^{2}}{\beta_{y}^{2}(\gamma_{2} - L)^{2}}\mathbb{E}[\|\bar{y}_{t} - y^{+}(\bar{\hat{x}}_{t}, \bar{\hat{y}}_{t})\|^{2}].$$

$$(47)$$

Proof.

$$\frac{1}{2}\mathbb{E}[\|\bar{\hat{y}}_{t} - \hat{y}^{+}(\bar{\hat{x}}_{t+1})\|^{2}] \\
\leq \mathbb{E}[\|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_{t}\|^{2}] + \mathbb{E}[\|\bar{\hat{y}}_{t+1} - \hat{y}^{+}(\bar{\hat{x}}_{t+1})\|^{2}] \\
\leq \mathbb{E}[\|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_{t}\|^{2}] + \mathbb{E}[\|\bar{\hat{y}}_{t} + \hat{\beta}_{y}\eta(\bar{y}_{t+1} - \bar{\hat{y}}_{t}) - \bar{\hat{y}}_{t} - \hat{\beta}_{y}\eta(y^{*}(\bar{\hat{x}}_{t}, \bar{\hat{y}}_{t}) - \bar{\hat{y}}_{t})\|^{2}] \\
= \mathbb{E}[\|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_{t}\|^{2}] + \hat{\beta}_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{y}_{t+1} - y^{*}(\bar{\hat{x}}_{t}, \bar{\hat{y}}_{t})\|^{2}] \\
\leq \mathbb{E}[\|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_{t}\|^{2}] + 2\hat{\beta}_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}] + 2\hat{\beta}_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{y}_{t} - y^{*}(\bar{\hat{x}}_{t}, \bar{\hat{y}}_{t})\|^{2}] \\
\leq \mathbb{E}[\|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_{t}\|^{2}] + 2\hat{\beta}_{y}^{2}\eta^{2}\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t}\|^{2}] + \frac{2\hat{\beta}_{y}^{2}}{\beta_{y}^{2}(\gamma_{2} - L)^{2}}\mathbb{E}[\|\bar{y}_{t} - y^{+}(\bar{\hat{x}}_{t}, \bar{\hat{y}}_{t})\|^{2}] . \tag{48}$$

**Lemma C.7.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\mathbb{E}[\|\bar{y}_{t} - y^{+}(\bar{x}_{t}, \bar{y}_{t})\|^{2}] \leq 4\beta_{y}^{2}\eta^{2}L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{x}_{t}\|^{2}] + 4\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{v}_{t}\|^{2}] + 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t}\|^{2}].$$

$$(49)$$

Proof.

```
1129
1130 \mathbb{E}[\|y^{+}(\bar{x}_{t}, \bar{y}_{t}) - \bar{y}_{t}\|^{2}]
1131 \leq 2\mathbb{E}[\|y^{+}(\bar{x}_{t}, \bar{y}_{t}) - \bar{y}_{t+1}\|^{2}] + 2\mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}]
1132 = 2\mathbb{E}[\|\bar{y}^{+}(\bar{x}_{t}, \bar{y}_{t}) - \bar{y}_{t+1}\|^{2}] + 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t}\|^{2}]
1133 = 2\mathbb{E}[\|\bar{y}_{t} + \beta_{y}\eta\nabla_{y}F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{y}_{t} - \beta_{y}\eta\bar{v}_{t}\|^{2}] + 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t}\|^{2}]
= 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\nabla_{y}F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{v}_{t}\|^{2}] + 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t}\|^{2}]
```

$$\leq 4\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\nabla_{y}F_{d}(\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t})\|^{2}]$$

$$+ 4\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{v}_{t}\|^{2}] + 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t}\|^{2}]$$

$$\leq 4\beta_{y}^{2}\eta^{2}L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{x}_{t}\|^{2}] + 4\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{v}_{t}\|^{2}] + 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t}\|^{2}] .$$

$$\leq 4\beta_{y}^{2}\eta^{2}L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{x}_{t}\|^{2}] + 4\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{v}_{t}\|^{2}] + 2\beta_{y}^{2}\eta^{2}\mathbb{E}[\|\bar{v}_{t}\|^{2}] .$$

$$\qquad \Box$$

**Lemma C.8.** Given Assumptions 3.1-3.4, by defining

$$\mathcal{P}_t = \mathbb{E}[F(\bar{x}_t, \bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t)] - 2\mathbb{E}[F_d(\bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t)] + 2\mathbb{E}[q(\bar{\hat{x}}_t)], \qquad (51)$$

by setting  $\eta \leq \frac{1}{\hat{\beta}_x}$ ,  $\eta \leq \frac{1}{\hat{\beta}_y}$ , and  $\beta_x \leq \min\{\frac{L^2}{120\gamma_1^3}, \frac{\sqrt{\mu(\gamma_1 - L)^3}(\gamma_2 - L)^2}{512\sqrt{6\gamma_1c_{\hat{\beta}_x}}\gamma_2c_{\hat{\beta}_y}}\}$ , then the following inequality holds:

$$\begin{array}{ll} & \mathcal{P}_{t+1} - \mathcal{P}_{t} \leq -\frac{\beta_{x}\eta}{4}\mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{\hat{y}}_{t})\|^{2}] - \frac{\beta_{y}\eta}{2}\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{\hat{y}}_{t})\|^{2}] \\ & + \frac{\beta_{x}\eta}{2}\mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{\hat{y}}_{t}) - \bar{u}_{t}\|^{2}] + A_{3}\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{\hat{y}}_{t}) - \bar{v}_{t}\|^{2}] \\ & + \left(4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2} - \frac{\beta_{x}\eta}{4}\right)\mathbb{E}[\|\bar{u}_{t}\|^{2}] \\ & + \left(\beta_{y}^{2}\eta^{2}L_{d} + \frac{3\beta_{y}\eta}{4} + \frac{\beta_{y}^{2}\eta^{2}(\gamma_{2} + L)}{2} + 4A_{1}\hat{\beta}_{y}^{2}\eta^{2}\beta_{y}^{2}\eta^{2} + 2A_{2}\beta_{y}^{2}\eta^{2} - \frac{7}{8}\beta_{y}\eta\right)\mathbb{E}[\|\bar{v}_{t}\|^{2}] \\ & + \left(2\gamma_{1}C_{x_{xx}} + \frac{\gamma_{1}}{6\hat{\beta}_{x}\eta} + 6\gamma_{1}\hat{\beta}_{x}\eta\left(10C_{x_{yx}}^{2}C_{y_{xx}}^{2}C_{y_{xx}}^{2} + \frac{4}{\gamma_{1} - L}\frac{2\gamma_{2}^{2}C_{y_{xx}}^{2}}{\mu}\right) - \frac{\gamma_{1}(2 - \hat{\beta}_{x}\eta)}{2\hat{\beta}_{x}\eta}\right)\mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] \\ & + \left(2A_{1} - \frac{\gamma_{2}(2 - \hat{\beta}_{y}\eta)}{2\hat{\beta}_{y}\eta}\right)\mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}], \end{array} \tag{52}$$

where

$$A_{1} = 6\gamma_{1}\hat{\beta}_{x}\eta \left(10C_{x_{y\hat{x}\hat{y}}^{1}\hat{x}\hat{y}}^{2}C_{y_{x\hat{y}}^{2}}^{2} + \frac{4}{\gamma_{1} - L}\frac{2\gamma_{2}^{2}}{\mu}\left(C_{y_{x\hat{y}}^{2}}^{2} + \frac{(1 - \hat{\beta}_{y}\eta)^{2}}{\hat{\beta}_{y}^{2}\eta^{2}}\right)\right),$$

$$A_{2} = 60\gamma_{1}\hat{\beta}_{x}\eta C_{x_{y\hat{x}\hat{y}}^{1}}^{2}\left(1 + \frac{1}{\beta_{y}^{2}\eta^{2}(\gamma_{2} - L)^{2}}\right) + A_{1}\frac{4\hat{\beta}_{y}^{2}}{\beta_{y}^{2}(\gamma_{2} - L)^{2}},$$

$$A_{3} = \beta_{y}\eta + 120\gamma_{1}\hat{\beta}_{x}\eta\beta_{y}^{2}\eta^{2}C_{x_{y\hat{x}\hat{y}}^{2}}^{2} + 4A_{2}\beta_{y}^{2}\eta^{2}.$$

$$(53)$$

and

$$\beta_{y} = \beta_{x} \underbrace{\frac{(\gamma_{1} - L)^{2}}{64L^{2}}}_{c_{\beta_{y}} = O(1)}, \quad \hat{\beta}_{x} = \beta_{x} \underbrace{\frac{(\gamma_{1} - L)^{4}(\gamma_{2} - L)^{2}\mu}{24 \times 64^{2}\gamma_{1}L^{2}\left(5\gamma_{1}^{2}\mu + 16\gamma_{2}^{2}(\gamma_{1} - L)\right)}_{c_{\hat{\beta}_{x}} = O(1/\kappa)},$$

$$\hat{\beta}_{y} = \beta_{x} \underbrace{\frac{(\gamma_{1} - L)^{4}(\gamma_{2} - L)^{4}}{64^{2} \times 480\gamma_{1}^{3}\gamma_{2}^{2}L^{2}}}_{c_{\hat{\beta}_{x}} = O(1)}.$$
(54)

*Proof.* Based on Lemmas C.1, C.2, C.3, we have

$$\begin{split} &\mathcal{P}_{t+1} - \mathcal{P}_{t} \leq -\frac{\beta_{x}\eta}{2}\mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t})\|^{2}] + \frac{\beta_{y}\eta}{2}\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t})\|^{2}] \\ &+ \frac{\beta_{x}\eta}{2}\mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t}) - \bar{u}_{t}\|^{2}] + \left(4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2} - \frac{\beta_{x}\eta}{4}\right)\mathbb{E}[\|\bar{u}_{t}\|^{2}] \\ &+ \left(\beta_{y}^{2}\eta^{2}L_{d} + \frac{3\beta_{y}\eta}{4} + \frac{\beta_{y}^{2}\eta^{2}(\gamma_{2} + L)}{2}\right)\mathbb{E}[\|\bar{v}_{t}\|^{2}] + \left(-\frac{\gamma_{1}(2 - \hat{\beta}_{x}\eta)}{2\hat{\beta}_{x}\eta}\right)\mathbb{E}[\|\bar{\hat{x}}_{t+1} - \bar{\hat{x}}_{t}\|^{2}] \end{split}$$

$$\begin{aligned} & + \left( -\frac{\gamma_{2}(2 - \hat{\beta}_{y} \eta)}{2 \hat{\beta}_{y} \eta} \right) \mathbb{E}[\|\hat{\bar{y}}_{t+1} - \hat{\bar{y}}_{t}\|^{2}] - 2\beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \\ & + 2\gamma_{1} \mathbb{E}[\langle \hat{\bar{x}}_{t+1} - \hat{\bar{x}}_{t}, x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - x^{*}(\bar{x}_{t}, \hat{y}^{*}(\bar{x}_{t+1})) \rangle] \,. \end{aligned}$$
 (55)
$$\begin{aligned} & + 2\gamma_{1} \mathbb{E}[\langle \hat{\bar{x}}_{t+1} - \bar{x}_{t}, x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - x^{*}(\bar{x}_{t}, \hat{y}^{*}(\bar{x}_{t+1})) \rangle] \,. \end{aligned}$$
 (55)
$$\begin{aligned} & + 2\gamma_{1} \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \\ & + 2\gamma_{1} \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] , \text{we have} \end{aligned}$$

$$- 2\beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \\ & + 2\beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \end{aligned}$$
 
$$= -2\beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \\ & = -2\beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \end{aligned}$$
 
$$= -2\beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \\ & = -2\beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \end{aligned}$$
 
$$= -2\beta_{y} \eta \mathbb{E}[\langle \nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \\ & -\beta_{y} \eta \mathbb{E}[\|\nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}), \bar{v}_{t} \rangle] \\ & + \beta_{y} \eta \mathbb{E}[\|\nabla_{y} F_{d}(\bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) \|^{2}] \\ & -\beta_{y} \eta \mathbb{E}[\|\nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) \|^{2}] - (1 - \nu)\beta_{y} \eta \mathbb{E}[\|\bar{v}_{t}\|^{2}] + \beta_{y} \eta \mathbb{E}[\|\nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{v}_{t}\|^{2}] \\ & -\beta_{y} \eta \mathbb{E}[\|\nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{x}_{t}\|^{2}] \\ & -\beta_{y} \eta \mathbb{E}[\|\nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t}) - \bar{x}_{t}\|^{2}] \\ & -\beta_{y} \eta \mathbb{E}[\|\nabla$$

where the third step holds due to Young's inequality  $2a^Tb \le \frac{1}{\nu}||a||^2 + \nu||b||^2$  with  $\nu > 0$  being a constant, and the last step holds due to the following inequality:

$$\mathbb{E}[\|\nabla_{y}F(x^{*}(\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t}),\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t}) - \nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t})\|^{2}]$$

$$= \mathbb{E}[\|\nabla_{y}f(x^{*}(\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t}),\bar{y}_{t}) - \nabla_{y}f(\bar{x}_{t},\bar{y}_{t})\|^{2}] \leq L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t}) - \bar{x}_{t}\|^{2}].$$
(57)

For 
$$2\gamma_{1}\mathbb{E}[\langle \bar{x}_{t+1} - \bar{x}_{t}, x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - x^{*}(\bar{x}_{t}, \hat{y}^{*}(\bar{x}_{t+1}))\rangle]$$
, we have 
$$2\gamma_{1}\mathbb{E}[\langle \bar{x}_{t+1} - \bar{x}_{t}, x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - x^{*}(\bar{x}_{t}, \hat{y}^{*}(\bar{x}_{t+1}))\rangle]$$

$$= 2\gamma_{1}\mathbb{E}[\langle \bar{x}_{t+1} - \bar{x}_{t}, x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - x^{*}(\bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1}))\rangle]$$

$$+ 2\gamma_{1}\mathbb{E}[\langle \bar{x}_{t+1} - \bar{x}_{t}, x^{*}(\bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1})) - x^{*}(\bar{x}_{t}, \hat{y}^{*}(\bar{x}_{t+1}))\rangle]$$

$$\leq \frac{\gamma_{1}}{6\hat{\beta}_{x}\eta}\mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] + 6\gamma_{1}\hat{\beta}_{x}\eta\mathbb{E}[\|x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - x^{*}(\bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1}))\|^{2}]$$

$$+ 2\gamma_{1}\mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|\|x^{*}(\bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1})) - x^{*}(\bar{x}_{t}, \hat{y}^{*}(\bar{x}_{t+1}))\|]$$

$$\leq \frac{\gamma_{1}}{6\hat{\beta}_{x}\eta}\mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] + 6\gamma_{1}\hat{\beta}_{x}\eta\mathbb{E}[\|x^{*}(\bar{y}_{t+1}; \bar{x}_{t+1}, \bar{y}_{t}) - x^{*}(\bar{x}_{t+1}, \hat{y}^{*}(\bar{x}_{t+1}))\|^{2}]$$

$$+ 2\gamma_{1}C_{x_{x\bar{x}\bar{y}}^{1}}\mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}], \qquad (58)$$

where the second step holds due to Young's inequality  $2a^Tb \leq \frac{1}{\nu}\|a\|^2 + \nu\|b\|^2$  with  $\nu = 6\hat{\beta}_x\eta$  and  $a^Tb \leq \|a\|\|b\|$ , and the last step holds due to Lemma B.2.

Then, by plugging Eq. (56) and Eq. (58) into Eq. (55) with  $\nu = \frac{1}{8}$ , we have

$$\mathcal{P}_{t+1} - \mathcal{P}_{t} \leq -\frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t})\|^{2}] - \frac{\beta_{y}\eta}{2} \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t})\|^{2}]$$

$$+ \frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t}) - \bar{u}_{t}\|^{2}] + \beta_{y}\eta \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t}) - \bar{v}_{t}\|^{2}]$$

$$+ \left(4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2} - \frac{\beta_{x}\eta}{4}\right) \mathbb{E}[\|\bar{u}_{t}\|^{2}] + \left(\beta_{y}^{2}\eta^{2}L_{d} + \frac{3\beta_{y}\eta}{4} + \frac{\beta_{y}^{2}\eta^{2}(\gamma_{2} + L)}{2} - \frac{7}{8}\beta_{y}\eta\right) \mathbb{E}[\|\bar{v}_{t}\|^{2}]$$

$$\begin{aligned} & + \left( 2\gamma_{1}C_{x_{xy}^{1}} + \frac{\gamma_{1}}{6\hat{\beta}_{x}\eta} - \frac{\gamma_{1}(2-\hat{\beta}_{x}\eta)}{2\hat{\beta}_{x}\eta} \right) \mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] + \left( -\frac{\gamma_{2}(2-\hat{\beta}_{y}\eta)}{2\hat{\beta}_{y}\eta} \right) \mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}] \\ & + 8\beta_{y}\eta L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{x}_{t}\|^{2}] + 6\gamma_{1}\hat{\beta}_{x}\eta\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}_{t}) - x^{*}(\bar{x}_{t+1}))\|^{2}] \\ & + 8\beta_{y}\eta L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{x}_{t}\|^{2}] + 6\gamma_{1}\hat{\beta}_{x}\eta\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}^{*}) - x^{*}(\bar{x}_{t+1}))\|^{2}] \\ & + 8\beta_{y}\eta L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}_{t}) - \bar{x}_{t}\|^{2}] + 6\gamma_{1}\hat{\beta}_{x}\eta\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}^{*}) - x^{*}(\bar{x}_{t+1})]^{2}] \\ & + 2246 \end{aligned}$$

$$\begin{aligned} & + 8\beta_{y}\eta L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}_{t}) - \bar{x}_{t}\|^{2}] + 6\gamma_{1}\hat{\beta}_{x}\eta\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}^{*}(\bar{x}_{t+1}))\|^{2}] \\ & + 2126 \end{aligned}$$

$$\begin{aligned} & + 8\beta_{y}\eta L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}_{t}) - \bar{x}_{t}\|^{2}] + 6\gamma_{1}\hat{\beta}_{x}\eta\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}^{*}(\bar{x}_{t+1}))\|^{2}] \\ & + 2126 \end{aligned}$$

$$\begin{aligned} & + 8\beta_{y}\eta L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}_{t}) - \bar{x}_{t}\|^{2}] + 6\gamma_{1}\hat{\beta}_{x}\eta\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}^{*}(\bar{x}_{t+1}))\|^{2}] \\ & + 8\beta_{y}\eta L^{2}\mathbb{E}[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}_{t}) - x^{*}(\bar{x}_{t+1})]^{2}] \\ & + 2126 \end{aligned}$$

$$& + 2[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}_{t}) - x^{*}(\bar{x}_{t+1},\bar{y}^{*}(\bar{x}_{t+1}))\|^{2}] \\ & + 22[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}_{t}) - x^{*}(\bar{x}_{t+1},\bar{y}^{*}(\bar{x}_{t+1}))\|^{2}] \\ & + 22[\|x^{*}(\bar{y}_{t+1};\bar{x}_{t+1},\bar{y}_{t}) - x^{*}(\bar{x}_{t+1},\bar{y}^{*}(\bar{x}_{t+1}))\|^{2}] \\ & + 22[\|x^{*}(\bar{x}_{t+1},\bar{y}^{*}(\bar{x}_{t+1})]^{2}] \\ & + 22\beta_{y}^{2}\eta^{2}C^{2}_{y_{x}^{2}}\mathbb{E}[\|\bar{x}_{t} - \bar{x}_{t+1}\|^{2}] \\ & + 22\beta_{y}^{2}\eta^{2}C^{2}_{y_{x}^{2}}\mathbb{E}[\|\bar{x}_{t} - \bar{x}_{t+1}\|^{2}] \\ & + 22\beta_{y}^{2}\eta^{2}C^{2}_{y_{x}^{2}}\mathbb{E}[\|\bar{x}_{t} - \bar{x}_{t+1}\|^{2}] \\ & + 22\beta_{y}^{2}\eta^{2}C^{2}_{y_{x}^{2}}\mathbb{E}[\|\bar{x}_{t} - \bar{x}_{t},\bar{y}_{t}\|^{2}] \\ & + 22\beta_{y}^{2}\eta^{2}C^{2}_{y_{x}^{2}}\mathbb{E}[\|\bar{x}_{t} - \bar{x}_{t}\|^{2}] \\ & + 22\beta_{y}$$

where the second step holds due to Lemma C.4 and Lemma C.5.

By plugging the above inequality into Eq. (59), we have

$$\mathcal{P}_{t+1} - \mathcal{P}_{t} \leq -\frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{\hat{y}}_{t})\|^{2}] - \frac{\beta_{y}\eta}{2} \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{\hat{y}}_{t})\|^{2}]$$

$$+ \frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{\hat{y}}_{t}) - \bar{u}_{t}\|^{2}] + \left(\beta_{y}\eta + 120\gamma_{1}\hat{\beta}_{x}\eta\beta_{y}^{2}\eta^{2}C_{x_{yxy}^{2}}^{2}\right) \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{\hat{y}}_{t}) - \bar{v}_{t}\|^{2}]$$

$$+ \left(4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2} - \frac{\beta_{x}\eta}{4}\right) \mathbb{E}[\|\bar{u}_{t}\|^{2}] + \left(\beta_{y}^{2}\eta^{2}L_{d} + \frac{3\beta_{y}\eta}{4} + \frac{\beta_{y}^{2}\eta^{2}(\gamma_{2} + L)}{2} - \frac{7}{8}\beta_{y}\eta\right) \mathbb{E}[\|\bar{v}_{t}\|^{2}]$$

$$+ \left(2\gamma_{1}C_{x_{xy}^{1}} + \frac{\gamma_{1}}{6\hat{\beta}_{x}\eta} + 6\gamma_{1}\hat{\beta}_{x}\eta\left(10C_{x_{yx}^{2}}^{2}C_{y_{xy}^{2}}^{2} + \frac{4}{\gamma_{1} - L}\frac{2\gamma_{2}^{2}C_{y_{xy}^{2}}^{2}}{\mu}\right) - \frac{\gamma_{1}(2 - \hat{\beta}_{x}\eta)}{2\hat{\beta}_{x}\eta}\right) \mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}]$$

$$+ \left(-\frac{\gamma_{2}(2 - \hat{\beta}_{y}\eta)}{2\hat{\beta}_{y}\eta}\right) \mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}] + \left(8\beta_{y}\eta L^{2} + 120\gamma_{1}\hat{\beta}_{x}\eta\beta_{y}^{2}\eta^{2}L^{2}C_{x_{yx}^{2}}^{2}\right) \mathbb{E}[\|x^{*}(\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{x}_{t}\|^{2}]$$

$$+ 60\gamma_{1}\hat{\beta}_{x}\eta C_{x_{yx}^{2}}^{2}\left(1 + \frac{1}{\beta_{y}^{2}\eta^{2}(\gamma_{2} - L)^{2}}\right) \mathbb{E}[\|y^{+}(\bar{x}_{t},\bar{y}_{t}) - \bar{y}_{t}\|^{2}]$$

$$+ 6\gamma_{1}\hat{\beta}_{x}\eta\left(10C_{x_{yx}^{2}}^{2}C_{y_{xy}^{2}}^{2} + \frac{4}{\gamma_{1} - L}\frac{2\gamma_{2}^{2}}{\mu}\left(C_{y_{x}^{2}}^{2} + \frac{(1 - \hat{\beta}_{y}\eta)^{2}}{\hat{\beta}_{y}^{2}\eta^{2}}\right)\right) \mathbb{E}[\|\bar{y}_{t} - \hat{y}^{+}(\bar{x}_{t+1})\|^{2}] .$$

$$+ 6\gamma_{1}\hat{\beta}_{x}\eta\left(10C_{x_{yx}^{2}}^{2}C_{y_{xx}^{2}}^{2} + \frac{4}{\gamma_{1} - L}\frac{2\gamma_{2}^{2}}{\mu}\left(C_{y_{x}^{2}}^{2} + \frac{(1 - \hat{\beta}_{y}\eta)^{2}}{\hat{\beta}_{y}^{2}\eta^{2}}\right)\right) \mathbb{E}[\|\bar{y}_{t} - \hat{y}^{+}(\bar{x}_{t+1})\|^{2}] .$$

$$+ 6\gamma_{1}\hat{\beta}_{x}\eta\left(10C_{x_{yx}^{2}}^{2}C_{y_{xx}^{2}}^{2} + \frac{4}{\gamma_{1} - L}\frac{2\gamma_{2}^{2}}{\mu}\left(C_{y_{x}^{2}}^{2} + \frac{(1 - \hat{\beta}_{y}\eta)^{2}}{\hat{\beta}_{y}^{2}\eta^{2}}\right)\right) \mathbb{E}[\|\bar{y}_{t} - \hat{y}^{+}(\bar{x}_{t+1})\|^{2}] .$$

Furthermore, based on Lemma C.6, we have

$$\mathcal{P}_{t+1} - \mathcal{P}_{t} \leq -\frac{\beta_{x}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t}, \bar{y}_{t}; \hat{\bar{x}}_{t}, \bar{\bar{y}}_{t})\|^{2}] - \frac{\beta_{y}\eta}{2} \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{\bar{x}}_{t}, \bar{\bar{y}}_{t})\|^{2}]$$

$$\frac{\beta_{s}\eta}{2} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{u}_{t}\|^{2}] + \left(\beta_{y}\eta + 120\gamma_{1}\beta_{x}\eta\beta_{y}^{2}\eta^{2}C_{x_{yyy}}^{2}\right) \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{v}_{t}\|^{2}]$$

$$+ \left(4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2} - \frac{\beta_{x}\eta}{4}\right) \mathbb{E}[\|\bar{u}_{t}\|^{2}]$$

$$+ \left(\beta_{y}^{2}\eta^{2}L_{d} + \frac{3\beta_{y}\eta}{6\beta_{x}\eta} + 6\gamma_{1}\beta_{x}\eta\left(10C_{x_{yy}}^{2}C_{y_{yy}}^{2} + L\right) + 4A_{1}\beta_{y}^{2}\eta^{2}\beta_{y}^{2}\eta^{2} - \frac{7}{8}\beta_{y}\eta\right) \mathbb{E}[\|\bar{v}_{t}\|^{2}]$$

$$+ \left(2\gamma_{1}C_{x_{yy}}^{1} + \frac{\gamma_{1}}{6\beta_{x}\eta} + 6\gamma_{1}\beta_{x}\eta\left(10C_{x_{yy}}^{2}C_{y_{yy}}^{2} + \frac{4}{\gamma_{1} - L}\frac{2\gamma_{2}^{2}C_{y_{yy}}^{2}}{\mu}\right) - \frac{\gamma_{1}(2 - \beta_{x}\eta)}{2\beta_{x}\eta}\right) \mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}]$$

$$+ \left(2A_{1} - \frac{\gamma_{2}(2 - \beta_{y}\eta)}{2\beta_{y}\eta}\right) \mathbb{E}[\|\bar{u}_{t+1} - \bar{u}_{t}\|^{2}]$$

$$+ \left(8\beta_{y}\eta L^{2} + 120\gamma_{1}\beta_{x}\eta\beta_{y}^{2}\eta^{2}L^{2}C_{x_{yyy}}^{2}} \mathbb{E}[\|x^{*}(y_{t};\bar{x}_{t},\bar{y}_{t}) - x_{t}\|^{2}]$$

$$+ \left(60\gamma_{1}\beta_{x}\etaC_{x_{yy}}^{2}\right) \mathbb{E}[\|\bar{u}_{t+1} - \bar{u}_{t}\|^{2}]$$

$$+ \left(60\gamma_{1}\beta_{x}\etaC_{x_{yy}}^{2}\right) \mathbb{E}[\|\bar{u}_{t+1} - \bar{u}_{t}\|^{2}]$$

$$+ \left(60\gamma_{1}\beta_{x}\etaC_{x_{yy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}^{*}(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{y}_{t}\|^{2}] , \qquad (62)$$

$$+ \left(60\gamma_{1}\beta_{x}\etaC_{x_{yy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}^{*}(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{y}_{t}\|^{2}] , \qquad (62)$$

$$+ \left(60\gamma_{1}\beta_{x}\etaC_{x_{yy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}^{*}(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{y}_{t}\|^{2}]$$

$$+ \left(60\gamma_{1}\beta_{x}\etaC_{x_{yy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right)$$

$$+ \left(60\gamma_{1}\beta_{x}\eta\beta_{y}^{2}\eta^{2}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right)$$

$$+ \left(62\gamma_{1}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{v}_{t}\|_{2} + \frac{1}{\eta^{2}}C_{x_{xy}}^{2}\right) \mathbb{E}[\|\bar{$$

$$+ \left(2\gamma_{1}C_{x_{\hat{x}\hat{y}}^{1}} + \frac{\gamma_{1}}{6\hat{\beta}_{x}\eta} + 6\gamma_{1}\hat{\beta}_{x}\eta\left(10C_{x_{y\hat{x}\hat{y}}^{1}}^{2}C_{y_{\hat{x}\hat{y}}^{1}}^{1} + \frac{4}{\gamma_{1} - L}\frac{2\gamma_{2}^{2}C_{y_{\hat{x}\hat{y}}^{1}}^{2}}{\mu}\right) - \frac{\gamma_{1}(2 - \hat{\beta}_{x}\eta)}{2\hat{\beta}_{x}\eta}\right)\mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] 
+ \left(2A_{1} - \frac{\gamma_{2}(2 - \hat{\beta}_{y}\eta)}{2\hat{\beta}_{y}\eta}\right)\mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}],$$
(64)

where  $A_3 = \beta_y \eta + 120 \gamma_1 \hat{\beta}_x \eta \beta_y^2 \eta^2 C_{x_{y\hat{x}\hat{y}}^2}^2 + 4 A_2 \beta_y^2 \eta^2$ .

Then, for  $\mathbb{E}[\|\nabla_x F(\bar{x}_t, \bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t)\|^2]$ , we set

$$\frac{(A_3 + 7\beta_y \eta)L^2}{(\gamma_1 - L)^2} - \frac{\beta_x \eta}{2}$$

$$= \frac{1}{(\gamma_1 - L)^2} \left( 8\beta_y \eta L^2 + 120\gamma_1 \hat{\beta}_x \eta \beta_y^2 \eta^2 L^2 C_{x_{y\hat{x}\hat{y}}}^2 + 4A_2 \beta_y^2 \eta^2 L^2 \right) - \frac{\beta_x \eta}{2} \le -\frac{\beta_x \eta}{4} .$$
(65)

Specifically, we enforce

$$\frac{8\beta_{y}\eta L^{2}}{(\gamma_{1} - L)^{2}} \leq \frac{\beta_{x}\eta}{8} ,$$

$$\frac{120\gamma_{1}\hat{\beta}_{x}\eta\beta_{y}^{2}\eta^{2}L^{2}C_{x_{y\hat{x}\hat{y}}}^{2}}{(\gamma_{1} - L)^{2}} \leq \frac{\beta_{x}\eta}{32 \times 16} ,$$

$$\frac{4\beta_{y}^{2}\eta^{2}L^{2}}{(\gamma_{1} - L)^{2}}A_{2} \leq \frac{\beta_{x}\eta}{32 \times 16} .$$
(66)

For the first inequality in Eq. (66), we set

$$\beta_y = \beta_x \underbrace{\frac{(\gamma_1 - L)^2}{64L^2}}_{c_{\beta_y} = O(1)}.$$
(67)

For the last inequality in Eq. (66), from the definition of  $A_1$  and  $A_2$ , we enforce

$$\frac{4\beta_{y}^{2}\eta^{2}L^{2}}{(\gamma_{1}-L)^{2}}60\gamma_{1}\hat{\beta}_{x}\eta C_{x_{y\hat{x}\hat{y}}}^{2} \leq \frac{\beta_{x}\eta}{32\times64},$$

$$\frac{4L^{2}}{(\gamma_{1}-L)^{2}}\left(\frac{60\gamma_{1}\hat{\beta}_{x}\eta C_{x_{y\hat{x}\hat{y}}}^{2}}{(\gamma_{2}-L)^{2}} + \frac{4\hat{\beta}_{y}^{2}\eta^{2}}{(\gamma_{2}-L)^{2}} \frac{24\gamma_{1}\hat{\beta}_{x}\eta}{\gamma_{1}-L} \frac{2\gamma_{2}^{2}}{\mu} \frac{(1-\hat{\beta}_{y}\eta)^{2}}{\hat{\beta}_{y}^{2}\eta^{2}}\right) \leq \frac{\beta_{x}\eta}{32\times64},$$

$$\frac{4L^{2}}{(\gamma_{1}-L)^{2}}60\gamma_{1}\hat{\beta}_{x}\eta C_{x_{y\hat{x}\hat{y}}}^{2}C_{y_{x\hat{y}}}^{2}\frac{4\hat{\beta}_{y}^{2}\eta^{2}}{(\gamma_{2}-L)^{2}} \leq \frac{\beta_{x}\eta}{32\times64},$$

$$\frac{4L^{2}}{(\gamma_{1}-L)^{2}}6\gamma_{1}\hat{\beta}_{x}\eta \frac{4}{\gamma_{1}-L} \frac{2\gamma_{2}^{2}}{\mu}C_{y_{x\hat{x}\hat{y}}}^{2}\frac{4\hat{\beta}_{y}^{2}\eta^{2}}{(\gamma_{2}-L)^{2}} \leq \frac{\beta_{x}\eta}{32\times64}.$$
(68)

To solve the first inequality in Eq. (68), since  $\hat{\beta}_x \eta \leq 1$  and  $\eta < 1$ , from  $C_{x_{\eta \hat{x} \hat{y}}^1} = \frac{\gamma_1}{\gamma_1 - L}$ , we obtain

$$\beta_x \le \frac{L^2}{120\gamma_1^3} \ . \tag{69}$$

Here, we have also shown that the second inequality in Eq. (66) holds.

Then, to address the second inequality in Eq. (68), note that since  $\hat{\beta}_y \eta \leq 1$ , it follows that  $1 - \hat{\beta}_y \eta \leq 1$ . Consequently, we obtain

$$\hat{\beta}_x = \beta_x \frac{(\gamma_1 - L)^4 (\gamma_2 - L)^2 \mu}{24 \times 64^2 \gamma_1 L^2 \left( 5\gamma_1^2 \mu + 16\gamma_2^2 (\gamma_1 - L) \right)}$$

$$c_{\beta_x} = O(1/\kappa)$$
(70)

Similarly, for the third inequality in Eq. (68), from  $C_{y^2_{\hat{x}\hat{y}}}=rac{\gamma_2}{\gamma_2-L}$ , we obtain

$$\hat{\beta}_y = \beta_x \underbrace{\frac{(\gamma_1 - L)^4 (\gamma_2 - L)^4}{64^2 \times 480 \gamma_1^3 \gamma_2^2 L^2}}_{c_{\hat{\beta}_y} = O(1)}.$$
(71)

Moreover, to solve the last inequality in Eq. (68), we obtain

$$\beta_x \le \frac{\sqrt{\mu(\gamma_1 - L)^3(\gamma_2 - L)^2}}{512\sqrt{6\gamma_1 c_{\hat{\beta}_x}}\gamma_2 c_{\hat{\beta}_y}} = O(1) . \tag{72}$$

Finally, by plugging Eq. (65) into Eq. (64), the proof is complete.

## D KEY LEMMAS RELATED TO THE DECENTRALIZED SETTING

## D.1 Consensus Errors

**Lemma D.1.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{p}_{t+1} - p_{t+1}^{(k)}\|^{2}]$$

$$\leq \lambda \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{p}_{t} - p_{t}^{(k)}\|^{2}] + 3\rho_{x}^{2} \eta^{4} \frac{1}{1 - \lambda} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|u_{t}^{(k)} - \nabla_{x} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]$$

$$+ \frac{9(L^{2} + \gamma_{1}^{2})}{1 - \lambda} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|x_{t+1}^{(k)} - x_{t}^{(k)}\|^{2}] + \frac{9L^{2}}{1 - \lambda} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|y_{t+1}^{(k)} - y_{t}^{(k)}\|^{2}]$$

$$+ \frac{9\gamma_{1}^{2}}{1 - \lambda} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\hat{x}_{t+1}^{(k)} - \hat{x}_{t}^{(k)}\|^{2}] + 3\rho_{x}^{2} \eta^{4} \sigma^{2} \frac{1}{1 - \lambda}.$$
(73)

Proof.

$$\frac{1}{K} \sum_{t=1}^{K} \mathbb{E}[\|\bar{p}_{t+1} - p_{t+1}^{(k)}\|^2]$$

1438 
$$K \underset{k=1}{\overset{\sum}} \mathbb{E}[\|P_{t+1} - P_{t+1}\|]$$
1439 
$$= \frac{1}{K} \mathbb{E}[\|\bar{P}_{t+1} - P_{t+1}\|_F^2]$$
1441

1441  
1442 
$$= \frac{1}{K} \mathbb{E}[\|\bar{P}_t - \bar{U}_t + \bar{U}_{t+1} - P_t W + U_t - U_{t+1}\|_F^2]$$
1443

$$\leq (1+a)\frac{1}{K}\mathbb{E}[\|\bar{P}_t - P_t W\|_F^2] + (1+1/a)\frac{1}{K}\mathbb{E}[\|-\bar{U}_t + \bar{U}_{t+1} + U_t - U_{t+1}\|_F^2]$$

$$\leq (1+a)\lambda^2 \frac{1}{K} \mathbb{E}[\|\bar{P}_t - P_t\|_F^2] + (1+1/a) \frac{1}{K} \mathbb{E}[\|U_t - U_{t+1}\|_F^2]$$

$$\leq \lambda \frac{1}{K} \mathbb{E}[\|\bar{P}_t - P_t\|_F^2] + \frac{1}{1 - \lambda} \frac{1}{K} \mathbb{E}[\|U_t - U_{t+1}\|_F^2], \qquad (74)$$

where  $a = \frac{1-\lambda}{\lambda}$ . Then, we have the following inequality to complete the proof:

1451  
1452 
$$\frac{1}{K}\mathbb{E}[\|U_t - U_{t+1}\|_F^2]$$

$$= \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|u_{t+1}^{(k)} - u_{t}^{(k)}\|^{2}]$$

1456
$$= \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|(1 - \rho_x \eta^2)(u_t^{(k)} - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \xi_{t+1}^{(k)}))$$

$$\begin{aligned} & + \nabla_x F^{(k)}(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \hat{x}_{t+1}^{(k)}, \hat{y}_{t+1}^{(k)}; \xi_{t+1}^{(k)}) - u_t^{(k)} \|^2] \\ & \leq 3 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| - \rho_x \eta^2 u_t^{(k)} + \rho_x \eta^2 \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) \|^2] \\ & + 3 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| - \rho_x \eta^2 \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) + \rho_x \eta^2 \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) \|^2] \\ & + 3 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \xi_{t+1}^{(k)}) + \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \xi_{t+1}^{(k)}) + \nabla_x F^{(k)}(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) \|^2] \\ & + 3 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) \|^2] + 3 \rho_x^2 \eta^4 \sigma^2 \\ & + 3 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \hat{x}_t^{(k)}) \|^2] + 3 \rho_x^2 \eta^4 \sigma^2 \\ & + 3 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) \|^2] + 3 \rho_x^2 \eta^4 \sigma^2 \\ & + 3 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| u_t^{(k)} - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) \|^2] + 3 \rho_x^2 \eta^4 \sigma^2 \\ & + 9 (L^2 + \gamma_1^2) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| x_t^{(k)} - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) \|^2] + 3 \rho_x^2 \eta^4 \sigma^2 \\ & + 9 (L^2 + \gamma_1^2) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| x_t^{(k)} - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) \|^2] + 3 \rho_x^2 \eta^4 \sigma^2 \\ & + 9 (L^2 + \gamma_1^2) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| x_t^{(k)} - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) \|^2] \\ & + 9 \gamma_1^2 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| x_t^{(k)} - \hat{x}_t^{(k)} \|^2], \\ & + 9 \gamma_1^2 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| x_t^{(k)} - \hat{x}_t^{(k)} \|^2], \\ & + 9 \gamma_1^2 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| x_t^{(k)} - \hat{x}_t^{(k)} \|^2], \\ & + 9 \gamma_1^2 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| x_t^{(k)} - \hat{x}_t^{(k)} \|^2], \\ & + 9 \gamma_1^2 \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\| x_t^{(k)} - \hat{x}_t^{(k)} \|^2], \\ & + 9 \gamma_1^2 \frac{1}{K} \sum_$$

where the last step holds due to the following inequality:

$$\mathbb{E}[\|-\nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)};\xi_{t+1}^{(k)})+\nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)},\hat{y}_{t+1}^{(k)};\xi_{t+1}^{(k)})\|^{2}] \\
&=\mathbb{E}[\|-\nabla_{x}f^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\xi_{t+1}^{(k)})-\gamma_{1}(x_{t}^{(k)}-\hat{x}_{t}^{(k)}) \\
&+\nabla_{x}f^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\xi_{t+1}^{(k)})+\gamma_{1}(x_{t+1}^{(k)}-\hat{x}_{t+1}^{(k)})\|^{2}] \\
&\leq 3\mathbb{E}[\|\nabla_{x}f^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\xi_{t+1}^{(k)})-\nabla_{x}f^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\xi_{t+1}^{(k)})\|^{2}] \\
&+3\gamma_{1}^{2}\mathbb{E}[\|x_{t+1}^{(k)}-x_{t}^{(k)}\|^{2}]+3\gamma_{1}^{2}\mathbb{E}[\|\hat{x}_{t+1}^{(k)}-\hat{x}_{t}^{(k)}\|^{2}] \\
&\leq 3(L^{2}+\gamma_{1}^{2})\mathbb{E}[\|x_{t+1}^{(k)}-x_{t}^{(k)}\|^{2}]+3L^{2}\mathbb{E}[\|y_{t+1}^{(k)}-y_{t}^{(k)}\|^{2}]+3\gamma_{1}^{2}\mathbb{E}[\|\hat{x}_{t+1}^{(k)}-\hat{x}_{t}^{(k)}\|^{2}]. \tag{76}$$

**Lemma D.2.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{q}_{t+1} - q_{t+1}^{(k)}\|^{2}]$$

$$\leq \lambda \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{q}_{t} - q_{t}^{(k)}\|^{2}] + 3\rho_{y}^{2} \eta^{4} \frac{1}{1 - \lambda} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]$$

$$+ \frac{9L^{2}}{1 - \lambda} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|x_{t+1}^{(k)} - x_{t}^{(k)}\|^{2}] + \frac{9(L^{2} + \gamma_{2}^{2})}{1 - \lambda} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|y_{t+1}^{(k)} - y_{t}^{(k)}\|^{2}]$$

$$+ \frac{9\gamma_{2}^{2}}{1 - \lambda} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\hat{y}_{t+1}^{(k)} - \hat{y}_{t}^{(k)}\|^{2}] + 3\rho_{y}^{2} \eta^{4} \sigma^{2} \frac{1}{1 - \lambda}.$$
(77)

This lemma can be proved by following Lemma D.1. Thus, we omit its proof.

**Lemma D.3.** Given Assumptions 3.1-3.4, when  $\hat{\beta}_x \leq \frac{1-\lambda}{4}$ , the following inequality holds:

$$\mathbb{E}[\|\hat{X}_{t+1} - \bar{\hat{X}}_{t+1}\|_F^2] \leq \left(1 - \frac{\eta(1 - \lambda^2)}{4}\right) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{\hat{x}}_t - \hat{x}_t^{(k)}\|^2] \\
+ \frac{4\eta \hat{\beta}_x^2}{1 - \lambda^2} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{x}_t - x_t^{(k)}\|^2] + \frac{4\eta \hat{\beta}_x^2}{1 - \lambda^2} \frac{2\eta \beta_x^2}{1 - \lambda^2} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{p}_t - p_t^{(k)}\|^2]. \tag{78}$$

Proof.

1521 
$$\|\hat{X}_{t+1} - \bar{X}_{t+1}\|_F^2$$
1522 
$$\|\hat{X}_{t+1} - \bar{X}_{t+1}\|_F^2$$
1523 
$$= \|\hat{X}_t + \eta(\hat{X}_{t+1} - \hat{X}_t) - \bar{X}_t - \eta\hat{\beta}_x(\bar{X}_{t+1} - \bar{X}_t)\|_F^2$$
1525 
$$= \|\hat{X}_t + \eta(\hat{X}_tW + \hat{\beta}_x(X_{t+1} - \hat{X}_t) - \hat{X}_t - \eta\hat{\beta}_x(\bar{X}_{t+1} - \bar{X}_t)\|_F^2$$
1526 
$$= \|(1 - \eta)(\hat{X}_t - \bar{X}_t) + \eta(\hat{X}_tW - \bar{X}_t) + \eta\hat{\beta}_x(X_{t+1} - \hat{X}_t) - \eta\hat{\beta}_x(\bar{X}_{t+1} - \bar{X}_t)\|_F^2$$
1528 
$$\leq (1 - \eta)\|\hat{X}_t - \hat{X}_t\|_F^2 + \eta\|\hat{X}_tW - \hat{X}_t + \hat{\beta}_x(X_{t+1} - \hat{X}_t) - \hat{\beta}_x(\bar{X}_{t+1} - \bar{X}_t)\|_F^2$$
1529 
$$\leq (1 - \eta)\|\hat{X}_t - \hat{X}_t\|_F^2 + (1 + c)\eta\|\hat{X}_tW - \hat{X}_t\|_F^2 + (1 + 1/c)\eta\hat{\beta}_x^2\|(X_{t+1} - \hat{X}_t) - (\bar{X}_{t+1} - \bar{X}_t)\|_F^2$$
1531 
$$\leq (1 - \eta)\|\hat{X}_t - \hat{X}_t\|_F^2 + (1 + c)\eta\hat{\lambda}^2\|\hat{X}_t - \hat{X}_t\|_F^2 + 2(1 + 1/c)\eta\hat{\beta}_x^2\|X_{t+1} - \bar{X}_{t+1}\|_F^2$$
1532 
$$+ 2(1 + 1/c)\eta\hat{\beta}_x^2\|\hat{X}_t - \hat{X}_t\|_F^2$$
1534 
$$\leq \left(1 - \frac{\eta(1 - \lambda^2)}{4}\right)\|\hat{X}_t - \hat{X}_t\|_F^2 + \frac{4\eta\hat{\beta}_x^2}{1 - \lambda^2}\|X_{t+1} - \bar{X}_{t+1}\|_F^2$$
1536 
$$\leq \left(1 - \frac{\eta(1 - \lambda^2)}{4}\right)\frac{1}{K}\sum_{k=1}^K \mathbb{E}[\|\hat{x}_t - \hat{x}_t^{(k)}\|^2]$$
1537 
$$+ \frac{4\eta\hat{\beta}_x^2}{1 - \lambda^2}\frac{1}{K}\sum_{k=1}^K \mathbb{E}[\|\bar{x}_t - x_t^{(k)}\|^2] + \frac{4\eta\hat{\beta}_x^2}{1 - \lambda^2}\frac{2\eta\beta_x^2}{1 - \lambda^2}\frac{1}{K}\sum_{k=1}^K \mathbb{E}[\|\bar{p}_t - p_t^{(k)}\|^2] ,$$
1540 
$$+ \frac{4\eta\hat{\beta}_x^2}{1 - \lambda^2}\frac{1}{K}\sum_{k=1}^K \mathbb{E}[\|\bar{x}_t - x_t^{(k)}\|^2] + \frac{4\eta\hat{\beta}_x^2}{1 - \lambda^2}\frac{2\eta\beta_x^2}{1 - \lambda^2}\frac{1}{K}\sum_{k=1}^K \mathbb{E}[\|\bar{p}_t - p_t^{(k)}\|^2] ,$$

where  $c=\frac{1-\lambda^2}{2\lambda^2}$  the second to last inequality holds due to  $\hat{\beta}_x \leq \frac{1-\lambda}{4}$ , and the last step holds due to Lemma D.5.

**Lemma D.4.** Given Assumptions 3.1-3.4, when  $\beta_{\hat{y}} \leq \frac{1-\lambda}{4}$ , the following inequality holds:

$$\mathbb{E}[\|\hat{Y}_{t+1} - \bar{\hat{Y}}_{t+1}\|_F^2] \le \left(1 - \frac{\eta(1-\lambda^2)}{4}\right) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{\hat{y}}_t - \hat{y}_t^{(k)}\|^2] \\
+ \frac{4\eta \hat{\beta}_y^2}{1-\lambda^2} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{y}_t - y_t^{(k)}\|^2] + \frac{4\eta \hat{\beta}_y^2}{1-\lambda^2} \frac{2\eta \beta_y^2}{1-\lambda^2} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{q}_t - q_t^{(k)}\|^2].$$
(80)

This lemma be proved by following Lemma D.3. Thus, we omit its proof.

**Lemma D.5.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{x}_{t+1} - x_{t+1}^{(k)}\|^{2}] = \frac{1}{K} \mathbb{E}[\|\bar{X}_{t+1} - X_{t+1}^{(k)}\|_{F}^{2}]$$

$$\leq \left(1 - \frac{\eta(1 - \lambda^{2})}{2}\right) \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{x}_{t} - x_{t}^{(k)}\|^{2}] + \frac{2\eta\beta_{x}^{2}}{1 - \lambda^{2}} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{p}_{t} - p_{t}^{(k)}\|^{2}]. \tag{81}$$

This lemma can be proved by following Lemma D.3. Thus, we omit its proof.

**Lemma D.6.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{y}_{t+1} - y_{t+1}^{(k)}\|^2]$$
(82)

$$\leq \left(1 - \frac{\eta(1-\lambda^2)}{2}\right) \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{y}_t - y_t^{(k)}\|^2] + \frac{2\eta\beta_y^2}{1-\lambda^2} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{q}_t - q_t^{(k)}\|^2]. \tag{83}$$

This lemma can be proved by following Lemma D.3. Thus, we omit its proof.

#### D.2 Gradient Estimation Errors

**Lemma D.7.** Given Assumptions 3.1-3.4, when  $\eta \leq \frac{1}{\sqrt{\rho_x}}$ , the following inequality holds:

$$\mathbb{E}\left[\left\|\frac{1}{K}\sum_{k=1}^{K}u_{t+1}^{(k)} - \frac{1}{K}\sum_{k=1}^{K}\nabla_{x}F^{(k)}(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \hat{x}_{t+1}^{(k)}, \hat{y}_{t+1}^{(k)})\right\|^{2}\right] \\
\leq (1 - \rho_{x}\eta^{2})\mathbb{E}\left[\left\|\frac{1}{K}\sum_{k=1}^{K}u_{t}^{(k)} - \frac{1}{K}\sum_{k=1}^{K}\nabla_{x}F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\right\|^{2}\right] + 2\rho_{x}^{2}\eta^{4}\sigma^{2}\frac{1}{K} \\
+ 4L^{2}\frac{1}{K^{2}}\sum_{k=1}^{K}\mathbb{E}\left[\left\|x_{t+1}^{(k)} - x_{t}^{(k)}\right\|^{2}\right] + 4L^{2}\frac{1}{K^{2}}\sum_{k=1}^{K}\mathbb{E}\left[\left\|y_{t+1}^{(k)} - y_{t}^{(k)}\right\|^{2}\right]. \tag{84}$$

$$\begin{aligned} & \text{Proof.} \\ & \mathbb{E}[\|\frac{1}{K}\sum_{k=1}^{K}u_{t+1}^{(k)} - \frac{1}{K}\sum_{k=1}^{K}\nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)},\hat{y}_{t+1}^{(k)})\|^{2}] \\ & = \mathbb{E}[\|(1-\rho_{x}\eta^{2})(\frac{1}{K}\sum_{k=1}^{K}U_{t}^{(k)} - \frac{1}{K}\sum_{k=1}^{K}\nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}))\|^{2}] \\ & + \mathbb{E}[\|(1-\rho_{x}\eta^{2})\frac{1}{K}\sum_{k=1}^{K}(\nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) - \nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)};\hat{x}_{t+1}^{(k)}) \\ & + \nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)},\hat{y}_{t+1}^{(k)};\xi_{t+1}^{(k)}) - \nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)}) \\ & + \nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)},y_{t+1}^{(k)};\xi_{t+1}^{(k)}) - \nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)},\hat{y}_{t+1}^{(k)};\hat{x}_{t+1}^{(k)}) \\ & + \rho_{x}\eta^{2}\frac{1}{K}\sum_{k=1}^{K}(\nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\xi_{t+1}^{(k)}) - \nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)},\hat{y}_{t+1}^{(k)}) \\ & + \rho_{x}\eta^{2}\frac{1}{K}\sum_{k=1}^{K}U_{t}^{(k)} - \frac{1}{K}\sum_{k=1}^{K}\nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) \\ & + 2(1-\rho_{x}\eta^{2})(\frac{1}{K}\sum_{k=1}^{K}U_{t}^{(k)} - \frac{1}{K}\sum_{k=1}^{K}\nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) \\ & + 2(1-\rho_{x}\eta^{2})^{2}\frac{1}{K^{2}}\sum_{k=1}^{K}\mathbb{E}[\|\nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) - \nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) \\ & + 2\rho_{x}\eta^{4}\frac{1}{K^{2}}\sum_{k=1}^{K}\mathbb{E}[\|\nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)}) - \nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) \\ & + 2\rho_{x}^{2}\eta^{4}\frac{1}{K^{2}}\sum_{k=1}^{K}\mathbb{E}[\|\nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)},\hat{y}_{t+1}^{(k)};\hat{x}_{t+1}^{(k)}) - \nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) \\ & + 2\rho_{x}^{2}\eta^{4}\frac{1}{K^{2}}\sum_{k=1}^{K}\mathbb{E}[\|\nabla_{x}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)},\hat{y}_{t+1}^{(k)};\hat{x}_{t+1}^{(k)}) - \nabla_{x}F^{(k)}(x_{t}^{(k)}$$

 where the last step holds due to the following inequality:

$$\mathbb{E}[\|\nabla_{x}F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)}) - \nabla_{x}F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)}; \xi_{t+1}^{(k)}) \\
+ \nabla_{x}F^{(k)}(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \hat{x}_{t+1}^{(k)}, \hat{y}_{t+1}^{(k)}; \xi_{t+1}^{(k)}) - \nabla_{x}F^{(k)}(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \hat{x}_{t+1}^{(k)}, \hat{y}_{t+1}^{(k)})\|^{2}] \\
= \mathbb{E}[\|\nabla_{x}f(x_{t}^{(k)}, y_{t}^{(k)}) + \gamma_{1}(x_{t}^{(k)} - \hat{x}_{t}^{(k)}) - \nabla_{x}f(x_{t}^{(k)}, y_{t}^{(k)}; \xi_{t+1}^{(k)}) - \gamma_{1}(x_{t}^{(k)} - \hat{x}_{t}^{(k)}) \\
+ \nabla_{x}f(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \xi_{t+1}^{(k)}) + \gamma_{1}(x_{t+1}^{(k)} - \hat{x}_{t+1}^{(k)}) - \nabla_{x}f(x_{t+1}^{(k)}, y_{t+1}^{(k)}) - \gamma_{1}(x_{t+1}^{(k)} - \hat{x}_{t+1}^{(k)})\|^{2}] \\
= \mathbb{E}[\|\nabla_{x}f(x_{t}^{(k)}, y_{t}^{(k)}) - \nabla_{x}f(x_{t}^{(k)}, y_{t}^{(k)}; \xi_{t+1}^{(k)}) \\
+ \nabla_{x}f(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \xi_{t+1}^{(k)}) - \nabla_{x}f(x_{t+1}^{(k)}, y_{t+1}^{(k)})\|^{2}] \\
\leq \mathbb{E}[\|\nabla_{x}f(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \xi_{t+1}^{(k)}) - \nabla_{x}f(x_{t}^{(k)}, y_{t}^{(k)}; \xi_{t+1}^{(k)})\|^{2}] \\
\leq 2L^{2}\mathbb{E}[\|x_{t+1}^{(k)} - x_{t}^{(k)}\|^{2}] + 2L^{2}\mathbb{E}[\|y_{t+1}^{(k)} - y_{t}^{(k)}\|^{2}]. \tag{86}$$

**Lemma D.8.** Given Assumptions 3.1-3.4, when  $\eta \leq \frac{1}{\sqrt{\rho_x}}$ , the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|u_{t+1}^{(k)} - \nabla_{x}F^{(k)}(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \hat{x}_{t+1}^{(k)}, \hat{y}_{t+1}^{(k)})\|^{2}]$$

$$\leq (1 - \rho_{x}\eta^{2}) \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|u_{t}^{(k)} - \nabla_{x}F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]$$

$$+ 4L^{2} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|x_{t+1}^{(k)} - x_{t}^{(k)}\|^{2}] + 4L^{2} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|y_{t+1}^{(k)} - y_{t}^{(k)}\|^{2}] + 2\rho_{x}^{2}\eta^{4}\sigma^{2} . \tag{87}$$

This lemma can be proved by following Lemma D.7. Thus, we omit its proof.

**Lemma D.9.** Given Assumptions 3.1-3.4, when  $\eta \leq \frac{1}{\sqrt{\rho_{\eta}}}$ , the following inequality holds:

$$\mathbb{E}\left[\left\|\frac{1}{K}\sum_{k=1}^{K}\nabla_{y}F^{(k)}(x_{t+1}^{(k)},y_{t+1}^{(k)};\hat{x}_{t+1}^{(k)},\hat{y}_{t+1}^{(k)}) - \frac{1}{K}\sum_{k=1}^{K}v_{t+1}^{(k)}\right\|^{2}\right] \\
\leq (1 - \rho_{y}\eta^{2})\mathbb{E}\left[\left\|\frac{1}{K}\sum_{k=1}^{K}v_{t}^{(k)} - \frac{1}{K}\sum_{k=1}^{K}\nabla_{y}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)})\right\|^{2}\right] \\
+ 4L^{2}\frac{1}{K^{2}}\sum_{k=1}^{K}\mathbb{E}\left[\left\|x_{t+1}^{(k)} - x_{t}^{(k)}\right\|^{2}\right] + 4L^{2}\frac{1}{K^{2}}\sum_{k=1}^{K}\mathbb{E}\left[\left\|y_{t+1}^{(k)} - y_{t}^{(k)}\right\|^{2}\right] + 2\rho_{y}^{2}\eta^{4}\sigma^{2}\frac{1}{K}. \tag{88}$$

This lemma can be proved by following Lemma D.7. Thus, we omit its proof.

**Lemma D.10.** Given Assumptions 3.1-3.4, when  $\eta \leq \frac{1}{\sqrt{\rho_y}}$ , the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\nabla_{y} F^{(k)}(x_{t+1}^{(k)}, y_{t+1}^{(k)}; \hat{x}_{t+1}^{(k)}, \hat{y}_{t+1}^{(k)}) - v_{t+1}^{(k)}\|^{2}]$$

$$\leq (1 - \rho_{y} \eta^{2}) \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\|^{2}]$$

$$+ 4L^{2} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|x_{t+1}^{(k)} - x_{t}^{(k)}\|^{2}] + 4L^{2} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|y_{t+1}^{(k)} - y_{t}^{(k)}\|^{2}] + 2\rho_{y}^{2} \eta^{4} \sigma^{2} . \tag{89}$$

Similarly, this lemma can be proved by following Lemma D.7. Thus, we omit its proof.

#### D.3 OTHER AUXILIARY LEMMAS

**Lemma D.11.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\hat{x}_{t+1}^{(k)} - \hat{x}_{t}^{(k)}\|^{2}] \leq 3\mathbb{E}[\|\bar{\hat{x}}_{t+1} - \bar{\hat{x}}_{t}\|^{2}] + 6\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{\hat{x}}_{t} - \hat{x}_{t}^{(k)}\|^{2}] 
+ \frac{12\eta\hat{\beta}_{x}^{2}}{1 - \lambda^{2}} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{x}_{t} - x_{t}^{(k)}\|^{2}] + \frac{12\eta\hat{\beta}_{x}^{2}}{1 - \lambda^{2}} \frac{2\eta\beta_{x}^{2}}{1 - \lambda^{2}} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{p}_{t} - p_{t}^{(k)}\|^{2}].$$
(90)

Proof.

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\hat{x}_{t+1}^{(k)} - \hat{x}_{t}^{(k)}\|^{2}]$$

$$= \frac{1}{K} \mathbb{E}[\|\hat{X}_{t+1} - \hat{X}_{t}\|_{F}^{2}]$$

$$\leq 3 \frac{1}{K} \mathbb{E}[\|\hat{X}_{t+1} - \hat{X}_{t+1}\|_{F}^{2}] + 3 \frac{1}{K} \mathbb{E}[\|\hat{X}_{t} - \hat{X}_{t}\|_{F}^{2}] + 3 \frac{1}{K} \mathbb{E}[\|\hat{X}_{t+1} - \hat{X}_{t}\|_{F}^{2}]$$

$$\leq 3 \mathbb{E}[\|\hat{x}_{t+1} - \hat{x}_{t}\|^{2}] + 6 \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\hat{x}_{t} - \hat{x}_{t}^{(k)}\|^{2}]$$

$$+ \frac{12\eta \hat{\beta}_{x}^{2}}{1 - \lambda^{2}} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{x}_{t} - x_{t}^{(k)}\|^{2}] + \frac{12\eta \hat{\beta}_{x}^{2}}{1 - \lambda^{2}} \frac{2\eta \beta_{x}^{2}}{1 - \lambda^{2}} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{p}_{t} - p_{t}^{(k)}\|^{2}], \tag{91}$$

where the last step holds due to Lemma D.3.

**Lemma D.12.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\hat{y}_{t+1}^{(k)} - \hat{y}_{t}^{(k)}\|^{2}] \leq 3\mathbb{E}[\|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_{t}\|^{2}] + 6\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{\hat{y}}_{t} - \hat{y}_{t}^{(k)}\|^{2}] \\
+ \frac{12\eta \hat{\beta}_{y}^{2}}{1 - \lambda^{2}} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{y}_{t} - y_{t}^{(k)}\|^{2}] + \frac{12\eta \hat{\beta}_{y}^{2}}{1 - \lambda^{2}} \frac{2\eta \beta_{y}^{2}}{1 - \lambda^{2}} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{q}_{t} - q_{t}^{(k)}\|^{2}].$$
(92)

This lemma can be proved by following Lemma D.11. Thus, we omit its proof.

**Lemma D.13.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|x_{t+1}^{(k)} - x_{t}^{(k)}\|^{2}] \leq 12\eta^{2} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{x}_{t} - x_{t}^{(k)}\|^{2}] 
+ 3\beta_{x}^{2} \eta^{2} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{p}_{t} - p_{t}^{(k)}\|^{2}] + 3\beta_{x}^{2} \eta^{2} \mathbb{E}[\|\bar{u}_{t}\|^{2}].$$
(93)

**Lemma D.14.** Given Assumptions 3.1-3.4, the following inequality holds:

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|y_{t+1}^{(k)} - y_{t}^{(k)}\|^{2}] \leq 12\eta^{2} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{y}_{t} - y_{t}^{(k)}\|^{2}] 
+ 3\beta_{y}^{2} \eta^{2} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{q}_{t} - q_{t}^{(k)}\|^{2}] + 3\beta_{y}^{2} \eta^{2} \mathbb{E}[\|\bar{v}_{t}\|^{2}].$$
(94)

Lemmas D.13, D.14 can be proved by following (Gao, 2022).

## E PROOF OF THEOREM 4.2

We first propose a novel potential function as follows:

$$\mathcal{L}_{t} = \mathcal{P}_{t} + c_{1} \mathbb{E}\left[\left\|\frac{1}{K} \sum_{k=1}^{K} u_{t}^{(k)} - \frac{1}{K} \sum_{k=1}^{K} \nabla_{x} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\right\|^{2}\right] \\
+ c_{2} \mathbb{E}\left[\left\|\frac{1}{K} \sum_{k=1}^{K} v_{t}^{(k)} - \frac{1}{K} \sum_{k=1}^{K} \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\right\|^{2}\right] \\
+ c_{3} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}\left[\left\|\bar{x}_{t} - x_{t}^{(k)}\right\|^{2}\right] + c_{4} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}\left[\left\|\bar{y}_{t} - y_{t}^{(k)}\right\|^{2}\right] + c_{5} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}\left[\left\|\bar{x}_{t} - \hat{x}_{t}^{(k)}\right\|^{2}\right] \\
+ c_{10} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}\left[\left\|\bar{y}_{t} - \hat{y}_{t}^{(k)}\right\|^{2}\right] + c_{6} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}\left[\left\|\bar{p}_{t} - p_{t}^{(k)}\right\|^{2}\right] + c_{7} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}\left[\left\|\bar{q}_{t} - q_{t}^{(k)}\right\|^{2}\right] \\
+ c_{8} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}\left[\left\|u_{t}^{(k)} - \nabla_{x} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\right\|^{2}\right] \\
+ c_{9} \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}\left[\left\|v_{t}^{(k)} - \nabla_{y} F^{(k)}(x_{t}^{(k)}, y_{t}^{(k)}; \hat{x}_{t}^{(k)}, \hat{y}_{t}^{(k)})\right\|^{2}\right], \tag{95}$$

where the coefficient  $\{c_i\}_{i=1}^9$  are positive.

Since

$$\mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\hat{\bar{x}}_{t},\bar{\hat{y}}_{t}) - \bar{u}_{t}\|^{2}]$$

$$\leq 2\mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\hat{\bar{x}}_{t},\bar{\hat{y}}_{t}) - \frac{1}{K}\sum_{k=1}^{K}\nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)})\|^{2}]$$

$$+ 2\mathbb{E}[\|\frac{1}{K}\sum_{k=1}^{K}\nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) - \bar{u}_{t}\|^{2}]$$

$$\leq 2L^{2}\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{x}_{t} - x_{t}^{(k)}\|^{2}] + 2L^{2}\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{y}_{t} - y_{t}^{(k)}\|^{2}]$$

$$+ 2\mathbb{E}[\|\frac{1}{K}\sum_{k=1}^{K}\nabla_{x}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) - \frac{1}{K}\sum_{k=1}^{K}u_{t}^{(k)}\|^{2}],$$
(96)

and

$$\mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{x}_{t},\bar{y}_{t}) - \bar{v}_{t}\|^{2}] \leq 2L^{2}\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{x}_{t} - x_{t}^{(k)}\|^{2}] + 2L^{2}\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{y}_{t} - y_{t}^{(k)}\|^{2}] + 2\mathbb{E}[\|\frac{1}{K}\sum_{k=1}^{K}\nabla_{y}F^{(k)}(x_{t}^{(k)},y_{t}^{(k)};\hat{x}_{t}^{(k)},\hat{y}_{t}^{(k)}) - \bar{v}_{t}\|^{2}],$$
(97)

we obtain

$$\mathcal{L}_{t+1} - \mathcal{L}_t \leq -\frac{\beta_x \eta}{4} \mathbb{E}[\|\nabla_x F(\bar{x}_t, \bar{y}_t; \hat{\bar{x}}_t, \bar{y}_t)\|^2] - \frac{\beta_y \eta}{2} \mathbb{E}[\|\nabla_y F(\bar{x}_t, \bar{y}_t; \hat{\bar{x}}_t, \bar{y}_t)\|^2]$$

$$+ (\beta_x \eta - c_1 \rho_x \eta^2) \mathbb{E}[\|\frac{1}{K} \sum_{k=1}^K \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) - \frac{1}{K} \sum_{k=1}^K u_t^{(k)}\|^2]$$

$$+ (2A_3 - \rho_y \eta^2 c_2) \mathbb{E}[\|\frac{1}{K} \sum_{k=1}^K \nabla_y f^{(k)}(x_t^{(k)}, y_t^{(k)}) - \frac{1}{K} \sum_{k=1}^K v_t^{(k)}\|^2]$$

$$+ (2A_3 - \rho_y \eta^2 c_2) \mathbb{E}[\|\frac{1}{K} \sum_{k=1}^K \nabla_y f^{(k)}(x_t^{(k)}, y_t^{(k)}) - \frac{1}{K} \sum_{k=1}^K v_t^{(k)}\|^2]$$

$$\begin{split} &+\left(4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2}-\frac{\beta_{x}\eta}{4}\right)\mathbb{E}[\|\bar{u}_{t}\|^{2}] \\ &+\left(\beta_{y}^{2}\eta^{2}L_{d}+\frac{3\beta_{y}\eta}{4}+\frac{\beta_{y}^{2}\eta^{2}(\gamma_{2}+L)}{2}+4A_{1}\hat{\beta}_{y}^{2}\eta^{2}\beta_{y}^{2}\eta^{2}+2A_{2}\beta_{y}^{2}\eta^{2}-\frac{7}{8}\beta_{y}\eta\right)\mathbb{E}[\|\bar{v}_{t}\|^{2}] \\ &+\left(2\gamma_{1}C_{x_{\frac{1}{2}y}}+\frac{\gamma_{1}}{6\hat{\beta}_{x}\eta}+6\gamma_{1}\hat{\beta}_{x}\eta\left(10C_{x_{\frac{1}{2}y}}^{2}C_{y_{\frac{1}{2}y}}^{2}+\frac{4}{\gamma_{1}-L}\frac{2\gamma_{2}^{2}C_{y_{\frac{1}{2}y}}^{2}}{\mu}\right)-\frac{\gamma_{1}(2-\hat{\beta}_{x}\eta)}{2\hat{\beta}_{x}\eta}\right)\mathbb{E}[\|\bar{x}_{t+1}-\bar{x}_{t}\|^{2}] \\ &+\left(2A_{1}-\frac{\gamma_{2}(2-\hat{\beta}_{y}\eta)}{2\hat{\beta}_{y}\eta}\right)\mathbb{E}[\|\hat{b}_{t+1}-\bar{b}_{t}\|^{2}] \\ &+\left(\frac{4L^{2}c_{1}}{K}+\frac{4L^{2}c_{2}}{K}+\frac{9(L^{2}+\gamma_{1}^{2})c_{6}}{1-\lambda}+\frac{9L^{2}c_{7}}{1-\lambda}+4L^{2}c_{8}+4L^{2}c_{9}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|x_{t+1}^{(k)}-x_{t}^{(k)}\|^{2}] \\ &+\left(\frac{4L^{2}c_{1}}{K}+\frac{4L^{2}c_{2}}{K}+\frac{9L^{2}c_{5}}{1-\lambda}+\frac{9(L^{2}+\gamma_{2}^{2})c_{7}}{1-\lambda}+4L^{2}c_{8}+4L^{2}c_{9}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|y_{t+1}^{(k)}-y_{t}^{(k)}\|^{2}] \\ &+\left(\frac{9\gamma_{1}^{2}}{1-\lambda}c_{6}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\hat{x}_{t+1}^{(k)}-\hat{x}_{t}^{(k)}\|^{2}]+\left(\frac{9\gamma_{2}^{2}}{1-\lambda}c_{7}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\hat{y}_{t+1}^{(k)}-\hat{y}_{t}^{(k)}\|^{2}] \\ &+\left(\beta_{x}\eta L^{2}+2L^{2}A_{3}+\frac{4\eta\hat{\beta}_{x}^{2}}{1-\lambda^{2}}c_{5}-\frac{\eta(1-\lambda^{2})}{2}c_{3}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{y}_{t}-y_{t}^{(k)}\|^{2}] \\ &+\left(\frac{2\eta\beta_{x}^{2}}{1-\lambda^{2}}c_{3}+\frac{8\eta^{2}\beta_{x}^{2}\hat{\beta}_{x}^{2}}{(1-\lambda^{2})^{2}}c_{10}-\frac{\eta(1-\lambda^{2})}{2}c_{4}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{y}_{t}-y_{t}^{(k)}\|^{2}] \\ &+\left(\frac{2\eta\beta_{x}^{2}}{2}c_{3}+\frac{8\eta^{2}\beta_{x}^{2}\hat{\beta}_{x}^{2}}{(1-\lambda^{2})^{2}}c_{10}-\frac{\eta(1-\lambda^{2})}{2}c_{4}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{y}_{t}-y_{t}^{(k)}\|^{2}] \\ &+\left(\frac{2\eta\beta_{x}^{2}}{2}c_{4}+\frac{8\eta^{2}\beta_{x}^{2}\hat{\beta}_{x}^{2}}{(1-\lambda^{2})^{2}}c_{10}-(1-\lambda)c_{6}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{y}_{t}-y_{t}^{(k)}\|^{2}] \\ &+\left(\frac{2\eta\beta_{x}^{2}}{2}c_{4}+\frac{8\eta^{2}\beta_{x}^{2}\hat{\beta}_{x}^{2}}{(1-\lambda^{2})^{2}}c_{10}-(1-\lambda)c_{7}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{y}_{t}-y_{t}^{(k)}\|^{2}] \\ &+\left(\frac{2\eta\beta_{x}^{2}}}{1-\lambda^{2}}c_{4}+\frac{8\eta^{2}\beta_{x}^{2}\hat{\beta}_{x}^{2}}{(1-\lambda^{2})^{2}}c_{10}-(1-\lambda)c_{7}\right)\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}[\|\bar{y}_{t}-y_{t}^{(k)}\|^{2}] \\ &+\left(\frac{2\eta\beta_{x}^{2}}{2}\eta^{4}\frac{1}{1-\lambda}c_{7}-\rho_{y}\eta^{2}c_{8}\right)\frac{1}{K}\sum_{k=1}^{K$$

By setting

$$\mathcal{X} = \frac{4L^2c_1}{K} + \frac{4L^2c_2}{K} + \frac{9(L^2 + \gamma_1^2)c_6}{1 - \lambda} + \frac{9L^2c_7}{1 - \lambda} + 4L^2c_8 + 4L^2c_9, 
\mathcal{Y} = \frac{4L^2c_1}{K} + \frac{4L^2c_2}{K} + \frac{9L^2c_6}{1 - \lambda} + \frac{9(L^2 + \gamma_2^2)c_7}{1 - \lambda} + 4L^2c_8 + 4L^2c_9,$$
(99)

and due to  $\lambda < 1$ , we obtain  $\frac{1}{1-\lambda^2} \leq \frac{1}{1-\lambda}$ , and further derive

$$\mathcal{L}_{t+1} - \mathcal{L}_t \leq -\frac{\beta_x \eta}{4} \mathbb{E}[\|\nabla_x F(\bar{x}_t, \bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t)\|^2] - \frac{\beta_y \eta}{2} \mathbb{E}[\|\nabla_y F(\bar{x}_t, \bar{y}_t; \bar{\hat{x}}_t, \bar{\hat{y}}_t)\|^2]$$

$$\begin{aligned} &+ (\beta_x \eta - c_1 \rho_x \eta^2) \mathbb{E}[\|\frac{1}{K} \sum_{k=1}^K \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)}) - \frac{1}{K} \sum_{k=1}^K u_t^{(k)}\|^2] \\ &+ \left(2\beta_y \eta + 240\gamma_1 \beta_x \eta \beta_y^2 \eta^2 C_{x_y^2 xy}^2 + 8A_2 \beta_y^2 \eta^2 - \rho_y \eta^2 c_2\right) \mathbb{E}[\|\frac{1}{K} \sum_{k=1}^K \nabla_y f^{(k)}(x_t^{(k)}, y_t^{(k)}) - \frac{1}{K} \sum_{k=1}^K v_t^{(k)}\|^2] \\ &+ \left(4\beta_y \eta \beta_x^2 \eta^2 L^2 + 3\beta_x^2 \eta^2 \mathcal{X} - \frac{\beta_x \eta}{4}\right) \mathbb{E}[\|\tilde{u}_t\|^2] \\ &+ \left(\beta_y^2 \eta^2 L_d + \frac{3\beta_y \eta}{4} + \frac{\beta_y^2 \eta^2 (\gamma_2 + L)}{2} + 4A_1 \hat{\beta}_y^2 \eta^2 \beta_y^2 \eta^2 + 2A_2 \beta_y^2 \eta^2 + 3\beta_y^2 \eta^2 \mathcal{Y} - \frac{7}{8}\beta_y \eta\right) \mathbb{E}[\|\tilde{v}_t\|^2] \\ &+ \left(2\gamma_1 C_{x_y^1 y} + \frac{\gamma_1}{6\beta_x \eta} + 6\gamma_1 \hat{\beta}_x \eta \left(10C_{x_y^2 x}^2 C_{y_y^2 x}^2 + \frac{4}{\gamma_1 - L} \frac{2\gamma_2^2 C_{y_y^2 x}^2}{\mu}\right) + \frac{27\gamma_1^2}{1 - \lambda} c_6 \\ &- \frac{\gamma_1 (2 - \hat{\beta}_x \eta)}{2\hat{\beta}_x \eta}\right) \mathbb{E}[\|\tilde{k}_{t+1} - \tilde{k}_t\|^2] \\ &+ \left(2A_1 + \frac{7\gamma_2}{1 - \lambda} c_7 - \frac{\gamma_2 (2 - \hat{\beta}_y \eta)}{2\beta_y \eta}\right) \mathbb{E}[\|\tilde{b}_{t+1} - \tilde{y}_t\|^2] \\ &+ \left(\beta_x \eta L^2 + 2L^2 A_3 + \frac{4\eta \beta_y^2}{1 - \lambda} c_5 + \frac{108\eta \beta_y^2 \gamma_2^2}{(1 - \lambda)^2} c_6 + 12\eta^2 \mathcal{X} - \frac{\eta (1 - \lambda^2)}{2} c_3\right) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\tilde{y}_t - y_t^{(k)}\|^2] \\ &+ \left(\beta_x \eta L^2 + 2L^2 A_3 + \frac{4\eta \beta_y^2}{1 - \lambda} c_1 + \frac{108\eta \beta_y^2 \gamma_2^2}{(1 - \lambda)^2} c_7 + 12\eta^2 \mathcal{Y} - \frac{\eta (1 - \lambda^2)}{2} c_4\right) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\tilde{y}_t - y_t^{(k)}\|^2] \\ &+ \left(\frac{2\eta \beta_y^2}{1 - \lambda} c_3 + \frac{8\eta^2 \beta_y^2 \beta_y^2}{(1 - \lambda)^2} c_5 + \frac{216\eta^2 \beta_y^2 \beta_y^2 \gamma_2^2}{(1 - \lambda)^3} c_6 + 3\beta_x^2 \eta^2 \mathcal{X} - (1 - \lambda) c_6\right) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\tilde{y}_t - y_t^{(k)}\|^2] \\ &+ \left(\frac{2\eta \beta_y^2}{1 - \lambda} c_4 - \frac{8\eta^2 \beta_y^2 \beta_y^2}{(1 - \lambda)^2} c_7 + \frac{12\eta^2 \mathcal{Y}}{2} c_7 + 3\beta_y^2 \eta^2 \mathcal{Y} - (1 - \lambda) c_6\right) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\tilde{y}_t - y_t^{(k)}\|^2] \\ &+ \left(\frac{3\rho_y^2 \eta^4}{1 - \lambda} c_6 - \rho_x \eta^2 c_8\right) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\tilde{y}_t - \hat{y}_t^{(k)}\|^2] \\ &+ \left(\frac{3\rho_y^2 \eta^4}{1 - \lambda} c_7 - \rho_y \eta^2 c_9\right) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\tilde{y}_t - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)})\|^2] \\ &+ \left(\frac{3\rho_y^2 \eta^4}{1 - \lambda} c_7 - \rho_y \eta^2 c_9\right) \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\tilde{y}_t^{(k)} - \nabla_y F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)})\|^2] \\ &+ \left(\frac{3\rho_y^2 \eta^4}{1 - \lambda} c_7 - \rho_y$$

To cancel out  $\mathbb{E}[\|\frac{1}{K}\sum_{k=1}^K \nabla_x F^{(k)}(x_t^{(k)},y_t^{(k)};\hat{x}_t^{(k)},\hat{y}_t^{(k)}) - \frac{1}{K}\sum_{k=1}^K u_t^{(k)}\|^2]$ , i.e.,

$$\beta_x \eta - \rho_x \eta^2 c_1 \le 0. \tag{101}$$

Then, we set

$$c_1 = \frac{\beta_x}{\rho_x \eta} \,. \tag{102}$$

To cancel out  $\mathbb{E}[\|\frac{1}{K}\sum_{k=1}^K \nabla_y f^{(k)}(x_t^{(k)},y_t^{(k)}) - \frac{1}{K}\sum_{k=1}^K v_t^{(k)}\|^2]$ , i.e.,

$$2\beta_y \eta + 240\gamma_1 \hat{\beta}_x \eta \beta_y^2 \eta^2 C_{x_{y\hat{x}\hat{y}}}^2 + 8A_2 \beta_y^2 \eta^2 - \rho_y \eta^2 c_2 \le 0.$$
 (103)

Specifically, since the second and last inequality in Eq. (66) holds, we have

$$240\gamma_1 \hat{\beta}_x \eta \beta_y^2 \eta^2 C_{x_y^2 \hat{y}}^2 \le \frac{2\beta_x \eta}{32 \times 16} \frac{(\gamma_1 - L)^2}{L^2} ,$$

$$8A_2 \beta_y^2 \eta^2 \le \frac{2\beta_x \eta}{32 \times 16} \frac{(\gamma_1 - L)^2}{L^2} . \tag{104}$$

Then, by the definition of  $c_{\beta_y}$ , i.e.,  $c_{\beta_y} = \frac{(\gamma_1 - L)^2}{64L^2}$ , we set

$$240\gamma_{1}\hat{\beta}_{x}\eta\beta_{y}^{2}\eta^{2}C_{x_{y\hat{x}\hat{y}}^{2}}^{2} \leq \frac{2\beta_{x}\eta}{32\times16}64c_{\beta_{y}} = \frac{1}{4}\beta_{y}\eta,$$

$$8A_{2}\beta_{y}^{2}\eta^{2} \leq \frac{2\beta_{x}\eta}{32\times16}64c_{\beta_{y}} = \frac{\beta_{y}\eta}{4}.$$
(105)

Therefore, we obtain

$$c_2 = \frac{5\beta_y}{2\rho_y \eta} \,. \tag{106}$$

To cancel out  $\frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|u_t^{(k)} - \nabla_x F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)})\|^2]$ , i.e.,

$$\frac{3\rho_x^2\eta^4}{1-\lambda}c_6 - \rho_x\eta^2c_8 \le 0. {107}$$

Here, because  $\rho_x \eta^2 < 1$ , we set

$$c_6 = \beta_x \eta (1 - \lambda) , \quad c_8 = 3\beta_x \eta .$$
 (108)

Similarly, to cancel out  $\frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|v_t^{(k)} - \nabla_y F^{(k)}(x_t^{(k)}, y_t^{(k)}; \hat{x}_t^{(k)}, \hat{y}_t^{(k)})\|^2]$ , i.e.,

$$\frac{3\rho_y^2\eta^4}{1-\lambda}c_7 - \rho_y\eta^2c_9 \le 0. {109}$$

Because  $\rho_y \eta^2 < 1$ , we set

$$c_7 = \beta_u \eta (1 - \lambda) , \quad c_9 = 3\beta_u \eta .$$
 (110)

To cancel out  $\frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{\hat{x}}_t - \hat{x}_t^{(k)}\|^2]$ , i.e.,

$$\frac{54\gamma_1^2}{1-\lambda}c_6 - \frac{\eta(1-\lambda^2)}{4}c_5 \le 0, \tag{111}$$

we set

$$c_5 = \frac{216\beta_x \gamma_1^2}{(1-\lambda)} \,. \tag{112}$$

To cancel out  $\frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\hat{\bar{y}}_t - \hat{y}_t^{(k)}\|^2]$ , i.e.,

$$\frac{54\gamma_2^2}{1-\lambda}c_7 - \frac{\eta(1-\lambda^2)}{4}c_{10} \le 0, \qquad (113)$$

we set

$$c_{10} = \frac{216\beta_y \gamma_2^2}{(1-\lambda)} \ . \tag{114}$$

To cancel out  $\frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{x}_t - x_t^{(k)}\|^2]$ , i.e.,

$$\beta_x \eta L^2 + 2L^2 A_3 + \frac{4\eta \hat{\beta}_x^2}{1-\lambda} c_5 + \frac{108\eta \hat{\beta}_x^2 \gamma_1^2}{(1-\lambda)^2} c_6 + 12\eta^2 \mathcal{X} - \frac{\eta (1-\lambda^2)}{2} c_3 \le 0.$$
 (115)

Firstly, from the definition of  $\mathcal{X}$ , we have

$$\mathcal{X} = \frac{4L^{2}c_{1}}{K} + \frac{4L^{2}c_{2}}{K} + \frac{9(L^{2} + \gamma_{1}^{2})c_{6}}{1 - \lambda} + \frac{9L^{2}c_{7}}{1 - \lambda} + 4L^{2}c_{8} + 4L^{2}c_{9} ,$$

$$= \frac{4L^{2}}{K} \frac{\beta_{x}}{\rho_{x}\eta} + \frac{4L^{2}}{K} \frac{5\beta_{y}}{2\rho_{y}\eta} + 9(L^{2} + \gamma_{1}^{2})\beta_{x}\eta + 9L^{2}\beta_{y}\eta + 12L^{2}\beta_{x}\eta + 12L^{2}\beta_{y}\eta ,$$

$$= \frac{4L^{2}}{K} \frac{\beta_{x}}{\rho_{x}\eta} + \frac{10L^{2}}{K} \frac{\beta_{y}}{\rho_{y}\eta} + (21L^{2} + 9\gamma_{1}^{2})\beta_{x}\eta + 21L^{2}\beta_{y}\eta .$$
(116)

Moreover, from the definition of  $A_3$  and Eq. (103), we have

$$\beta_y \eta + 120 \gamma_1 \hat{\beta}_x \eta \beta_y^2 \eta^2 C_{x_{y \hat{x} \hat{y}}}^2 + 4A_2 \beta_y^2 \eta^2 \le \frac{5}{4} \beta_y \eta \tag{117}$$

Therefore, we set

$$\beta_{x}\eta L^{2} + 2L^{2}A_{3} + \frac{4\eta\hat{\beta}_{x}^{2}}{1-\lambda}c_{5} + \frac{108\eta\hat{\beta}_{x}^{2}\gamma_{1}^{2}}{(1-\lambda)^{2}}c_{6} + 12\eta^{2}\mathcal{X}$$

$$\leq \beta_{x}\eta L^{2} + \frac{5}{2}\beta_{y}\eta L^{2} + \frac{4\eta\hat{\beta}_{x}^{2}}{1-\lambda}\frac{216\beta_{x}\gamma_{1}^{2}}{(1-\lambda)} + \frac{108\gamma_{1}^{2}\eta\hat{\beta}_{x}^{2}}{(1-\lambda)}\beta_{x}\eta$$

$$+ 12\eta^{2}\left(\frac{4L^{2}}{K}\frac{\beta_{x}}{\rho_{x}\eta} + \frac{10L^{2}}{K}\frac{\beta_{y}}{\rho_{y}\eta} + (21L^{2} + 9\gamma_{1}^{2})\beta_{x}\eta + 21L^{2}\beta_{y}\eta\right)$$

$$\leq \beta_{x}\eta L^{2} + \beta_{x}\eta L^{2}\frac{5\beta_{y}}{2\beta_{x}} + \beta_{x}\eta\hat{\beta}_{x}^{2}\frac{864\gamma_{1}^{2}}{(1-\lambda)^{2}} + \beta_{x}\eta\hat{\beta}_{x}\frac{108\gamma_{1}^{2}}{(1-\lambda)}$$

$$+ 12\eta\left(\frac{4L^{2}}{\rho_{x}K}\beta_{x} + \frac{10L^{2}}{\rho_{y}K}\beta_{y} + (21L^{2} + 9\gamma_{1}^{2})\beta_{x} + 21L^{2}\beta_{y}\right)$$

$$= \beta_{x}\eta L^{2} + \beta_{x}\eta L^{2}\frac{5}{2}c_{\beta_{y}} + \beta_{x}^{3}\eta c_{\hat{\beta}_{x}}^{2}\frac{864\gamma_{1}^{2}}{(1-\lambda)^{2}} + \beta_{x}^{2}\eta c_{\hat{\beta}_{x}}\frac{108\gamma_{1}^{2}}{(1-\lambda)}$$

$$+ 12\beta_{x}\eta\left(\frac{4L^{2}}{\rho_{x}K} + \frac{10L^{2}}{\rho_{y}K}c_{\beta_{y}} + (21L^{2} + 9\gamma_{1}^{2}) + 21L^{2}c_{\beta_{y}}\right)$$

$$\leq \frac{\eta(1-\lambda)}{2}c_{3}, \tag{118}$$

where the second step holds due to  $\hat{\beta}_x \eta < 1$  and  $\eta < 1$ , the fourth step holds due to Eq. (54). By solving this inequality, we obtain

$$c_{3} \geq \frac{2\beta_{x}}{(1-\lambda)} \left( \frac{48L^{2}}{\rho_{x}K} + \frac{120L^{2}}{\rho_{y}K} c_{\beta_{y}} + 253L^{2} + 108\gamma_{1}^{2} + 255L^{2}c_{\beta_{y}} + \beta_{x}^{2}c_{\hat{\beta}_{x}}^{2} \frac{864\gamma_{1}^{2}}{(1-\lambda)^{2}} + \beta_{x}c_{\hat{\beta}_{x}} \frac{108\gamma_{1}^{2}}{(1-\lambda)} \right). \tag{119}$$

Then, we set

$$c_{3} = \frac{2\beta_{x}}{(1-\lambda)} \left( \underbrace{\frac{48L^{2}}{\rho_{x}K} + \frac{120L^{2}}{\rho_{y}K} c_{\beta_{y}} + 253L^{2} + 108\gamma_{1}^{2} + 302L^{2}c_{\beta_{y}}}_{c_{3,1}} \right) + \underbrace{\frac{2\beta_{x}^{3}}{(1-\lambda)^{3}}}_{c_{3,2}} \underbrace{\frac{864\gamma_{1}^{2}c_{\hat{\beta}_{x}}^{2}}{(1-\lambda)^{2}}}_{c_{3,2}} + \underbrace{\frac{2\beta_{x}^{2}}{(1-\lambda)^{2}}}_{c_{3,3}} \underbrace{\frac{108\gamma_{1}^{2}c_{\hat{\beta}_{x}}}{c_{3,3}}}_{c_{3,3}}.$$

$$(120)$$

Here, it is easy to know that  $c_{3,1}=O(1)$  when  $\rho_x=O(1/K)$  and  $\rho_y=O(1/K)$ ,  $c_{3,2}=O(1/\kappa^2)$  and  $c_{3,3}=O(1/\kappa)$  due to  $c_{\hat{\beta}_x}=O(1/\kappa)$ .

To cancel out  $\frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{y}_t - y_t^{(k)}\|^2]$ , i.e.,

$$\beta_x \eta L^2 + 2L^2 A_3 + \frac{4\eta \hat{\beta}_y^2}{1-\lambda} c_{10} + \frac{108\eta \hat{\beta}_y^2 \gamma_2^2}{(1-\lambda)^2} c_7 + 12\eta^2 \mathcal{Y} - \frac{\eta (1-\lambda^2)}{2} c_4 \le 0.$$
 (121)

Firstly, from the definition of  $\mathcal{Y}$ , we have

$$\mathcal{Y} = \frac{4L^2c_1}{K} + \frac{4L^2c_2}{K} + \frac{9L^2c_6}{1-\lambda} + \frac{9(L^2 + \gamma_2^2)c_7}{1-\lambda} + 4L^2c_8 + 4L^2c_9$$

$$= \frac{4L^2}{K} \frac{\beta_x}{\rho_x \eta} + \frac{10L^2}{K} \frac{\beta_y}{\rho_y \eta} + 21L^2\beta_x \eta + (9L^2 + 21\gamma_2^2)\beta_y \eta . \tag{122}$$

Therefore, we set

$$\beta_{x}\eta L^{2} + 2L^{2}A_{3} + \frac{4\eta\hat{\beta}_{y}^{2}}{1-\lambda}c_{10} + \frac{108\eta\hat{\beta}_{y}^{2}\gamma_{2}^{2}}{(1-\lambda)^{2}}c_{7} + 12\eta^{2}\mathcal{Y}$$

$$\leq \beta_{x}\eta L^{2} + \beta_{x}\eta L^{2}\frac{5}{2}c_{\beta_{y}} + \beta_{x}^{3}\eta c_{\beta_{y}}c_{\hat{\beta}_{x}}^{2}\frac{864\gamma_{2}^{2}}{(1-\lambda)^{2}} + \beta_{x}^{2}\eta c_{\beta_{y}}c_{\hat{\beta}_{x}}\frac{108\gamma_{2}^{2}}{(1-\lambda)}$$

$$+ 12\beta_{x}\eta\left(\frac{4L^{2}}{\rho_{x}K} + \frac{10L^{2}}{\rho_{y}K}c_{\beta_{y}} + 21L^{2} + (21L^{2} + 9\gamma_{2}^{2})c_{\beta_{y}}\right)$$

$$\leq \frac{\eta(1-\lambda)}{2}c_{4}, \qquad (123)$$

where the second step holds due to  $\hat{\beta}_y \eta < 1$  and  $\eta < 1$ , the fourth step holds due to Eq. (54). By solving this inequality, we set

$$c_{4} = \frac{2\beta_{x}}{(1-\lambda)} \left( \underbrace{\frac{48L^{2}}{\rho_{x}K} + \frac{120L^{2}}{\rho_{y}K} c_{\beta_{y}} + 253L^{2} + 255L^{2}c_{\beta_{y}} + 108\gamma_{2}^{2}c_{\beta_{y}}}_{c_{4,1}} \right) + \underbrace{\frac{2\beta_{x}^{3}}{(1-\lambda)^{3}}}_{c_{4,2}} \underbrace{\frac{864\gamma_{2}^{2}c_{\beta_{y}}c_{\hat{\beta}_{x}}^{2}}{(1-\lambda)^{2}}}_{c_{4,3}} \underbrace{\frac{108\gamma_{2}^{2}c_{\beta_{y}}c_{\hat{\beta}_{x}}}{c_{4,3}}}_{(124)}.$$
(124)

Similarly, it is easy to know that  $c_{4,1}=O(1)$  when  $\rho_x=O(1/K)$  and  $\rho_y=O(1/K)$ ,  $c_{4,2}=O(1/\kappa^2)$  and  $c_{4,3}=O(1/\kappa)$  due to  $c_{\hat{\beta}_x}=O(1/\kappa)$ .

To cancel out  $\frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\bar{p}_t - p_t^{(k)}\|^2]$ , i.e.,

$$\frac{2\eta\beta_x^2}{1-\lambda}c_3 + \frac{8\eta^2\beta_x^2\hat{\beta}_x^2}{(1-\lambda)^2}c_5 + \frac{216\eta^2\beta_x^2\hat{\beta}_x^2\gamma_1^2}{(1-\lambda)^3}c_6 + 3\beta_x^2\eta^2\mathcal{X} - (1-\lambda)c_6 \le 0.$$
 (125)

Firstly, we enforce

$$\frac{216\eta^2 \beta_x^2 \hat{\beta}_x^2 \gamma_1^2}{(1-\lambda)^3} c_6 \le \frac{(1-\lambda)}{4} c_6 \ . \tag{126}$$

Then, based on Eq. (54), we obtain

$$\beta_x \le \frac{(1-\lambda)}{6\sqrt{\gamma_1 c_{\hat{\beta}_x}}} \ . \tag{127}$$

Then, we enforce

$$c_{3} \frac{2\eta \beta_{x}^{2}}{1-\lambda} \leq \frac{\beta_{x}\eta}{4} (1-\lambda)^{2} ,$$

$$c_{5} \frac{8\eta^{2} \beta_{x}^{2} \hat{\beta}_{x}^{2}}{(1-\lambda)^{2}} \leq \frac{\beta_{x}\eta}{4} (1-\lambda)^{2} ,$$

$$3\beta_{x}^{2} \eta^{2} \mathcal{X} \leq \frac{\beta_{x}\eta}{16} (1-\lambda)^{2} .$$
(128)

To solve the first inequality in Eq. (128), we enforce

$$\frac{2\eta \beta_x^2}{1-\lambda} \frac{2\beta_x}{(1-\lambda)} c_{3,1} \le \frac{\beta_x \eta}{12} (1-\lambda)^2 ,$$

2052
2053
$$\frac{2\eta\beta_x^2}{1-\lambda} \frac{2\beta_x^3}{(1-\lambda)^3} c_{3,2} \le \frac{\beta_x\eta}{12} (1-\lambda)^2 , \qquad (129)$$
2054
2055
$$\frac{2\eta\beta_x^2}{1-\lambda} \frac{2\beta_x^2}{(1-\lambda)^2} c_{3,3} \le \frac{\beta_x\eta}{12} (1-\lambda)^2 . \qquad (129)$$

Therefore, we obtain

$$\beta_x \le \min \left\{ \frac{(1-\lambda)^2}{4\sqrt{3c_{3,1}}} \,, \frac{(1-\lambda)^{3/2}}{2(3c_{3,2})^{1/4}} \,, \frac{(1-\lambda)^{5/3}}{2(6c_{3,3})^{1/3}} \right\} \,. \tag{130}$$

To solve the second inequality in Eq. (128), we obtain

$$\beta_x \le \frac{(1-\lambda)^{5/4}}{12\sqrt{\gamma_1 c_{\hat{\beta}_x}}} \,. \tag{131}$$

To solve the last inequality in Eq. (128), we enforce

$$3\beta_x \eta \frac{4L^2}{K} \frac{\beta_x}{\rho_x \eta} \le \frac{1}{16 \times 4} (1 - \lambda)^2 , \quad 3\beta_x \eta \frac{10L^2}{K} \frac{\beta_y}{\rho_y \eta} \le \frac{1}{16 \times 4} (1 - \lambda)^2 ,$$
$$3\beta_x \eta (21L^2 + 9\gamma_1^2)\beta_x \eta \le \frac{1}{16 \times 4} (1 - \lambda)^2 , \quad 3\beta_x \eta 21L^2 \beta_y \eta \le \frac{1}{16 \times 4} (1 - \lambda)^2 . \tag{132}$$

We obtain

$$\beta_x \le \left\{ \frac{\sqrt{\rho_x K} (1 - \lambda)}{16\sqrt{3}L}, \frac{\sqrt{\rho_y K} (1 - \lambda)}{8L\sqrt{30c_{\beta_y}}}, \frac{(1 - \lambda)}{8\sqrt{3(21L^2 + 9\gamma_1^2)}}, \frac{(1 - \lambda)}{24L\sqrt{7c_{\beta_y}}} \right\}. \tag{133}$$

To cancel out  $\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{q}_t - q_t^{(k)}\|^2]$ , i.e.,

$$\frac{2\eta\beta_y^2}{1-\lambda}c_4 + \frac{8\eta^2\beta_y^2\hat{\beta}_y^2}{(1-\lambda)^2}c_{10} + \frac{216\eta^2\beta_y^2\hat{\beta}_y^2\gamma_2^2}{(1-\lambda)^3}c_7 + 3\beta_y^2\eta^2\mathcal{Y} - (1-\lambda)c_7 \le 0.$$
 (134)

Firstly, we enforce

$$\frac{216\eta^2 \beta_y^2 \hat{\beta}_y^2 \gamma_2^2}{(1-\lambda)^3} c_7 \le \frac{(1-\lambda)}{6} c_7. \tag{135}$$

Then, based on Eq. (54), we obtain

$$\beta_x \le \frac{(1-\lambda)}{6\sqrt{\gamma_2 c_{\beta_y} c_{\hat{\beta}_y}}} \,. \tag{136}$$

Then, we enforce

$$c_{4} \frac{2\eta \beta_{y}^{2}}{1-\lambda} \leq \frac{\beta_{y} \eta}{4} (1-\lambda)^{2} ,$$

$$c_{10} \frac{8\eta^{2} \beta_{y}^{2} \hat{\beta}_{y}^{2}}{(1-\lambda)^{2}} \leq \frac{\beta_{y} \eta}{4} (1-\lambda)^{2} ,$$

$$3\beta_{y}^{2} \eta^{2} \mathcal{Y} \leq \frac{\beta_{y} \eta}{96} (1-\lambda)^{2} .$$
(137)

To solve the first inequality in Eq. (137), we enforce

$$\frac{2\eta\beta_y^2}{1-\lambda} \frac{2\beta_x}{(1-\lambda)} c_{4,1} \le \frac{\beta_y \eta}{12} (1-\lambda)^2 ,$$

$$\frac{2\eta\beta_y^2}{1-\lambda} \frac{2\beta_x^3}{(1-\lambda)^3} c_{4,2} \le \frac{\beta_y \eta}{12} (1-\lambda)^2 ,$$
(138)

$$\frac{2\eta\beta_y^2}{1-\lambda}\frac{2\beta_x^2}{(1-\lambda)^2}c_{4,3} \leq \frac{\beta_y\eta}{12}(1-\lambda)^2$$

2109 Therefore, we obtain

$$\beta_x \le \min \left\{ \frac{(1-\lambda)^2}{4\sqrt{3c_{\beta_y}c_{4,1}}} , \frac{(1-\lambda)^{3/2}}{2(3c_{\beta_y}c_{4,2})^{1/4}} , \frac{(1-\lambda)^{5/3}}{2(6c_{\beta_y}c_{4,3})^{1/3}} \right\} . \tag{139}$$

To solve the second inequality in Eq. (137), we obtain

$$\beta_x \le \frac{(1-\lambda)^{5/4}}{12\sqrt{\gamma_2 c_{\beta_y} c_{\hat{\beta}_y}}} \,. \tag{140}$$

To solve the last inequality in Eq. (137), we enforce

$$3\beta_{y}\eta \frac{4L^{2}}{K} \frac{\beta_{x}}{\rho_{x}\eta} \leq \frac{1}{96 \times 4} (1 - \lambda)^{2} , \quad 3\beta_{y}\eta \frac{10L^{2}}{K} \frac{\beta_{y}}{\rho_{y}\eta} \leq \frac{1}{96 \times 4} (1 - \lambda)^{2} ,$$

$$3\beta_{y}\eta 21L^{2}\beta_{x}\eta \leq \frac{1}{96 \times 4} (1 - \lambda)^{2} , \quad 3\beta_{y}\eta (21L^{2} + 9\gamma_{2}^{2})\beta_{y}\eta \leq \frac{1}{96 \times 4} (1 - \lambda)^{2} . \tag{141}$$

We obtain

$$\beta_x \le \left\{ \frac{\sqrt{\rho_x K} (1 - \lambda)}{48\sqrt{2c_{\beta_y}} L}, \frac{\sqrt{\rho_y K} (1 - \lambda)}{48c_{\beta_y} L\sqrt{5}}, \frac{(1 - \lambda)}{24L\sqrt{42c_{\beta_y}}}, \frac{(1 - \lambda)}{24c_{\beta_y} \sqrt{2(21L^2 + 9\gamma_2^2)}} \right\}. \tag{142}$$

For  $\mathbb{E}[\|\bar{\hat{x}}_{t+1} - \bar{\hat{x}}_t\|^2]$ , by setting

$$2\gamma_1 C_{x_{\hat{x}\hat{y}}^1} + 6\gamma_1 \hat{\beta}_x \eta \left( 10 C_{x_{\hat{y}\hat{x}\hat{y}}^1}^2 C_{y_{\hat{x}\hat{y}}^1}^2 + \frac{4}{\gamma_1 - L} \frac{2\gamma_2^2 C_{y_{\hat{x}\hat{y}}^1}^2}{\mu} \right) + \frac{27\gamma_1^2}{1 - \lambda} c_6 - \frac{\gamma_1}{3\hat{\beta}_x \eta} \le -\frac{\gamma_1}{4\hat{\beta}_x \eta} . \quad (143)$$

Specifically, we enforce

$$2\gamma_{1}C_{x_{\hat{x}\hat{y}}^{1}} \leq \frac{\gamma_{1}}{36\hat{\beta}_{x}\eta} ,$$

$$6\gamma_{1}\hat{\beta}_{x}\eta \left(10C_{x_{\hat{y}\hat{x}\hat{y}}^{1}}^{2}C_{y_{\hat{x}\hat{y}}^{1}}^{2} + \frac{4}{\gamma_{1} - L} \frac{2\gamma_{2}^{2}C_{y_{\hat{x}\hat{y}}^{1}}^{2}}{\mu}\right) \leq \frac{\gamma_{1}}{36\hat{\beta}_{x}\eta} ,$$

$$\frac{27\gamma_{1}^{2}}{1 - \lambda}c_{6} \leq \frac{\gamma_{1}}{36\hat{\beta}_{x}\eta} .$$
(144)

Since  $C_{x^1_{\hat{x}\hat{y}}}=rac{\gamma_1}{\gamma_1-L}$ ,  $C_{x^1_{y\hat{x}\hat{y}}}=rac{\gamma_1}{\gamma_1-L}$ , and  $C_{y^1_{\hat{x}\hat{y}}}=rac{\gamma_1}{\gamma_2-L}$ , we obtain

$$\beta_x \le \min \left\{ \frac{\gamma_1 - L}{72c_{\hat{\beta}_x}\gamma_1}, \frac{(\gamma_1 - L)(\gamma_2 - L)\sqrt{\mu}}{6\gamma_1c_{\hat{\beta}_x}\sqrt{6(10\gamma_1^2\mu + 8\gamma_2^2(\gamma_1 - L))}}, \frac{1}{18\sqrt{3c_{\hat{\beta}_x}\gamma_1}} \right\}.$$
(145)

For  $\mathbb{E}[\|\bar{\hat{y}}_{t+1} - \bar{\hat{y}}_t\|^2]$ , by setting

$$2A_1 + \frac{27\gamma_2^2}{1-\lambda}c_7 - \frac{\gamma_2}{2\hat{\beta}_n\eta} \le -\frac{\gamma_2}{4\hat{\beta}_n\eta} \tag{146}$$

Specifically, from the definition of  $A_1$ , we enforce

$$12\gamma_{1}\hat{\beta}_{x}\eta C_{y_{\hat{x}\hat{y}}^{2}}^{2}\left(10C_{x_{y\hat{x}\hat{y}}^{2}}^{2} + \frac{8\gamma_{2}^{2}}{(\gamma_{1} - L)\mu}\right) \leq \frac{\gamma_{2}}{16\hat{\beta}_{y}\eta},$$

$$12\gamma_{1}\hat{\beta}_{x}\eta \frac{8\gamma_{2}^{2}}{(\gamma_{1} - L)\mu} \frac{(1 - \hat{\beta}_{y}\eta)^{2}}{\hat{\beta}_{y}^{2}\eta^{2}} \leq \frac{\gamma_{2}}{16\hat{\beta}_{y}\eta},$$

$$\frac{27\gamma_{2}^{2}}{1 - \lambda}c_{7} \leq \frac{\gamma_{2}}{8\hat{\beta}_{y}\eta}.$$
(147)

To solve the second inequality, we use the second inequality in Eq. (68) to obtain the following:

$$\frac{4}{(\gamma_2 - L)^2} \frac{24\gamma_1 \hat{\beta}_x \eta}{\gamma_1 - L} \frac{2\gamma_2^2}{\mu} \frac{(1 - \hat{\beta}_y \eta)^2}{\hat{\beta}_y^2 \eta^2} \le \frac{\beta_x \eta}{32 \times 64} \frac{(\gamma_1 - L)^2}{4L^2} . \tag{148}$$

Then, it is easy to derive

$$12\gamma_1 \hat{\beta}_x \eta \frac{8\gamma_2^2}{(\gamma_1 - L)\mu} \frac{(1 - \hat{\beta}_y \eta)^2}{\hat{\beta}_y^2 \eta^2} \le \frac{\beta_x \eta}{32 \times 64} \frac{(\gamma_1 - L)^2}{4L^2} \frac{(\gamma_2 - L)^2}{2}$$
(149)

Therefore, it leads us to solve

$$\frac{\beta_x \eta}{32 \times 64} \frac{(\gamma_1 - L)^2}{4L^2} \frac{(\gamma_2 - L)^2}{2} \le \frac{\gamma_2}{16\hat{\beta}_u \eta}$$
 (150)

and we obtain

$$\beta_x \le \frac{32L}{\sqrt{c_{\hat{\beta}_n}}(\gamma_1 - L)(\gamma_2 - L)} \tag{151}$$

Finally, to solve the first and last inequality in Eq. (147), from  $C_{x_{y\hat{x}\hat{y}}^1} = \frac{\gamma_1}{\gamma_1 - L}$  and  $C_{y_{\hat{x}\hat{y}}^2} = \frac{\gamma_2}{\gamma_2 - L}$ , we obtain

$$\beta_x \le \min \left\{ \frac{\sqrt{\mu}(\gamma_1 - L)(\gamma_2 - L)}{8\sqrt{3c_{\hat{\beta}_x}c_{\hat{\beta}_y}\gamma_1\gamma_2(10\gamma_1^2\mu + 8\gamma_2^2(\gamma_1 - L))}}, \frac{1}{6\sqrt{6\gamma_2c_{\beta_y}c_{\hat{\beta}_y}}} \right\}$$
(152)

By setting

$$c_{1} = \frac{\beta_{x}}{\rho_{x}\eta}, \quad c_{2} = \frac{5\beta_{y}}{2\rho_{y}\eta},$$

$$c_{3} \triangleq \frac{2\beta_{x}}{(1-\lambda)}c_{3,1} + \frac{2\beta_{x}^{3}}{(1-\lambda)^{3}}c_{3,2} + \frac{2\beta_{x}^{2}}{(1-\lambda)^{2}}c_{3,3},$$

$$c_{4} \triangleq \frac{2\beta_{x}}{(1-\lambda)}c_{4,1} + \frac{2\beta_{x}^{3}}{(1-\lambda)^{3}}c_{4,2} + \frac{2\beta_{x}^{2}}{(1-\lambda)^{2}}c_{4,3},$$

$$c_{5} = \frac{216\beta_{x}\gamma_{1}^{2}}{(1-\lambda)}, \quad c_{10} = \frac{216\beta_{y}\gamma_{2}^{2}}{(1-\lambda)}$$

$$c_{6} = \beta_{x}\eta(1-\lambda), \quad c_{7} = \beta_{y}\eta(1-\lambda), \quad c_{8} = 3\beta_{x}\eta, \quad c_{9} = 3\beta_{y}\eta, \quad (153)$$

we obtain

2198
2199
$$\mathcal{L}_{t+1} - \mathcal{L}_{t} \leq -\frac{\beta_{x}\eta}{4} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t})\|^{2}] - \frac{\beta_{y}\eta}{2} \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t},\bar{y}_{t};\bar{\hat{x}}_{t},\bar{\hat{y}}_{t})\|^{2}]$$
2200
$$+ \left(4\beta_{y}\eta\beta_{x}^{2}\eta^{2}L^{2} + 3\beta_{x}^{2}\eta^{2}\mathcal{X} - \frac{\beta_{x}\eta}{4}\right) \mathbb{E}[\|\bar{u}_{t}\|^{2}]$$
2201
$$+ \left(\beta_{y}^{2}\eta^{2}L_{d} + \frac{3\beta_{y}\eta}{4} + \frac{\beta_{y}^{2}\eta^{2}(\gamma_{2} + L)}{2} + 4A_{1}\hat{\beta}_{y}^{2}\eta^{2}\beta_{y}^{2}\eta^{2} + 2A_{2}\beta_{y}^{2}\eta^{2} + 3\beta_{y}^{2}\eta^{2}\mathcal{Y} - \frac{7}{8}\beta_{y}\eta\right) \mathbb{E}[\|\bar{v}_{t}\|^{2}]$$
2205
$$- \frac{\gamma_{1}}{4\hat{\beta}_{x}\eta} \mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] - \frac{\gamma_{2}}{4\hat{\beta}_{y}\eta} \mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}]$$
2206
$$+ 2c_{1}\rho_{x}^{2}\eta^{4}\sigma^{2}\frac{1}{K} + 2c_{2}\rho_{y}^{2}\eta^{4}\sigma^{2}\frac{1}{K} + 3c_{6}\rho_{x}^{2}\eta^{4}\sigma^{2}\frac{1}{1 - \lambda} + 3c_{7}\rho_{y}^{2}\eta^{4}\sigma^{2}\frac{1}{1 - \lambda} + 2c_{8}\rho_{x}^{2}\eta^{4}\sigma^{2} + 2c_{9}\rho_{y}^{2}\eta^{4}\sigma^{2}.$$
2208
2209
$$+ 2c_{1}\rho_{x}^{2}\eta^{4}\sigma^{2}\frac{1}{K} + 2c_{2}\rho_{y}^{2}\eta^{4}\sigma^{2}\frac{1}{K} + 3c_{6}\rho_{x}^{2}\eta^{4}\sigma^{2}\frac{1}{1 - \lambda} + 3c_{7}\rho_{y}^{2}\eta^{4}\sigma^{2}\frac{1}{1 - \lambda} + 2c_{8}\rho_{x}^{2}\eta^{4}\sigma^{2} + 2c_{9}\rho_{y}^{2}\eta^{4}\sigma^{2}.$$
(154)

For  $\mathbb{E}[\|\bar{u}_t\|^2]$ , we enforce

$$4\beta_y \eta \beta_x^2 \eta^2 L^2 + 3\beta_x^2 \eta^2 \mathcal{X} - \frac{\beta_x \eta}{4} \le -\frac{\beta_x \eta}{8} . \tag{155}$$

Specifically, we enforce

$$4\beta_y \eta \beta_x^2 \eta^2 L^2 \le \frac{\beta_x \eta}{16}$$
$$3\beta_x^2 \eta^2 \mathcal{X} \le \frac{\beta_x \eta}{16} . \tag{156}$$

To solve the first inequality, we obtain

$$\beta_x \le \frac{1}{8L\sqrt{c_{\beta_n}}} \tag{157}$$

To solve the last inequality, we use the last inequality in Eq. (128) along with the fact that  $1 - \lambda < 1$ , from which it is straightforward to show that the inequality holds.

For  $\mathbb{E}[\|\bar{v}_t\|^2]$ , we enforce

$$\beta_y^2 \eta^2 L_d + \frac{3\beta_y \eta}{4} + \frac{\beta_y^2 \eta^2 (\gamma_2 + L)}{2} + 4A_1 \hat{\beta}_y^2 \eta^2 \beta_y^2 \eta^2 + 2A_2 \beta_y^2 \eta^2 + 3\beta_y^2 \eta^2 \mathcal{Y} - \frac{7}{8} \beta_y \eta \le -\frac{1}{32} \beta_y \eta. \tag{158}$$

Firstly, from Eq. (105) and the definition of  $A_2$ , we obtain  $2A_2\beta_y^2\eta^2 \leq \frac{\beta_y\eta}{4\times 4}$ , and

$$2\beta_y^2 \eta^2 A_1 \frac{4\hat{\beta}_y^2}{\beta_y^2 (\gamma_2 - L)^2} \le \frac{\beta_y \eta}{4 \times 4}$$
(159)

By reformulating the above inequality, we obtain

$$4A_1\hat{\beta}_y^2\eta^2\beta_y^2\eta^2 \le \frac{\beta_y^3\eta^3(\gamma_2 - L)^2}{4 \times 8} \tag{160}$$

Therefore, we enforce

$$\beta_y^2 \eta^2 L_d + \frac{\beta_y^2 \eta^2 (\gamma_2 + L)}{2} + \frac{\beta_y^3 \eta^3 (\gamma_2 - L)^2}{32} + 3\beta_y^2 \eta^2 \mathcal{Y} \le \frac{1}{32} \beta_y \eta \tag{161}$$

Specifically, we enforce

$$\beta_y^2 \eta^2 L_d + \frac{\beta_y^2 \eta^2 (\gamma_2 + L)}{2} \le \frac{1}{96} \beta_y \eta ,$$

$$\frac{\beta_y^3 \eta^3 (\gamma_2 - L)^2}{32} \le \frac{1}{96} \beta_y \eta ,$$

$$3\beta_y^2 \eta^2 \mathcal{Y} \le \frac{1}{96} \beta_y \eta ,$$
(162)

To solve the first and second inequality, we obtain

$$\beta_x \le \min\left\{\frac{1}{48c_{\beta_y}(2L_d + \gamma_2 + L)}, \frac{1}{\sqrt{3}c_{\beta_y}(\gamma_2 - L)}\right\}$$
 (163)

To solve the last inequality, we use the last inequality in Eq. (137) along with the fact that  $1 - \lambda < 1$ , from which it is straightforward to show that the inequality holds.

In summary, by setting

$$\beta_x \leq \min \left\{ \frac{(1-\lambda)}{6\sqrt{\gamma_1 c_{\hat{\beta}_x}}}, \frac{(1-\lambda)^2}{4\sqrt{3c_{3,1}}}, \frac{(1-\lambda)^{3/2}}{2(3c_{3,2})^{1/4}}, \frac{(1-\lambda)^{5/3}}{2(6c_{3,3})^{1/3}}, \frac{(1-\lambda)^{5/4}}{12\sqrt{\gamma_1 c_{\hat{\beta}_x}}}, \frac{\sqrt{\rho_x K}(1-\lambda)}{16\sqrt{3}L}, \frac{\sqrt{\rho_y K}(1-\lambda)}{8L\sqrt{30c_{\beta_y}}}, \frac{(1-\lambda)}{8\sqrt{3(21L^2+9\gamma_1^2)}}, \frac{(1-\lambda)}{6\sqrt{\gamma_2 c_{\beta_y} c_{\hat{\beta}_y}}}, \frac{(1-\lambda)^2}{4\sqrt{3c_{\beta_y} c_{4,1}}}, \frac{(1-\lambda)^{3/2}}{2(3c_{\beta_y} c_{4,2})^{1/4}}, \frac{(1-\lambda)^2}{2(3c_{\beta_y} c_{$$

$$\begin{split} &\frac{(1-\lambda)^{5/3}}{2(6c_{\beta_y}c_{4,3})^{1/3}}, \frac{(1-\lambda)^{5/4}}{12\sqrt{\gamma_2c_{\beta_y}}c_{\hat{\beta}_y}}, \frac{\sqrt{\rho_xK}(1-\lambda)}{48\sqrt{2c_{\beta_y}}L}, \frac{\sqrt{\rho_yK}(1-\lambda)}{48c_{\beta_y}L\sqrt{5}}, \frac{(1-\lambda)}{24L\sqrt{42c_{\beta_y}}}, \\ &\frac{(1-\lambda)}{24c_{\beta_y}\sqrt{2(21L^2+9\gamma_2^2)}}, \frac{\gamma_1-L}{72c_{\hat{\beta}_x}\gamma_1}, \frac{(\gamma_1-L)(\gamma_2-L)\sqrt{\mu}}{6\gamma_1c_{\hat{\beta}_x}\sqrt{6(10\gamma_1^2\mu+8\gamma_2^2(\gamma_1-L))}}, \frac{1}{18\sqrt{3c_{\hat{\beta}_x}\gamma_1}}, \\ &\frac{4L}{\sqrt{c_{\hat{\beta}_y}}(\gamma_1-L)(\gamma_2-L)}, \frac{\sqrt{\mu}(\gamma_1-L)(\gamma_2-L)}{8\sqrt{3c_{\hat{\beta}_x}c_{\hat{\beta}_y}\gamma_1\gamma_2(10\gamma_1^2\mu+8\gamma_2^2(\gamma_1-L))}}, \frac{1}{6\sqrt{6\gamma_2c_{\beta_y}c_{\hat{\beta}_y}}}, \\ &\frac{L^2}{120\gamma_1^3}, \frac{\sqrt{\mu}(\gamma_1-L)^3(\gamma_2-L)^2}{512\sqrt{6\gamma_1c_{\hat{\beta}_x}\gamma_2c_{\hat{\beta}_y}}}, \frac{1}{8L\sqrt{c_{\beta_y}}}, \frac{1}{48c_{\beta_y}(2L_d+\gamma_2+L)}, \frac{1}{\sqrt{3}c_{\beta_y}(\gamma_2-L)} \Big\} \\ &\eta \leq \min \left\{ \frac{1}{\sqrt{\rho_x}}, \frac{1}{\hat{\beta}_x}, \frac{1}{\hat{\beta}_x}, \frac{1}{\hat{\beta}_x}, \frac{1}{2\beta_x(\gamma_1+L)} \right\}, \end{split}$$

we obtain

$$\mathcal{L}_{t+1} - \mathcal{L}_{t} \leq -\frac{\beta_{x}\eta}{4} \mathbb{E}[\|\nabla_{x}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] - \frac{\beta_{x}c_{\beta_{y}}\eta}{2} \mathbb{E}[\|\nabla_{y}F(\bar{x}_{t}, \bar{y}_{t}; \bar{x}_{t}, \bar{y}_{t})\|^{2}] 
- \gamma_{1}c_{\hat{\beta}_{x}}\frac{\beta_{x}\eta}{4} \mathbb{E}[\|\bar{x}_{t+1} - \bar{x}_{t}\|^{2}] - \gamma_{2}c_{\hat{\beta}_{y}}\frac{\beta_{x}\eta}{4} \mathbb{E}[\|\bar{y}_{t+1} - \bar{y}_{t}\|^{2}] 
+ 2\beta_{x}\rho_{x}\eta^{3}\sigma^{2}\frac{1}{K} + 5\beta_{x}c_{\beta_{y}}\rho_{y}\eta^{3}\sigma^{2}\frac{1}{K} + 9\beta_{x}\rho_{x}^{2}\eta^{5}\sigma^{2} + 9c_{\beta_{y}}\beta_{x}\rho_{y}^{2}\eta^{5}\sigma^{2} .$$
(165)

Because

$$\|\nabla_{x} f(\bar{x}_{t}, \bar{y}_{t})\|^{2} \leq 2\|\nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \bar{\hat{x}}_{t}, \bar{y}_{t})\|^{2} + 2\gamma_{1}^{2}\|\bar{x}_{t+1} - \bar{\hat{x}}_{t}\|^{2},$$
  
$$\|\nabla_{y} f(\bar{x}_{t}, \bar{y}_{t})\|^{2} \leq 2\|\nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \bar{\hat{x}}_{t}, \bar{y}_{t})\|^{2} + 2\gamma_{2}^{2}\|\bar{y}_{t+1} - \bar{\hat{y}}_{t}\|^{2},$$
(166)

we obtain

$$\frac{1}{T} \sum_{t=0}^{T-1} \left( \mathbb{E}[\|\nabla_{x} f(\bar{x}_{t}, \bar{y}_{t})\|^{2}] + \kappa \mathbb{E}[\|\nabla_{y} f(\bar{x}_{t}, \bar{y}_{t})\|^{2}] \right) 
\leq \frac{1}{T} \sum_{t=0}^{T-1} \left( 2\mathbb{E}[\|\nabla_{x} F(\bar{x}_{t}, \bar{y}_{t}; \hat{\bar{x}}_{t}, \bar{\hat{y}}_{t})\|^{2}] + 2\kappa \mathbb{E}[\|\nabla_{y} F(\bar{x}_{t}, \bar{y}_{t}; \hat{\bar{x}}_{t}, \bar{\hat{y}}_{t})\|^{2}] + 2\gamma_{1}^{2} \mathbb{E}[\|\bar{x}_{t+1} - \bar{\hat{x}}_{t}\|^{2}] \right) 
+ 2\kappa \gamma_{2}^{2} \mathbb{E}[\|\bar{y}_{t+1} - \hat{\bar{y}}_{t}\|^{2}] \right) 
\leq \max \left\{ \frac{8}{\beta_{x}\eta}, \frac{8\kappa}{\beta_{x}\eta c_{\beta_{y}}}, \frac{8\gamma_{1}}{\beta_{x}\eta c_{\hat{\beta}_{x}}}, \frac{8\kappa\gamma_{2}}{\beta_{x}\eta c_{\hat{\beta}_{y}}} \right\} \left( \frac{\mathcal{L}_{0} - \mathcal{L}_{T}}{T} + 2\beta_{x}\rho_{x}\eta^{3}\sigma^{2} \frac{1}{K} + 5\beta_{x}c_{\beta_{y}}\rho_{y}\eta^{3}\sigma^{2} \frac{1}{K} + 9\beta_{x}\rho_{x}^{2}\eta^{5}\sigma^{2} + 9c_{\beta_{y}}\beta_{x}\rho_{y}^{2}\eta^{5}\sigma^{2} \right).$$
(167)

By setting  $\gamma_1 = O(L)$ ,  $\gamma_2 = O(L)$ , we obtain

$$c_{\beta_y} = O(1) , \quad c_{\hat{\beta}_x} = O\left(\frac{1}{L^2 \kappa}\right) , \quad c_{\hat{\beta}_y} = O(1) .$$
 (168)

Because

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{p}_{0} - p_{0}^{(k)}\|^{2}]$$

$$= \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\frac{1}{K} \sum_{j=1}^{K} \nabla_{x} F^{(j)}(x_{0}, y_{0}; \hat{x}_{0}, \hat{y}_{0}; \xi_{0}^{(j)}) - \nabla_{x} F^{(k)}(x_{0}, y_{0}; \hat{x}_{0}, \hat{y}_{0}; \xi_{0}^{(k)})\|^{2}]$$

$$\leq 3 \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\frac{1}{K} \sum_{j=1}^{K} \nabla_{x} F^{(j)}(x_{0}, y_{0}; \hat{x}_{0}, \hat{y}_{0}; \xi_{0}^{(j)}) - \frac{1}{K} \sum_{j=1}^{K} \nabla_{x} F^{(j)}(x_{0}, y_{0}; \hat{x}_{0}, \hat{y}_{0})\|^{2}]$$

$$+3\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}\left[\left\|\frac{1}{K}\sum_{j=1}^{K}\nabla_{x}F^{(j)}(x_{0},y_{0};\hat{x}_{0},\hat{y}_{0}) - \nabla_{x}F^{(k)}(x_{0},y_{0};\hat{x}_{0},\hat{y}_{0})\right\|^{2}\right]$$

$$+3\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}\left[\left\|\nabla_{x}F^{(k)}(x_{0},y_{0};\hat{x}_{0},\hat{y}_{0}) - \nabla_{x}F^{(k)}(x_{0},y_{0};\hat{x}_{0},\hat{y}_{0};\xi_{0}^{(k)})\right\|^{2}\right]$$

$$\leq 6\sigma^{2} + 6\frac{1}{K}\sum_{k=1}^{K}\mathbb{E}\left[\left\|\nabla_{x}f^{(k)}(x_{0},y_{0})\right\|^{2}\right],$$
(169)

and

$$\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{q}_0 - q_0^{(k)}\|^2] \le 6\sigma^2 + 6\frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\nabla_y f^{(k)}(x_0, y_0)\|^2],$$
 (170)

we have

$$\mathcal{L}_{0} = \mathcal{P}_{0} + \frac{\beta_{x}}{\rho_{x}\eta} \mathbb{E}[\|\frac{1}{K} \sum_{k=1}^{K} u_{0}^{(k)} - \frac{1}{K} \sum_{k=1}^{K} \nabla_{x} F^{(k)}(x_{0}^{(k)}, y_{0}^{(k)}; \hat{x}_{0}^{(k)}, \hat{y}_{0}^{(k)})\|^{2}] \\
+ \frac{5\beta_{y}}{2\rho_{y}\eta} \mathbb{E}[\|\frac{1}{K} \sum_{k=1}^{K} v_{0}^{(k)} - \frac{1}{K} \sum_{k=1}^{K} \nabla_{y} F^{(k)}(x_{0}^{(k)}, y_{0}^{(k)}; \hat{x}_{0}^{(k)}, \hat{y}_{0}^{(k)})\|^{2}] \\
+ \beta_{x}\eta(1 - \lambda) \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{p}_{0} - p_{0}^{(k)}\|^{2}] + \beta_{y}\eta(1 - \lambda) \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\bar{q}_{0} - q_{0}^{(k)}\|^{2}] \\
+ 3\beta_{x}\eta \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|u_{0}^{(k)} - \nabla_{x} F^{(k)}(x_{0}^{(k)}, y_{0}^{(k)}; \hat{x}_{0}^{(k)}, \hat{y}_{0}^{(k)})\|^{2}] \\
+ 3\beta_{y}\eta \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|v_{0}^{(k)} - \nabla_{y} F^{(k)}(x_{0}^{(k)}, y_{0}^{(k)}; \hat{x}_{0}^{(k)}, \hat{y}_{0}^{(k)})\|^{2}] \\
\leq \mathcal{P}_{0} + \frac{\beta_{x}}{\rho_{x}\eta} \frac{\sigma^{2}}{B} + \frac{5\beta_{y}}{2\rho_{y}\eta} \frac{\sigma^{2}}{B} + 9\beta_{x}\eta\sigma^{2} + 9\beta_{y}\eta\sigma^{2} + 6\beta_{x}\eta \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\nabla_{x} f^{(k)}(x_{0}, y_{0})\|^{2}] \\
+ 6\beta_{y}\eta \frac{1}{K} \sum_{k=1}^{K} \mathbb{E}[\|\nabla_{y} f^{(k)}(x_{0}, y_{0})\|^{2}]. \tag{171}$$

Then, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \left( \mathbb{E}[\|\nabla_x f(\bar{x}_t, \bar{y}_t)\|^2] + \kappa \mathbb{E}[\|\nabla_y f(\bar{x}_t, \bar{y}_t)\|^2] \right) \\
\leq O\left(\frac{\kappa \mathcal{P}_0}{\beta_x \eta T}\right) + O\left(\frac{\kappa}{T} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\nabla_x f^{(k)}(x_0, y_0)\|^2]\right) + O\left(\frac{\kappa}{T} \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\|\nabla_y f^{(k)}(x_0, y_0)\|^2]\right) \\
+ O\left(\frac{\kappa \sigma^2}{\rho_x \eta^2 T B}\right) + O\left(\frac{\kappa \sigma^2}{\rho_y \eta^2 T B}\right) + O\left(\frac{\kappa \sigma^2}{T}\right) + O\left(\frac{\kappa \rho_x \eta^2 \sigma^2}{K}\right) + O\left(\frac{\kappa \rho_y \eta^2 \sigma^2}{K}\right) \\
+ O(\kappa \rho_x^2 \eta^4 \sigma^2) + O(\kappa \rho_y^2 \eta^4 \sigma^2). \tag{172}$$

By setting 
$$\beta_x = O((1-\lambda)^2)$$
,  $\eta = O(\frac{K\epsilon}{\kappa^{1/2}})$ ,  $\rho_x = O(\frac{1}{K})$ ,  $\rho_y = O(\frac{1}{K})$ ,  $B = O(\frac{\kappa^{1/2}}{\epsilon})$ ,  $T = O(\frac{\kappa^{3/2}}{K(1-\lambda)^2\epsilon^3})$ , we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \left( \mathbb{E}[\|\nabla_x f(\bar{x}_t, \bar{y}_t)\|^2] + \kappa \mathbb{E}[\|\nabla_y f(\bar{x}_t, \bar{y}_t)\|^2] \right) \le O(\epsilon^2) . \tag{173}$$