
A Stochastic Linearized Augmented Lagrangian Method for Decentralized Bilevel Optimization

Songtao Lu[†] Siliang Zeng[‡] Xiaodong Cui[†] Mark S. Squillante[†]
Lior Horesh[†] Brian Kingsbury[†] Jia Liu^{*} Mingyi Hong[‡]

[†]IBM Research, Thomas J. Watson Research Center, Yorktown Heights, NY 10598
songtao@ibm.com, {cuix, mss, lhoresh, bedk}@us.ibm.com

[‡]Dept. of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455
{zeng0176, mhong}@umn.edu

^{*}Dept. of Electrical and Computer Engineering, The Ohio State University, Columbus, OH 43210
liu@ece.osu.edu

Abstract

Bilevel optimization has been shown to be a powerful framework for formulating multi-task machine learning problems, e.g., reinforcement learning (RL) and meta-learning, where the decision variables are coupled in both levels of the minimization problems. In practice, the learning tasks would be located at different computing resource environments, and thus there is a need for deploying a decentralized training framework to implement multi-agent and multi-task learning. We develop a stochastic linearized augmented Lagrangian method (SLAM) for solving general nonconvex bilevel optimization problems over a graph, where both upper and lower optimization variables are able to achieve a consensus. We also establish that the theoretical convergence rate of the proposed SLAM to the Karush-Kuhn-Tucker (KKT) points of this class of problems is on the same order as the one achieved by the classical distributed stochastic gradient descent for only single-level nonconvex minimization problems. Numerical results tested on multi-agent RL problems showcase the superiority of SLAM compared with the benchmarks.

1 Introduction

In this paper, we consider the following general decentralized bilevel optimization (DBO) framework with applications to machine learning problems. Suppose that there are n nodes over a connected graph $\mathcal{G} = \{\mathcal{E}, \mathcal{V}\}$, where \mathcal{E} and \mathcal{V} represent the edges and vertices. Let \mathcal{N}_i denote the set of neighboring nodes for node i . Then the goal of DBO is to have these nodes jointly minimize two levels of optimization problems. More formally, DBO is expressed as

$$\min_{\mathbf{x}_1, \dots, \mathbf{x}_n} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}_i, \mathbf{y}_{i,1}^*(\mathbf{x}_i), \dots, \mathbf{y}_{i,m}^*(\mathbf{x}_i)) \quad (1a)$$

$$\text{s.t. } \mathbf{x}_i = \mathbf{x}_j, j \in \mathcal{N}_i, \forall i \in [n] \quad (1b)$$

$$\mathbf{y}_k^*(\mathbf{x}) = \arg \min_{\mathbf{y}_{1,k}, \dots, \mathbf{y}_{n,k}} \frac{1}{n} \sum_{i=1}^n g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) \quad \text{s.t. } \mathbf{y}_{i,k} = \mathbf{y}_{j,k}, j \in \mathcal{N}_i, \forall k \in [m], \quad (1c)$$

where vector \mathbf{x}_i is the upper level (UL) optimization variable at each node i , vector $\mathbf{y}_{i,k}$ denotes the lower level (LL) decision variable for the k th learning task at node i , $f_i(\cdot)$ is a (smooth) UL loss function and possibly nonconvex with respect to (w.r.t.) both the UL and LL variables, $g_{i,k}(\cdot)$ denotes the LL objective function of the k th task at node i , m represents the total number of LL

optimization problems, the consensus constraints $\mathbf{x}_i = \mathbf{x}_j, \mathbf{y}_{i,k} = \mathbf{y}_{j,k}, j \in \mathcal{N}_i, \forall i \in [n], \forall k \in [m]$, enforce the model agreements at each level of the problems and for each LL learning task, and $\mathbf{y}_k^* = [\mathbf{y}_{1,k}^*, \dots, \mathbf{y}_{n,k}^*]^T$ is the optimal solutions of the k th LL problem under the consensus constraints.

Applications of Bilevel Optimization. Many machine learning problems can be formulated mathematically as a form of bilevel optimization or, more precisely, a special case of problem (1), e.g., meta-learning or meta reinforcement learning (RL), actor-critic (AC) schemes in RL, hyperparameter optimization (HPO), and so on.

Classical bilevel optimization is referred to as the case where there is no consensus constraint but with only two levels of the minimization subproblems, i.e., $\min_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}^*(\mathbf{x}))$, s.t. $\mathbf{y}^*(\mathbf{x}) = \arg \min_{\mathbf{y}} g(\mathbf{x}, \mathbf{y})$, which is also known as Stackelberg games [1] with the UL decision variable as the leader and the LL decision variable as the follower. It turns out that this class of optimization problems is useful in formulating a wide range of hierarchical or nested structured machine learning problems. For example, one of the most popular domain adaption learning models, model-agnostic meta-learning (MAML) [2, 3], can be written as a special case of bilevel programming [4], where the UL model provides a good initialization for accelerating learning procedures by implementing the LL algorithms. The idea behind the model design is that the UL model is considered as the meta learner that searches for a permutation-invariant subspace over multiple task-specific learners at the LL so that the performance of the MAML model can be generalized well for unseen or testing data samples. The theoretical analysis of the generalization performance of this class of bilevel problems has shown that MAML can indeed decrease the generalization error as the number of tasks increases, at least for strongly convex loss functions [5]. Subsequently, a thorough ablation study from the latent representation perspective shows that feature reuse is the actual dominant factor in improving the generalization performance of MAML [6], and the authors propose a neural network-oriented algorithm with almost no inner loop (ANIL) that splits the neural network parameters into two parts corresponding to the UL and LL optimization problems, respectively. Extensive numerical experiments illustrate that ANIL achieves almost the same accuracy as the classical MAML but with significant computational savings. This example further strengthens the necessity of variable splitting in the learning structure by optimizing two levels of objective functions to enhance the generalization performance. Beyond the traditional supervised meta-learning scenarios, MAML has also been applied to increasing the generalization ability of agents in RL problems by replacing the (stochastic) gradient with the (natural) policy gradient (PG) [3] under the same two-level structure.

Besides meta-learning problems, AC structure in RL is another class of common learning frameworks that can be formulated by a bilevel optimization problem in nature [7, 8, 9], where the actor step at the UL aims at optimizing the policy while the critic step at the LL is responsible for value function evaluation. In addition, as the expressiveness of neural networks increased sharply over the past decades, the reuse of large models with adaptation to multi-task learning problems presents promising solutions by leveraging the pre-train and fine-tune strategy, such as in applications of HPO [10, 11] where the hyperparameters are trained at the UL problem so that the downstream learning tasks are learned with low costs including the expense of both computation and memory.

Applications of Multi-agent Settings. When multiple computational resources are available and connected, it is well motivated that exploring them solves distributed large-scale problems with a reduced amount of training time or performs multi-task learning. The bilevel structure of the meta-learning (ML) is a good fit in this scenario as either UL/LL or both levels may need to access the networked data samples rather than local ones. For example, a federated learning setting of MAML [12] and bilevel optimization [13] have been built up over multiple nodes recently, where the meta/UL learner finds an initial shared model while the local/LL learners leverage it for adapting data distributions of individual users. In such a way, the federated MAML model can realize personalized learning without sharing heterogeneous data over numerous clients. Once there is no central controller for coordinating the model aggregation, a diffusion-based MAML (Dif-MAML) [14] is proposed by spreading the model parameters over a network, where the UL parameter is updated by one step of stochastic gradient descent (SGD) based on a combination of the parameters of neighbors as the initialization for local model updates.

Decentralized hierarchical structured learning is even more stringent in the multi-agent RL (MARL) setting [15] as the learning tasks are essentially located at scattered sensors and/or controllers. Under this setting, MARL problem becomes a multi-objective optimization problem under provided (approximate) value functions, where the policy of each agent needs to be learned locally by certain

efficient iterative methods, such as multi-agent deep deterministic policy gradient (MADDPG) [16], trust region methods [17], optimal baseline based variance reduced policy gradient [18], and/or improved by more advanced techniques, e.g., constrained policy optimization [19] and large sequence models [20]. In such a way, the total reward can be maximized over the distributed agents through optimizing the networked policy. In a fully collaborative setting, the team-based value function is even required to be shared over all the agents such that each agent is able to improve its policy based on the estimated total reward. For example, the decentralized AC (DAC) scheme has been investigated widely [15, 21, 22], where each agent uses the actor step to optimize its policy while the critic step performs one step [23] or multiple steps of temporal difference learning with mini-batch sampling (MDAC) [22, 24] and communications so that the team-based reward over the network is obtained by each agent. It turns out that DAC can be formulated as a special case of problem (1) as there is no consensus at the UL. Recently, it has been revealed that if there exists homogeneity of the state and action spaces, decentralized policy consensus (or a partial policy parameter sharing strategy) provides significant merits to the centralized training and decentralized execution paradigm in terms of learning scalability and efficiency [23, 25], which motivates the consensus process at both UL and LL DBO problems.

Related Theoretical Works. Given the fruitful results across these many applications, the corresponding theoretical analysis has been developing very fast as well for variants of bilevel optimization problems. For example, the convergence behaviors of classical inexact MAML (iMAML) methods have been quantified for both convex [26, 27] and nonconvex [28] cases of the UL loss function, where the LL algorithm only performs one step of stochastic gradient descent (SGD) based on the LL objective functions as the adaptation step. Moreover, the iteration complexity of ANIL with multiple iterations for minimizing the LL problems have been studied in [29], which justifies the significant computational advantages of ANIL compared with MAML in theory. Furthermore, the finite-time analysis of AC algorithms has shown [30] that, once the learning rates at both the actor and critic sides are chosen properly, a two timescale AC algorithm can achieve an $\mathcal{O}(\epsilon^{-2.5})$ iteration complexity for finding the first-order stationary points (FOSPs) of general nonconvex reward functions.

Besides these theoretical analyses in a specific learning setting, the algorithm design and corresponding convergence analysis for general bilevel optimization solvers have been recently advancing at a rapid pace under certain assumptions that the UL objective function is general nonconvex while the LL objective functions are strongly convex, which covers the existing convergence results shown for AC algorithms. The typical algorithms include those with double-loop structure, those with two timescale or single timescale but single-loop, and those with error-correction or accelerated/variance-reduction. To be more specific, double-loop algorithms, such as bilevel stochastic approximation (BSA) methods [31] and stochastic bilevel optimizers (stoBiO) [32], mainly request an inner loop to solve the LL problem up to a certain error tolerance or with a certain number of iterations and then switch back to optimize the UL problem, which can achieve an $\mathcal{O}(\epsilon^{-2})$ convergence rate to the ϵ -FOSPs. In practice, single-loop algorithms are implemented more efficiently in terms of computational complexity and hyperparameter tuning compared to double-loop algorithms. A two-timescale stochastic approximation (TTSA) was analyzed in [33], but it is shown that TTSA needs $\mathcal{O}(\epsilon^{-2.5})$ number of iterations to achieve the ϵ -FOSPs. Later, an error correction method, named the Single-Timescale stochastic BiLevel optimization (STABLE) method [34], improves the convergence rate of the single-loop algorithm to $\mathcal{O}(\epsilon^{-2})$ and a tighter analysis for ALternating Stochastic gradient dEscenT (ALSET) shows that the single-loop algorithm can also achieve a convergence of $\mathcal{O}(\epsilon^{-2})$ without the error correction technique. When more advanced momentum-assisted or variance reduction methods are adopted in the algorithm design, the subsequent works, such as the momentum-based recursive bilevel optimizer (MRBO) [35] and the single-timescale double-momentum stochastic approximation (SUSTAIN) [36] and the variance reduced BiAdam (VR-BiAdam) [37], can sharpen the convergence rate of bilevel algorithms to $\mathcal{O}(\epsilon^{-1.5})$.

For the theoretical works on MAML/MARL, it is shown in [22, 24] that when the critic side is allowed the consensus step at each agent to approximate the networked rewards, MDAC algorithms can achieve an $\mathcal{O}(\epsilon^{-2})$ convergence rate to FOSPs, but both of them require an inner loop procedure for the LL problem which makes the algorithms double loop. Dif-MAML [14] is able to perform the UL consensus-based meta learning, but iMAML considered in Dif-MAML is only a very special case of bilevel. Thus, the applicability of Dif-MAML is restrictive. One of the closest works to ours is coordinated AC (CAC) [23], which can realize the consensus on both UL and LL problems with $\mathcal{O}(\epsilon^{-2.5})$ number of iterations and is only for DAC problems. A theoretical comparison between our

Table 1: A comparison with closely related prior work on (decentralized) bilevel optimization learning. ‘‘Comm.’’ refer to whether the algorithm only needs one round of communication at either UL or LL per iteration; ‘‘Alg.’’ refs to the types of the basic stochastic algorithms adopted in the method.

Prior work	Consensus		Method	Rate	Comm.	Alg.	Setting
	UL	LL					
Ghadimi et al. [31]			BSA	$\mathcal{O}(1/\epsilon^2)$	-	SGD	bilevel
Hong et al.[33]			TTSA	$\mathcal{O}(1/\epsilon^{2.5})$	-	SGD	bilevel
Chen et al. [43]			ALSET	$\mathcal{O}(1/\epsilon^2)$	-	SGD	bilevel
Khanduri et al. [36]			SUSTAIN	$\mathcal{O}(1/\epsilon^{1.5})$	-	Momentum	bilevel
Kayaalp et al. [14]	✓		Dif-MAML	$\mathcal{O}(1/\epsilon^2)$	✓	SGD	iMAML
Kaiqing et al. [15]		✓	DAC	-	✓	PG	MARL
Chen et al. [22]		✓	MDAC	$\mathcal{O}(1/\epsilon^2)$		PG	MARL
Hairi et al. [24]		✓	MDAC	$\mathcal{O}(1/\epsilon^2)$		PG	MARL
Zeng et al. [23]	✓	✓	CAC	$\mathcal{O}(1/\epsilon^{2.5})$	✓	PG	MARL
This work	✓	✓	SLAM	$\mathcal{O}(1/(n\epsilon^2))$	✓	SGD/PG	bilevel

work and closely related previous works on bilevel programming is shown in Table 1. There is a line of independent work on decentralized optimization. But the existing works are only suitable for single-level minimization of only nonconvex problems, such as distributed SGD [38, 39], stochastic gradient tracking [40, 41] and stochastic primal dual algorithm [42], which can achieve an $\mathcal{O}(1/(n\epsilon^2))$ convergence rate to FOSPs for general nonconvex objective function optimization problems.

Main Contributions of This Work. In this work, we consider a very general DBO setting, where both UL and LL problems can include a consensus constraint for model parameter sharing and there would be multiple LL problems coupled with the UL problem. To solve this problem efficiently in a fully decentralized way, we propose a Stochastic Linearized Augmented Lagrangian Method (SLAM) for dealing with both of the two levels of the optimization processes and the consensus constraints at each level. Leveraging the linearized augmented Lagrangian function as a surrogate, the design of SLAM is simple and easily implemented as it is a single-loop algorithm with only step sizes to be tuned for convergence. We make the standard assumptions on Lipschitz continuity and convexity for both the UL and LL optimization problems as shown in the existing literature. We establish the conditions of SLAM w.r.t. convergence to ϵ -Karush-Kuhn-Tucker (KKT) points of problem (1) at a rate of $\mathcal{O}(1/(n\epsilon^2))$, matching the standard convergence rate achieved by decentralized SGD type of algorithm to FOSPs for only single-level nonconvex minimization problems. Remarkably, through numerical experiments on MARL problems, it is observed that SLAM can converge faster than the existing MARL methods and even achieve higher rewards in most cases.

To summarize, the main contributions of this work are highlighted as follows:

- ▶ Our proposed SLAM algorithm is generic, and thus generalizes the single agent-based bilevel algorithms to the multi-agent setting and is amenable to be specialized to solve multiple consensus-based DBO problems.
- ▶ SLAM is a single-timescale and single-loop algorithm that can find the ϵ -KKT points at a rate of $\mathcal{O}(1/(n\epsilon^2))$, which shows a linear speedup w.r.t. the number of nodes. To the best of our knowledge, this is the first work that shows a decentralized stochastic algorithm can achieve this rate under the constraints where any level or both levels of the DBO problem requires the consensus process.
- ▶ Numerical results that illustrate the proposed SLAM outperforms the state-of-the-art MARL algorithms over heterogeneous networks in terms of both convergence speed and achievable rewards.

Due to space limitations, all technical proofs are deferred to the supplement.

2 Decentralized Bilevel Optimization Framework

Problem formulation of DBO. One of the main motivations for performing decentralized joint learning is dealing with large-scale dataset or scattered data samples. At each node, the loss function can be written as $f_i(\mathbf{x}_i, \mathbf{y}_{i,1}^*(\mathbf{x}_i), \dots, \mathbf{y}_{i,m}^*(\mathbf{x}_i)) \triangleq \mathbb{E}_{\xi \in \mathcal{D}_i^U} [F_i(\mathbf{x}_i, \mathbf{y}_{i,1}^*(\mathbf{x}_i), \dots, \mathbf{y}_{i,m}^*(\mathbf{x}_i); \xi)]$, where \mathcal{D}_i^U denotes the local data distributions at the UL optimization problem, and $F_i(\mathbf{x}_i, \mathbf{y}_{i,1}^*(\mathbf{x}_i), \dots, \mathbf{y}_{i,m}^*(\mathbf{x}_i); \xi)$ represents the estimation error of the UL learning model on

data $\xi \in \mathcal{D}_i^U$. Similarly, the LL learning tasks also include randomly sampled data from a local distribution $\mathcal{D}_{i,k}^L$ for task k , so the LL cost function at each node can be expressed as $g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) \triangleq \mathbb{E}_{\zeta \in \mathcal{D}_{i,k}^L} [G_{i,k}(\mathbf{x}_i, \mathbf{y}_i; \zeta)]$, $\forall k$, where $G_{i,k}$ denotes the estimation error of the LL learning model over $\mathbf{y}_{k,i}$ on data $\zeta \in \mathcal{D}_{i,k}^L$. It is well known that SGD is one of the most efficient algorithms for tackling large amounts of data samples. Before showing the algorithm design, we first reformulate problem (1) in a concise and compact way from a global view of the variables. Let $\mathbf{x} \triangleq [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$ and $\mathbf{y}_k \triangleq [\mathbf{y}_{1,k}, \dots, \mathbf{y}_{n,k}]^T$. Then, problem (1) can be rewritten by concatenated variables as

$$\min_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}_i)) \quad (2a)$$

$$\text{s.t. } \mathbf{A}\mathbf{x} = 0, \quad (2b)$$

$$\mathbf{y}_k^*(\mathbf{x}) = \arg \min_{\mathbf{y}_k} g_k(\mathbf{x}, \mathbf{y}_k) \triangleq \frac{1}{n} \sum_{i=1}^n g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) \quad \text{s.t. } \mathbf{A}\mathbf{y}_k = 0, \forall k \in [m], \quad (2c)$$

where $g_k(\mathbf{x}, \mathbf{y}_k)$ denotes the k th LL loss function, $\mathbf{A} \in \mathbb{R}^{|\mathcal{E}| \times n}$ represents the incidence matrix¹ and $f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}_i))$ abbreviates $f_i(\mathbf{x}_i, \mathbf{y}_{i,1}^*(\mathbf{x}_i), \dots, \mathbf{y}_{i,m}^*(\mathbf{x}_i))$ for notational brevity.

Algorithm Design. Towards this end, it is straightforward to construct a variant of the classical augmented Lagrangian function for the UL optimization problem as

$$\mathcal{L}_{\rho\gamma}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) + \gamma \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} \rangle + \frac{\rho\gamma}{2} \|\mathbf{A}\mathbf{x}\|^2, \quad (3)$$

where $\boldsymbol{\lambda}$ denotes the dual variable (Lagrangian multiplier) for the consensus constraint, $\rho > 0$, and γ is a scaling factor (which will be determined later).

Motivated by the primal-dual optimization framework [44], one step of gradient descent based on the linearized objective function with a following gradient ascent step is sufficient for the minimization of the general nonconvex loss function under the linear constraints, which means that there is no need to solve an inner optimization problem before updating the Lagrangian multiplier as is done in the classical augmented Lagrangian method.

When both the UL and LL objective functions are differentiable and the inverse of the Hessian matrix at the LL problem exists, i.e., $\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x}))$ is invertible, then there exists a closed form for $\nabla f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}_i))$. Following the existing works on bilevel algorithm designs, replacing $\mathbf{y}_{i,k}^*(\mathbf{x}_i)$ by $\mathbf{y}_{i,k}$ in the gradient of $f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}_i))$ w.r.t. \mathbf{x}_i can serve as an efficient surrogate for the stochastic gradient estimate. However, in the decentralized setting, only individual loss functions are observable at each agent, therefore, the local UL implicit gradient is computed through replacing $g_k(\mathbf{x}, \mathbf{y}_k)$ by $g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k})$, denoted as $\bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}_{i,k})$. Let $\mathbf{h}_{g,k}^r$ and \mathbf{h}_f^r respectively denote the distributed stochastic gradient estimate of the LL and UL objective functions at points $(\mathbf{x}^r, \mathbf{y}_k^r)$ and $(\mathbf{x}^r, \mathbf{y}_k^{r+1})$, $\forall k$, w.r.t. \mathbf{y}_k and \mathbf{x} , where r represents the index of iterations. Thus, our proposed SLAM can be expressed as

$$\mathbf{y}_k^{r+1} = \arg \min_{\mathbf{y}_k} (\mathbf{h}_{g,k}^r + \gamma \mathbf{A}^T (\boldsymbol{\omega}_k^r + \rho \mathbf{A} \mathbf{y}_k^r), \mathbf{y}_k - \mathbf{y}_k^r) + \frac{\beta}{2} \|\mathbf{y}_k - \mathbf{y}_k^r\|^2, \quad \forall k, \quad (4a)$$

$$\boldsymbol{\omega}_k^{r+1} = \boldsymbol{\omega}_k^r + \frac{\rho}{\gamma} \mathbf{A} \mathbf{y}_k^{r+1}, \quad \forall k, \quad (4b)$$

$$\mathbf{x}^{r+1} = \arg \min_{\mathbf{x}} (\mathbf{h}_f^r + \gamma \mathbf{A}^T (\boldsymbol{\lambda}^r + \rho \mathbf{A} \mathbf{x}^r), \mathbf{x} - \mathbf{x}^r) + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}^r\|^2, \quad (4c)$$

$$\boldsymbol{\lambda}^{r+1} = \boldsymbol{\lambda}^r + \frac{\rho}{\gamma} \mathbf{A} \mathbf{x}^{r+1}, \quad (4d)$$

where $\boldsymbol{\omega}_k$ is the dual variable for ensuring the LL consensus process for each learning task, α and β are the parameters of the quadratic penalization terms, and ρ/γ here is the step-size for the updates of the dual variables.

Implementation of SLAM. Noting that the objective functions in each subproblem, i.e., (4a) and (4c), are quadratic, we can easily have the updates of both UL and LL optimization variables as

¹Here, we assume the problem dimension is 1, without loss of generality, to simplify the notation.

$$\mathbf{y}_k^{r+1} = \mathbf{y}_k^r - \frac{1}{\beta} (\mathbf{h}_{g,k}^r + \gamma \mathbf{A}^T \boldsymbol{\omega}_k^r + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbf{y}_k^r), \forall k, \quad (5a)$$

$$\mathbf{x}^{r+1} = \mathbf{x}^r - \frac{1}{\alpha} (\mathbf{h}_f^r + \gamma \mathbf{A}^T \boldsymbol{\lambda}^r + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbf{x}^r), \quad (5b)$$

where $1/\alpha$ and $1/\beta$ serve as the step-sizes of updating both UL and LL learning models. Subtracting the equality with the same one from the previous iteration for both (5a) and (5b) ends up with efficient model updates of both the UL and LL learning problems as follows:

$$\mathbf{y}_k^{r+1} = 2\mathbf{W}_g \mathbf{y}_k^r - \mathbf{W}'_g \mathbf{y}_k^{r-1} - \frac{1}{\beta} (\mathbf{h}_{g,k}^r - \mathbf{h}_{g,k}^{r-1}), \quad \forall k, \quad (6a)$$

$$\mathbf{x}^{r+1} = 2\mathbf{W}_f \mathbf{x}^r - \mathbf{W}'_f \mathbf{x}^{r-1} - \frac{1}{\alpha} (\mathbf{h}_f^r - \mathbf{h}_f^{r-1}), \quad (6b)$$

where the mixing matrices, with $\tau_g = \beta/\gamma$ and $\tau_f = \alpha/\gamma$, are defined as

$$\mathbf{W}_g \triangleq \mathbf{I} - \frac{(1 + \gamma^{-1})\rho}{2\tau_g} \mathbf{A}^T \mathbf{A}, \quad \mathbf{W}'_g \triangleq \mathbf{I} - \frac{\rho}{\tau_g} \mathbf{A}^T \mathbf{A}, \quad (7a)$$

$$\mathbf{W}_f \triangleq \mathbf{I} - \frac{(1 + \gamma^{-1})\rho}{2\tau_f} \mathbf{A}^T \mathbf{A}, \quad \mathbf{W}'_f \triangleq \mathbf{I} - \frac{\rho}{\tau_f} \mathbf{A}^T \mathbf{A}. \quad (7b)$$

According to (6a) and (6b), it can be readily observed that SLAM is amenable to a fully decentralized implementation. The detailed algorithm description is provided in Algorithm 1 from a local view of the model update, where $[\mathbf{W}]_{ij}$ denotes the ij th entry of matrix \mathbf{W} , $[\mathbf{h}_g^r]_{i,k}$ is the gradient estimate of $\nabla g_{i,k}(\mathbf{x}_i^r, \mathbf{y}_{i,k}^r)$ (i.e., $\mathbf{h}_{g,k}^r = [[\mathbf{h}_g^r]_{1,k}, \dots, [\mathbf{h}_g^r]_{n,k}]^T$), and similarly $[\mathbf{h}_f^r]_i$ is the local gradient estimate of $\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^{r+1})$ (i.e., $\mathbf{h}_f^r = [[\mathbf{h}_f^r]_1, \dots, [\mathbf{h}_f^r]_n]^T$).

Algorithm 1 Decentralized implementation of SLAM

Initialization: $\alpha, \beta, \gamma, \mathbf{x}_i^1, \mathbf{y}_{i,k}^1, \forall i, k$, and set $\boldsymbol{\lambda}^1 = \boldsymbol{\omega}_k^1 = 0, \forall k$;
1: **for** $r = 1, 2, \dots, T$ **do**
2: **for** $i = 1, 2, \dots, n$ in parallel over the network **do**
3: Estimate gradient $\nabla g_{i,k}(\mathbf{x}_i^r, \mathbf{y}_{i,k}^r)$ for each task and $\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^{r+1})$ locally
4: $\mathbf{y}_{i,k}^{r+1} = 2 \sum_{j \in \mathcal{N}_i} [\mathbf{W}_g]_{ij} \mathbf{y}_{j,k}^r - [\mathbf{W}'_g]_{ij} \mathbf{y}_{j,k}^{r-1} - \beta^{-1} ([\mathbf{h}_g^r]_{i,k} - [\mathbf{h}_g^{r-1}]_{i,k}) \triangleright$ LL models
5: $\mathbf{x}_i^{r+1} = 2 \sum_{j \in \mathcal{N}_i} [\mathbf{W}_f]_{ij} \mathbf{x}_j^r - [\mathbf{W}'_f]_{ij} \mathbf{x}_j^{r-1} - \alpha^{-1} ([\mathbf{h}_f^r]_i - [\mathbf{h}_f^{r-1}]_i) \triangleright$ UL model
6: **end for**
7: **end for**

Besides, if there is a consensus requirement at only one level of the optimization problem, then the problem at the other level becomes one with multiple objective functions. Our proposed SLAM can also be applied for solving any of these problems by a minor revision of the generic SLAM formulation. To be more specific, we provide the following discussion.

A Special Case of DBO (1) (with only consensus in the LL problems). If there is only a need for consensus of LL model parameters, then problem (2) reduces to the following DBO problem. For example, in solving multi-agent actor-critic RL problems, the UL optimization problem consists of improving the policy for each agent while the LL problem requires all the agents to jointly evaluate the value function over the whole network. The DBO problem is then expressed as

$$\min_{\mathbf{x}_i} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}_i)), \quad \forall i \in [n] \quad (8a)$$

$$\text{s.t. } \mathbf{y}_k^*(\mathbf{x}) = \arg \min_{\mathbf{y}_k} g_k(\mathbf{x}, \mathbf{y}_k) \triangleq \frac{1}{n} \sum_{i=1}^n g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) \quad \text{s.t. } \mathbf{A} \mathbf{y}_k = 0, \forall k \in [m]. \quad (8b)$$

The major difference between problem (2) and (8) is that the UL optimization problem includes multiple objectives over the model parameters $\mathbf{x}_i, \forall i \in [n]$. In this case, the updating rule of variable \mathbf{x} in (6b) reduces to $\mathbf{x}^{r+1} = \mathbf{x}^r - \mathbf{h}_f^r/\alpha$ by forgoing the dual update w.r.t. $\boldsymbol{\lambda}$. The detailed implementation is summarized in Algorithm 2, where we name this special case of SLAM by SLAM-L as the LL consensus process is the main feature in this setting.

Algorithm 2 Decentralized implementation of SLAM-L

Initialization: $\alpha, \beta, \gamma, \mathbf{x}_i^1, \mathbf{y}_{i,k}^1, \forall i, k$, and set $\omega_k^1 = 0, \forall k$;

- 1: **for** $r = 1, 2, \dots, T$ **do**
 - 2: **for** $i = 1, 2, \dots, n$ in parallel over the network **do**
 - 3: Estimate gradient $\nabla g_{i,k}(\mathbf{x}_i^r, \mathbf{y}_{i,k}^r)$ for each task and $\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^{r+1})$ locally
 - 4: $\mathbf{y}_{i,k}^{r+1} = 2 \sum_{j \in \mathcal{N}_i} [\mathbf{W}_g]_{ij} \mathbf{y}_{j,k}^r - [\mathbf{W}'_g]_{ij} \mathbf{y}_{j,k}^{r-1} - \beta^{-1} ([\mathbf{h}_g^r]_{i,k} - [\mathbf{h}_g^{r-1}]_{i,k})$ \triangleright LL models
 - 5: $\mathbf{x}_i^{r+1} = \mathbf{x}_i^r - \alpha^{-1} [\mathbf{h}_f^r]_i, \forall i$ \triangleright UL models
 - 6: **end for**
 - 7: **end for**
-

A Special Case of DBO (1) (with only consensus in the UL problem). The other special is analogous to the first one with the difference being the absence of the LL consensus process in comparison to (2), which is written as follows:

$$\min_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}_i)) \quad (9a)$$

$$\text{s.t. } \mathbf{A}\mathbf{x} = 0, \quad \mathbf{y}_{i,k}^*(\mathbf{x}_i) = \arg \min_{\mathbf{y}_{i,k}} g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}), \forall i \in [n], \forall k \in [m], \quad (9b)$$

where there are multiple objectives in the LL optimization problems. Problem (9) also covers a wide range of applications in machine learning, e.g., multi-task and/or personalized learning, and so on. In this case, the update of variable \mathbf{y}_k shown in (5a) is changed to $\mathbf{y}_k^{r+1} = \mathbf{y}_k^r - \mathbf{h}_g^r / \beta$ as there is no consensus constraint involved. Analogous to the previous case, the implementation of this algorithm is presented in Algorithm 3 and termed as SLAM-U.

Algorithm 3 Decentralized implementation of SLAM-U

Initialization: $\alpha, \beta, \gamma, \mathbf{x}_i^1, \mathbf{y}_{i,k}^1, \forall i, k$, and set $\lambda^1 = 0, \forall k$;

- 1: **for** $r = 1, 2, \dots, T$ **do**
 - 2: **for** $i = 1, 2, \dots, n$ in parallel over the network **do**
 - 3: Estimate gradient $\nabla g_{i,k}(\mathbf{x}_i^r, \mathbf{y}_{i,k}^r)$ for each task and $\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^{r+1})$ locally
 - 4: $\mathbf{y}_{i,k}^{r+1} = \mathbf{y}_{i,k}^r - \beta^{-1} [\mathbf{h}_g^r]_{i,k}$ \triangleright LL models
 - 5: $\mathbf{x}_i^{r+1} = 2 \sum_{j \in \mathcal{N}_i} [\mathbf{W}_f]_{ij} \mathbf{x}_j^r - [\mathbf{W}'_f]_{ij} \mathbf{x}_j^{r-1} - \alpha^{-1} ([\mathbf{h}_f^r]_i - [\mathbf{h}_f^{r-1}]_i)$ \triangleright UL model
 - 6: **end for**
 - 7: **end for**
-

3 Theoretical Convergence Results

Before showing the theoretical results about the convergence guarantees of SLAM, we first need five main classes of assumptions used in showing the descent of some quantifiable function so that SLAM can reach the ϵ -KKT points of the DBO problems. More detailed definitions and properties regarding these assumptions are deferred to the supplement.

3.1 Assumptions

Our theoretical results are based on the following assumptions on the properties of the loss functions in both the UL and LL optimization problems, which are mainly related to the continuity of the objective function and stochasticity of the gradient estimates.

- A1. (Lipschitz continuity of both UL and LL objective functions) Assume that functions $f_i(\cdot), \nabla f_i(\cdot), \nabla g_{i,k}(\cdot), \nabla^2 g_{i,k}(\cdot), \forall i$, are (block-wise) Lipschitz continuous with constants $L_{f,0}, L_{f,1}, L_{g,1}, L_{g,2}$ for both \mathbf{x} and $\mathbf{y}_k, \forall k$, and $\nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\cdot), \forall i$ are bounded by C_{xy} .
- A2. (Connectivity of graph \mathcal{G}) The communication graph \mathcal{G} is assumed to be well-connected, i.e., $\mathbb{1}^T \mathbf{L} = 0$ where $\mathbf{L} = \mathbf{A}^T \mathbf{A}$, and the second-smallest eigenvalue of \mathbf{L} is assumed to be strictly positive, i.e., $\tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A}) > 0$.

- A3. (Quality of the stochastic gradient estimate) The stochastic estimates of $\nabla f_i(\mathbf{x}_i, \mathbf{y}_{i,k})$, $\nabla_{\mathbf{y}_i} g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k})$, $\forall i, k$, are unbiased and their variances are bounded by σ_f^2, σ_g^2 .
- A4. Assume that the UL objective functions $f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}_i))$, $\forall i, k$ are lower bounded.
- A5. (Strong convexity of $g_{i,k}(\cdot)$ w.r.t. $\mathbf{y}_{i,k}$) Function $g_{i,k}(\cdot)$ is μ_g -strongly convex w.r.t. $\mathbf{y}_{i,k}$, $\forall i, k$.

Note that these assumptions are commonly used in the convergence analysis for bilevel and decentralized optimization algorithms. Given these assumptions, we are now in a position to provide the following theoretical convergence guarantees.

3.2 Convergence Rates of SLAM

Theorem 1. (Convergence rate of SLAM to ϵ -KKT points) Suppose that A1-A5 hold and assume $\|\nabla_{\mathbf{y}_i}^2 g_{i,k}(\cdot, \mathbf{y}_i) - n^{-1} \sum_{i=1}^n \nabla_{\mathbf{y}_i}^2 g_{i,k}(\cdot, \mathbf{y}'_i)\| \leq L_g \|\mathbf{y}_i - \mathbf{y}'_i\|$, $\forall i, k$ if $\nabla^2 g_{i,k}(\cdot)$, $\forall i, k$ are required in computing the UL implicit gradient. When step-sizes are chosen as $1/\alpha \sim 1/\beta \sim \mathcal{O}(\sqrt{n/T})$, $\tau_f, \tau_g \geq \mathcal{O}(\rho \sigma_{\max}(\mathbf{A}^T \mathbf{A}))$, the mini-batch size of \mathbf{h}_f^r is $\mathcal{O}(\log(nT))$, then the iterates $\{\mathbf{x}^r, \boldsymbol{\lambda}^r, \mathbf{y}_k^r, \boldsymbol{\omega}_k^r, \forall k, r\}$ generated by SLAM satisfy

$$\text{UL: } \frac{1}{T} \sum_{r=1}^T \mathbb{E}[\|\nabla f(\mathbb{1}\bar{\mathbf{x}}^r, \mathbf{y}_1^*(\mathbb{1}\bar{\mathbf{x}}^r), \dots, \mathbf{y}_m^*(\mathbb{1}\bar{\mathbf{x}}^r))\|^2] \sim \frac{1}{T} \sum_{r=1}^T \mathbb{E}[\|\mathbf{A}\mathbf{x}^r\|^2] \sim \mathcal{O}(1/\sqrt{nT}), \quad (10a)$$

$$\text{LL: } \frac{1}{T} \sum_{r=1}^T \mathbb{E}[\|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2] \sim \frac{1}{T} \sum_{r=1}^T \mathbb{E}[\|\mathbf{A}\mathbf{y}_k^r\|^2] \sim \mathcal{O}(1/\sqrt{nT}), \quad \forall k, \quad (10b)$$

where $\bar{\mathbf{x}} = n^{-1} \mathbb{1}^T \mathbf{x}$, and T denotes the total number of iterations.

Remark 1. It is noted in Theorem 1 that the convergence rate achieved by SLAM to find the ϵ -approximate KKT points of (1) (including both the first-order stationarity of the solutions and the violation of constraints) is on the order of $1/(n\epsilon^2)$. Therefore, it follows that a linear speedup w.r.t. the number of learners can be achieved by SLAM for DBO, matching the classical results of distributed SGD for only single-level general nonconvex problems.

Remark 2. In comparison with existing bilevel algorithms, SLAM is a single timescale algorithm since the learning rates can be chosen as $1/\alpha \sim 1/\beta$, which is consistent with ALSET [43].

Remark 3. The major novelty of obtaining these theoretical results relies on the developed variant of the augmented Lagrangian function and subsequently derived recursion of the successive dual variables, which quantify well the consensus errors resulting from both UL and LL optimization processes in terms of primal variables. Note that this is distinct from the existing theoretical analysis of stochastic algorithms, such as distributed SGD [38, 39], stochastic gradient tracking [40, 41], stochastic primal-dual algorithms [42, 45], etc.

Corollary 1. (Convergence rate of SLAM-L to ϵ -KKT points) Suppose that A1-A5 hold and assume $\|\nabla_{\mathbf{y}_i}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_i) - \nabla_{\mathbf{y}_i}^2 g_{i,k}(\mathbf{x}, \mathbf{y}'_i)\| \leq L_g \|\mathbf{y}_i - \mathbf{y}'_i\|$, $\forall i, k$ if $\nabla^2 g_{i,k}(\cdot)$, $\forall i, k$ are required in computing the UL implicit gradient. When step-sizes are chosen as $1/\alpha \sim \mathcal{O}(1/\sqrt{T})$, $1/\beta \sim \mathcal{O}(\sqrt{n/T})$, $\tau_f, \tau_g \geq \mathcal{O}(\rho \sigma_{\max}(\mathbf{A}^T \mathbf{A}))$, $\rho \geq n$, the mini-batch size of \mathbf{h}_f^r is $\mathcal{O}(\log(nT))$, the iterates $\{\mathbf{x}^r, \mathbf{y}_k^r, \boldsymbol{\omega}_k^r, \forall k, r\}$ generated by SLAM-L satisfy

$$\text{UL: } \frac{1}{T} \sum_{r=1}^T \mathbb{E}[\|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,1}^*(\mathbf{x}_i^r), \dots, \mathbf{y}_{i,m}^*(\mathbf{x}_i^r))\|^2], \forall i \sim \mathcal{O}(n/\sqrt{T}) \quad \text{and} \quad \text{LL: } (10b).$$

Remark 4. Different from Theorem 1, the stationarity of the UL model parameters requires the shrinkage of the gradient size over each individual UL problem as shown in Corollary 1, so there is no speedup on the convergence rate guarantee at UL.

Corollary 2. (Convergence rate of SLAM-U to ϵ -KKT points) Suppose that A1-A5 hold. Given the conditions on $1/\alpha, 1/\beta, \tau_f, \tau_g$ and the mini-batch size of \mathbf{h}_f^r shown in Theorem 1, the iterates $\{\mathbf{x}^r, \boldsymbol{\lambda}^r, \mathbf{y}_k^r, \forall k, r\}$ generated by SLAM-U satisfy

$$\text{UL: } (10a) \quad \text{and} \quad \text{LL: } \frac{1}{T} \sum_{r=1}^T \mathbb{E}[\|\mathbf{y}_k^r - \mathbf{y}_k^*(\mathbf{x}^r)\|^2] \sim \mathcal{O}(1/\sqrt{nT}), \forall k.$$

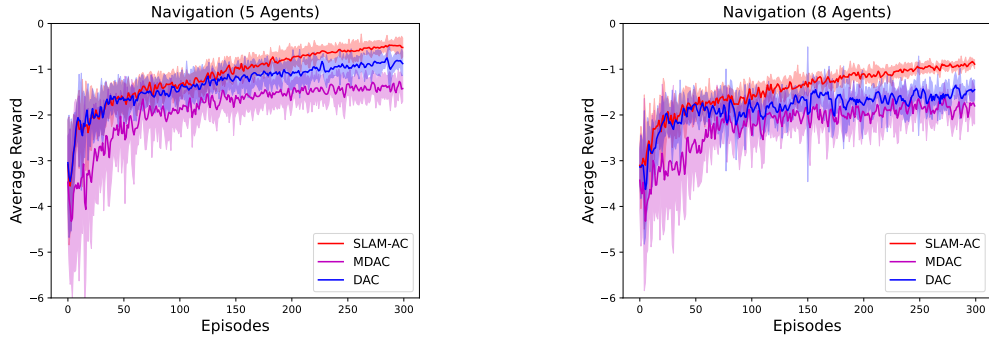


Figure 1: The averaged reward versus the learning process on the cooperative navigation task.

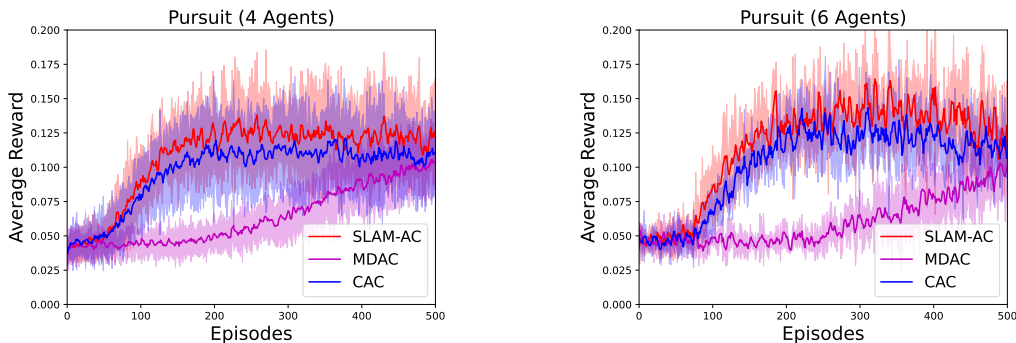


Figure 2: The averaged reward versus the learning process on the pursuit-evasion game. (Consensus with one layer of the actor neural nets and all layers of the critic neural nets.)

4 Numerical Results

In this section, we evaluate our proposed algorithm using two MARL environments: 1) the cooperative navigation task [16], which is built on the OpenAI Gym platform [46]; and 2) the pursuit-evasion game [47], which is built on the PettingZoo platform [48]. Detailed experimental settings and additional numerical results are provided in the supplement.

Cooperative Navigation Task. In this game, we consider that the n agents are aiming to jointly reach n different landmarks as soon as possible, where the Erdos Renyi Graph is used. We assume that each agent can observe the global state and has 5 possible actions: stay, left, right, up, and down. This task consists of a shared common goal of avoiding collision among the agents while they navigate to the targeting landmarks. In the simulations, each agent locally maintains two fully connected neural networks as the actor network (at UL w.r.t. \mathbf{x}_i) and the critic network (at LL w.r.t. \mathbf{y}_i), respectively. Moreover, each agent shares its critic network with its neighbors to cooperatively estimate the global value function and independently train its actor network to complete its local task.

We compare the performance of our proposed SLAM with application to the DAC setting, named SLAM-AC, with two benchmark algorithms: DAC [15] and mini-batch DAC (MDAC) [22]. Theoretically, MDAC needs an $\mathcal{O}(\epsilon^{-1} \ln \epsilon^{-1})$ batch size in its inner loop to update critic parameters before each update in policy parameters, which is not practical. Here, we set a small batch $B = 10$ in the inner loop for MDAC to achieve fast convergence. The simulation results on this coordination game are presented in Figure 1, where the performance is averaged over 5 independent Monte Carlo (MC) trials for each algorithm.

Pursuit-Evasion Game. In the pursuit-evasion game, there are two groups of nodes: pursuers (agents) and evaders. The agents are connected through a ring graph. Pursuers could observe the global state of the video game. An evader is considered caught if two pursuers simultaneously arrive at the evader’s location. As each pursuer should learn to cooperate with other pursuers to catch the

evaders, the pursuers share certain similarities with each other since they need to follow similar strategies to achieve their local tasks: simultaneously catching an evader with other pursuers.

We follow the experimental set up in [23], where all agents partially share their actor networks with neighbors for collaborations in their policy spaces and fully share their critic network to cooperatively learn the global value function. In Figure 2, we compare SLAM-AC with two benchmarks, CAC [23] and MDAC [22], with 5 MC trials again. To ensure a fair comparison, all algorithms use the same parameter sharing scheme mentioned above. Note that CAC [23] is a variant of DAC [15] and the only difference is that CAC can partially share its policy parameters while the policy parameters are not shared in DAC. In the experiment, each agent maintains two convolutional neural networks, one for the actor and one for the critic (Please refer to the supplement for detailed structures).

5 Concluding Remark

In this paper, we studied a generic form of the DBO problem, which is shown to have three major variants that formulate multiple hierarchical machine learning problems. Targeting these DBO problems, we proposed SLAM – a simple and elegant algorithm to solve DBO in a fully decentralized way. Under mild conditions, we establish theoretical results showing that our proposed SLAM is able to find the ϵ -KKT points with a convergence rate of $\mathcal{O}(1/(n\epsilon^2))$, which matches the standard convergence rate achieved by the classical distributed SGD algorithms for solving only single-level general nonconvex optimization problems. We tested the performance of SLAM numerically on a MARL scenario and found that SLAM outperformed the traditional AC algorithms w.r.t. convergence speed and (in most cases) achievable rewards.

Societal impact. To the best of our knowledge, we do not see any ethical or negative immediate societal consequence of this work.

Acknowledgments

M. Hong and S. Zeng are partially supported by NSF grants CIF-1910385 and CMMI-1727757.

References

- [1] V. Stackelberg, Heinrich, Von, and S. Heinrich, *The Theory of the Market Economy*. Oxford University Press, 1952.
- [2] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *Proceedings of International Conference on Machine Learning (ICML)*, pp. 1126–1135, 2017.
- [3] H. Liu, R. Socher, and C. Xiong, “Taming MAML: Efficient unbiased meta-reinforcement learning,” in *Proceedings of International Conference on Machine Learning (ICML)*, pp. 4061–4071, 2019.
- [4] A. Rajeswaran, C. Finn, S. M. Kakade, and S. Levine, “Meta-learning with implicit gradients,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 32, 2019.
- [5] A. Fallah, A. Mokhtari, and A. Ozdaglar, “Generalization of model-agnostic meta-learning algorithms: Recurring and unseen tasks,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [6] A. Raghu, M. Raghu, S. Bengio, and O. Vinyals, “Rapid learning or feature reuse? towards understanding the effectiveness of MAML,” in *Proceedings of International Conference on Learning Representations (ICLR)*, 2020.
- [7] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT Press, 2018.
- [8] Z. Yang, Y. Chen, M. Hong, and Z. Wang, “Provably global convergence of actor-critic: A case for linear quadratic regulator with ergodic cost,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 32, 2019.
- [9] T. Xu, Z. Wang, and Y. Liang, “Improving sample complexity bounds for (natural) actor-critic algorithms,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, pp. 4358–4369, 2020.
- [10] A. Shaban, C.-A. Cheng, N. Hatch, and B. Boots, “Truncated back-propagation for bilevel optimization,” in *Proceedings of International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 1723–1732, 2019.
- [11] R. Ranjan, V. M. Patel, and R. Chellappa, “Hyperfacer: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 1, pp. 121–135, 2017.

- [12] A. Fallah, A. Mokhtari, and A. Ozdaglar, “Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [13] J. Li, F. Huang, and H. Huang, “Local stochastic bilevel optimization with momentum-based variance reduction,” *arXiv preprint arXiv:2205.01608*, 2022.
- [14] M. Kayaalp, S. Vlaski, and A. H. Sayed, “Dif-MAML: Decentralized multi-agent meta-learning,” *IEEE Open Journal of Signal Processing*, vol. 3, pp. 71–93, 2022.
- [15] K. Zhang, Z. Yang, H. Liu, T. Zhang, and T. Basar, “Fully decentralized multi-agent reinforcement learning with networked agents,” in *Proceedings of International Conference on Machine Learning (ICML)*, pp. 5872–5881, PMLR, 2018.
- [16] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, “Multi-agent actor-critic for mixed cooperative-competitive environments,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017.
- [17] J. G. Kuba, R. Chen, M. Wen, Y. Wen, F. Sun, J. Wang, and Y. Yang, “Trust region policy optimisation in multi-agent reinforcement learning,” *arXiv preprint arXiv:2109.11251*, 2021.
- [18] J. G. Kuba, M. Wen, L. Meng, H. Zhang, D. Mguni, J. Wang, and Y. Yang, “Settling the variance of multi-agent policy gradients,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, pp. 13458–13470, 2021.
- [19] S. Gu, J. G. Kuba, M. Wen, R. Chen, Z. Wang, Z. Tian, J. Wang, A. Knoll, and Y. Yang, “Multi-agent constrained policy optimisation,” *arXiv preprint arXiv:2110.02793*, 2021.
- [20] M. Wen, J. G. Kuba, R. Lin, W. Zhang, Y. Wen, J. Wang, and Y. Yang, “Multi-agent reinforcement learning is a sequence modeling problem,” *arXiv preprint arXiv:2205.14953*, 2022.
- [21] Y. Lin, K. Zhang, Z. Yang, Z. Wang, T. Başar, R. Sandhu, and J. Liu, “A communication-efficient multi-agent actor-critic algorithm for distributed reinforcement learning,” in *Proceedings of IEEE 58th Conference on Decision and Control (CDC)*, pp. 5562–5567, 2019.
- [22] Z. Chen, Y. Zhou, R.-R. Chen, and S. Zou, “Sample and communication-efficient decentralized actor-critic algorithms with finite-time analysis,” in *Proceedings of International Conference on Machine Learning (ICML)*, pp. 3794–3834, 2022.
- [23] S. Zeng, T. Chen, A. Garcia, and M. Hong, “Learning to coordinate in multi-agent systems: A coordinated actor-critic algorithm and finite-time guarantees,” in *Proceedings of the 4th Annual Learning for Dynamics and Control Conference*, pp. 278–290, 2022.
- [24] F. Hairi, J. Liu, and S. Lu, “Finite-time convergence and sample complexity of multi-agent actor-critic reinforcement learning with average reward,” in *Proceedings of International Conference on Learning Representations (ICLR)*, 2022.
- [25] D. Chen, Y. Li, and Q. Zhang, “Communication-efficient actor-critic methods for homogeneous markov games,” in *Proceedings of International Conference on Learning Representations (ICLR)*, 2022.
- [26] C. Finn, A. Rajeswaran, S. Kakade, and S. Levine, “Online meta-learning,” in *Proceedings of International Conference on Machine Learning (ICML)*, pp. 1920–1930, 2019.
- [27] M.-F. Balcan, M. Khodak, and A. Talwalkar, “Provable guarantees for gradient-based meta-learning,” in *Proceedings of International Conference on Machine Learning (ICML)*, pp. 424–433, 2019.
- [28] A. Fallah, A. Mokhtari, and A. Ozdaglar, “On the convergence theory of gradient-based model-agnostic meta-learning algorithms,” in *Proceedings of International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 1082–1092, 2020.
- [29] K. Ji, J. D. Lee, Y. Liang, and H. V. Poor, “Convergence of meta-learning with task-specific adaptation over partial parameters,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, pp. 11490–11500, 2020.
- [30] Y. F. Wu, W. Zhang, P. Xu, and Q. Gu, “A finite-time analysis of two time-scale actor-critic methods,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, pp. 17617–17628, 2020.
- [31] S. Ghadimi and M. Wang, “Approximation methods for bilevel programming,” *arXiv preprint arXiv:1802.02246*, 2018.
- [32] K. Ji, J. Yang, and Y. Liang, “Bilevel optimization: Convergence analysis and enhanced design,” in *Proceedings of International Conference on Machine Learning (ICML)*, pp. 4882–4892, 2021.
- [33] M. Hong, H.-T. Wai, Z. Wang, and Z. Yang, “A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic,” *arXiv preprint arXiv:2007.05170*, 2020.
- [34] T. Chen, Y. Sun, and W. Yin, “Tighter analysis of alternating stochastic gradient method for stochastic nested problems,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

- [35] J. Yang, K. Ji, and Y. Liang, “Provably faster algorithms for bilevel optimization,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021.
- [36] P. Khanduri, S. Zeng, M. Hong, H.-T. Wai, Z. Wang, and Z. Yang, “A near-optimal algorithm for stochastic bilevel optimization via double-momentum,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [37] F. Huang and H. Huang, “Biadam: Fast adaptive bilevel optimization methods,” *arXiv preprint arXiv:2106.11396*, 2021.
- [38] X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu, “Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017.
- [39] H. Tang, X. Lian, M. Yan, C. Zhang, and J. Liu, “D²: Decentralized training over decentralized data,” in *Proceedings of International Conference on Machine Learning (ICML)*, pp. 4848–4856, 2018.
- [40] S. Lu, X. Zhang, H. Sun, and M. Hong, “GNSD: a gradient-tracking based nonconvex stochastic algorithm for decentralized optimization,” in *Proceedings of IEEE Data Science Workshop (DSW)*, pp. 315–321, 2019.
- [41] A. Koloskova, T. Lin, and S. U. Stich, “An improved analysis of gradient tracking for decentralized machine learning,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021.
- [42] X. Yi, S. Zhang, T. Yang, T. Chai, and K. H. Johansson, “A primal-dual SGD algorithm for distributed nonconvex optimization,” *IEEE/CAA Journal of Automatica Sinica*, vol. 9, no. 5, pp. 812–833, 2022.
- [43] T. Chen, Y. Sun, and W. Yin, “Tighter analysis of alternating stochastic gradient method for stochastic nested problems,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [44] J. Nocedal and S. Wright, *Numerical Optimization*. Springer Science & Business Media, 2006.
- [45] G. Lan, S. Lee, and Y. Zhou, “Communication-efficient algorithms for decentralized and stochastic optimization,” *Mathematical Programming*, vol. 180, no. 1, pp. 237–284, 2020.
- [46] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, “OpenAI Gym,” *arXiv preprint arXiv:1606.01540*, 2016.
- [47] J. K. Gupta, M. Egorov, and M. Kochenderfer, “Cooperative multi-agent control using deep reinforcement learning,” in *Proceedings of International Conference on Autonomous Agents and Multiagent Systems*, pp. 66–83, Springer, 2017.
- [48] J. Terry, B. Black, N. Grammel, M. Jayakumar, A. Hari, R. Sullivan, L. S. Santos, C. Dieffendahl, C. Horsch, R. Perez-Vicente, *et al.*, “PettingZoo: Gym for multi-agent reinforcement learning,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, pp. 15032–15043, 2021.
- [49] Y. Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*, vol. 87. Springer Science & Business Media, 2003.
- [50] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, “Policy gradient methods for reinforcement learning with function approximation,” in *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, vol. 12, 1999.

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- ▶ Did you include the license to the code and datasets? **[Yes]** See
- ▶ Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- ▶ Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]**
 - (b) Did you describe the limitations of your work? **[Yes]**
 - (c) Did you discuss any potential negative societal impacts of your work? **[Yes]**
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - (b) Did you include complete proofs of all theoretical results? **[Yes]**
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[No]**
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]**
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[Yes]**
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]**
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? **[Yes]**
 - (b) Did you mention the license of the assets? **[Yes]**
 - (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **[N/A]**
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[N/A]**
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]**

Supplementary Material

A Preliminaries

In this section, we provide some technical preliminaries for the proofs of the lemmas and theorems claimed in the main body of this paper, including parameter definitions and supporting results. First, let us define the filtrations

$$\mathcal{F}^r = \{\mathbf{x}^r, \mathbf{y}^r, \boldsymbol{\lambda}^r, \boldsymbol{\omega}_k^r, \xi^{r-1}, \zeta^{r-1}, \dots, x^0, \boldsymbol{\lambda}^0\}, \quad (11)$$

$$\mathcal{F}^{r'} = \{\mathbf{x}^r, \mathbf{y}^{r+1}, \boldsymbol{\lambda}^r, \boldsymbol{\omega}_k^r, \xi^{r-1}, \zeta^r, \dots, x^0, \boldsymbol{\lambda}^0\}, \quad (12)$$

which will often be used in the proofs when taking conditional expectation.

Inequalities used in the proof include

1. Quadrilateral identity:

$$\langle \mathbf{x}^{r+1} - \mathbf{x}^r, \mathbf{x}^r - \mathbf{x}^{r-1} \rangle = \frac{1}{2} (\|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 + \|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 - \|\mathbf{w}^{r+1}\|^2) \quad (13)$$

where

$$\mathbf{w}^{r+1} \triangleq \mathbf{x}^{r+1} - \mathbf{x}^r - (\mathbf{x}^r - \mathbf{x}^{r-1}). \quad (14)$$

2. Young's inequality with parameter $\theta > 0$:

$$\langle \mathbf{x}, \mathbf{y} \rangle \leq \frac{1}{2\theta} \|\mathbf{x}\|^2 + \frac{\theta}{2} \|\mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y}. \quad (15)$$

3. Given vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$, the convexity of norm $\|\cdot\|^2$ and a trivial application of Jensen's inequality yields the following inequality

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right\|^2 \leq \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2. \quad (16)$$

4. Given that \mathbf{x} is a random vector, then

$$\mathbb{E} [\|\mathbf{x}\|^2] = \|\mathbb{E}[\mathbf{x}]\|^2 + \mathbb{E} [\|\mathbf{x} - \mathbb{E}[\mathbf{x}]\|^2]. \quad (17)$$

Lemma 1.

Under A1, A2 and A5, the gradient of $f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}))$ is given by

$$\begin{aligned} \nabla f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) &= \nabla_{\mathbf{x}_i} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) \\ &\quad - \sum_{k=1}^m \nabla_{\mathbf{x}_i \mathbf{y}_k}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_k^*(\mathbf{x})) [\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x}))]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})), \end{aligned} \quad (18)$$

where $\nabla_{\mathbf{x}_i} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}))$ denotes the gradient of the objective function w.r.t. the first argument.

Proof. In order to remove the ambiguity of the notation, we first define $F_i(\mathbf{x}) \triangleq f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}))$. Following the classical proving steps [43, Proposition 1], we obtain the closed form of the implicit gradient by the chain rule:

$$\nabla_{\mathbf{x}_i} F_i(\mathbf{x}) = \nabla_{\mathbf{x}_i} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) + \sum_{k=1}^m \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}^*}^T \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})). \quad (19)$$

Based on the definition of $\mathbf{y}_k^*(\mathbf{x})$, it follows that

$$\nabla_{\mathbf{y}_k} g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) = 0, \quad \mathbf{A} \mathbf{y}_k^*(\mathbf{x}) = 0, \quad (20)$$

and thus we have

$$\nabla_{\mathbf{x}_i} \left(\frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{y}_k} g_{i,k}(\mathbf{x}, \mathbf{y}_{i,k}^*(\mathbf{x})) \right) = 0, \quad \mathbf{A} \mathbf{y}_k^*(\mathbf{x}) = 0. \quad (21)$$

Therefore, we obtain

$$\nabla_{\mathbf{x}_i \mathbf{y}_k}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_k^*(\mathbf{x})) + \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}^*}^T \nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) = 0. \quad (22)$$

According to A5, the inverse of the Hessian matrix exists. Substituting (22) back into (19) directly yields the desired result. \square

Note that $\bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}_{i,k})$ denotes the surrogate gradient at UL through replacing $\mathbf{y}_{i,k}^*(\mathbf{x})$ in (18) $\mathbf{y}_{i,k}$ with the local loss function, i.e.,

$$\begin{aligned} \bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}) &= \nabla_{\mathbf{x}_i} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}) \\ &\quad - \sum_{k=1}^m \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}). \end{aligned} \quad (23)$$

It has been shown in [31, Lemma 3.2] and [33, Lemma 1] that, when the gradient estimator is constructed in a certain way, we can have

$$\|\bar{\nabla} f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^{r+1}) - \mathbb{E} \mathbf{h}_{f,i}\| \triangleq b_{r,i} \leq \mathcal{O}(b^{m_b}), \quad \forall i, \quad (24)$$

due to the independence among the LL tasks, where $0 < b < 1$ and $m_b \geq r$ denotes the mini-batch size. Therefore, we only need to choose $m_b = \mathcal{O}(\log(nT))$ to obtain $b_r^2 \leq \mathcal{O}(1/\sqrt{nT})$, where $b_r \triangleq \sum_{i=1}^n b_{r,i}$.

Lemma 2.

Under A1 and A5, $\mathbf{y}_k^(\mathbf{x})$ is Lipschitz continuous, namely*

$$\|\bar{\mathbf{y}}_k^*(\mathbf{x}) - \bar{\mathbf{y}}_k^*(\mathbf{x}')\| \leq L_y \|\mathbf{x} - \mathbf{x}'\|, \quad \forall k, \quad (25)$$

where $L_y \triangleq \frac{C_{xy}}{\mu_g}$, and $\nabla \mathbf{y}_k^(\mathbf{x})$ is also Lipschitz continuous, namely*

$$\|\nabla \bar{\mathbf{y}}_k^*(\mathbf{x}) - \nabla \bar{\mathbf{y}}_k^*(\mathbf{x}')\| \leq L_{xy} \|\mathbf{x} - \mathbf{x}'\|, \quad \forall k, \quad (26)$$

where $L_{xy} \triangleq \frac{\sqrt{2}L_{g,2}}{\mu_g}(1 + L_y + C_{xy}(1 + L_y)\mu_g^{-1})$.

Proof. First Part. According to (20) and (22), we have

$$\|\nabla_{\mathbf{x}_i} \mathbf{y}_{i,k}^*(\mathbf{x})\| \leq \|\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x}))^{-1} \nabla_{\mathbf{x}_i \mathbf{y}_k} g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x}))^T\| \leq \frac{C_{xy}}{\mu_g}, \quad (27)$$

where we use A1 and A5. Therefore, we obtain $\|\nabla_{\mathbf{x}} \bar{\mathbf{y}}_k^*(\mathbf{x})\| \leq L_{g,2}/\mu_g$ and

$$\|\bar{\mathbf{y}}_k^*(\mathbf{x}) - \bar{\mathbf{y}}_k^*(\mathbf{x}')\| \leq \frac{C_{xy}}{\mu_g} \|\mathbf{x} - \mathbf{x}'\|. \quad (28)$$

Second Part. Next, we can have

$$\begin{aligned} &\|\nabla_{\mathbf{x}} \bar{\mathbf{y}}_k^*(\mathbf{x}) - \nabla_{\mathbf{x}} \bar{\mathbf{y}}_k^*(\mathbf{x}')\| \\ &= \|\nabla_{\mathbf{x} \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) [\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x}))]^{-1} - \nabla_{\mathbf{x} \mathbf{y}_k}^2 g_k(\mathbf{x}', \mathbf{y}_k^*(\mathbf{x}')) [\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}', \mathbf{y}_k^*(\mathbf{x}'))]^{-1}\| \quad (29) \\ &\leq \frac{1}{\mu_g} \|\nabla_{\mathbf{x} \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) - \nabla_{\mathbf{x} \mathbf{y}_k}^2 g_k(\mathbf{x}', \mathbf{y}_k^*(\mathbf{x}'))\| + \frac{C_{xy}}{\mu_g^2} \|\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) - \nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}', \mathbf{y}_k^*(\mathbf{x}'))\| \end{aligned}$$

$$\stackrel{(a)}{\leq} \frac{\sqrt{2}L_{g,2}}{\mu_g} (\|\mathbf{x} - \mathbf{x}'\| + \|\mathbf{y}_k^*(\mathbf{x}) - \mathbf{y}_k^*(\mathbf{x}')\|) + \frac{\sqrt{2}C_{xy}L_{g,2}}{\mu_g^2} (\|\mathbf{x} - \mathbf{x}'\| + \|\mathbf{y}_k^*(\mathbf{x}) - \mathbf{y}_k^*(\mathbf{x}')\|) \quad (30)$$

$$\leq \frac{\sqrt{2}L_{g,2}}{\mu_g} \left(1 + L_y + \frac{C_{xy}(1 + L_y)}{\mu_g} \right) \|\mathbf{x} - \mathbf{x}'\| \quad (31)$$

where in (a) we use

$$\|\nabla_{\mathbf{x} \mathbf{y}_k}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}_i)) - \nabla_{\mathbf{x} \mathbf{y}_k}^2 g_{i,k}(\mathbf{x}'_i, \mathbf{y}_{i,k}^*(\mathbf{x}'_i))\| \leq L_{g,2} (\|\mathbf{x}_i - \mathbf{x}'_i\| + \|\mathbf{y}_{i,k}^*(\mathbf{x}_i) - \mathbf{y}_{i,k}^*(\mathbf{x}'_i)\|) \quad (32a)$$

$$\|\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x}_i)) - \nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_{i,k}(\mathbf{x}'_i, \mathbf{y}_{i,k}^*(\mathbf{x}'_i))\| \leq L_{g,2} (\|\mathbf{x}_i - \mathbf{x}'_i\| + \|\mathbf{y}_{i,k}^*(\mathbf{x}_i) - \mathbf{y}_{i,k}^*(\mathbf{x}'_i)\|) \quad (32b)$$

by directly applying A1 and (25). \square

Lemma 3.

Under A1 and A5, there exists a constant $L_{f,y}$ such that function $\|\bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}) - \bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}'_{i,k})\|$ is upper bounded by the sum of $\|\mathbf{y}_{i,k} - \mathbf{y}'_{i,k}\|$, namely

$$\|\bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}_{1,k}, \dots, \mathbf{y}_{m,k}) - \bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}'_{1,k}, \dots, \mathbf{y}'_{m,k})\| \leq L_{f,y} \sum_{k=1}^m \|\mathbf{y}_{i,k} - \mathbf{y}'_{i,k}\|, \quad (33)$$

and there also exists a constant $L_{f,x}$ such that $\|\bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}) - \bar{\nabla} f_i(\mathbf{x}'_i, \mathbf{y}_{i,k})\|$ is upper bounded by $\|\mathbf{x} - \mathbf{x}'\|$, namely

$$\|\bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}_{1,k}, \dots, \mathbf{y}_{m,k}) - \bar{\nabla} f_i(\mathbf{x}'_i, \mathbf{y}_{1,k}, \dots, \mathbf{y}_{m,k})\| \leq mL_{f,x} \|\mathbf{x}_i - \mathbf{x}'_i\|, \quad (34)$$

where $L_{f,x}, L_{f,y}$ are only dependent on the parameters defined in A1 and A5.

Proof. First Part. Suppose assumptions A1 and A5 hold. From (23) and [31, Lemma 2.2.], we have

$$\begin{aligned} & \|\nabla_{\mathbf{x}_i, \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) \left[\nabla_{\mathbf{y}_{i,k}, \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}) \\ & \quad - \nabla_{\mathbf{x}_i, \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}'_{i,k}) \left[\nabla_{\mathbf{y}_{i,k}, \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}'_{i,k}) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}'_{i,k})\| \end{aligned} \quad (35)$$

$$\leq \left(\frac{L_{f,1} C_{xy}}{\mu_g} + L_{f,0} \left(\frac{L_{g,2}}{\mu_g} + \frac{L_{g,2} C_{xy}}{\mu_g^2} \right) \right) \|\mathbf{y}_{i,k} - \mathbf{y}'_{i,k}\|. \quad (36)$$

Based on the block-wise gradient Lipschitz continuity, we have

$$\|\nabla_{\mathbf{x}_i} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}) - \nabla_{\mathbf{x}_i} f_i(\mathbf{x}_i, \mathbf{y}'_{i,k})\| \leq \sum_{k=1}^m L_{f,1} \|\mathbf{y}_{i,k} - \mathbf{y}'_{i,k}\|. \quad (37)$$

Combing (23) with (36) and (37) gives the desired result immediately.

Second Part. Similarly, under A1 and A5, we can also get

$$\begin{aligned} & \|\bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}) - \bar{\nabla} f_i(\mathbf{x}'_i, \mathbf{y}_{i,k})\| \\ & \leq L_{f,1} \|\mathbf{x}_i - \mathbf{x}'_i\| + \frac{m}{\mu_g} \left(L_{f,1} C_{xy} + L_{f,0} \left(L_{g,2} + \frac{L_{g,2} C_{xy}}{\mu_g} \right) \right) \|\mathbf{x}_i - \mathbf{x}'_i\| \end{aligned} \quad (38)$$

From (2a), we can get the desired result by requiring $L_{f,x} \geq (L_{f,1} + m\mu_g^{-1}(L_{f,1} C_{xy} + L_{f,0}(L_{g,2} + L_{g,2} C_{xy} \mu_g^{-1}))) / m$. \square

B Convergence Analysis

We now present the proofs, related results, and technical details that establish the lemmas and theorems of our convergence analysis.

We first show that the difference between two successive iterates in Lemma 4 is upper bounded on the order of $1/T$. Then, we begin to quantify the process where one step of the variable updates would make: 1) Lemma 5 measures the closeness from the iterates $\bar{\mathbf{y}}^r$ to its optimal counterpart given the UL variable fixed; 2) Lemma 6 essentially gives the upper bound of the consensus violations from both UL and LL sides or the maximum ascent achieved by the dual update. As the byproducts of Lemma 6, recursion functions in terms of the successive differences of the UL and LL variables are derived in Lemma 7 and Lemma 8 based on the optimality conditions of the UL and LL optimization problems respectively, which serve as the critical role of establishing the potential functions that can evaluate the process of the sequence generated by SLAM to KKT points. Finally, all the above properties are used in the proof of Theorem 1.

B.1 Upper Bounds of Successive Primal Variables

Lemma 4.

Under A1, A3, A4, suppose that the iterates $\{\mathbf{x}^r, \mathbf{y}_k^r, \forall k, r\}$ are generated by (5a) and (5b) and the step sizes are chosen to be $\alpha = \alpha_0 \sqrt{T}$ and $\beta = \beta_0 \sqrt{T}$. Then, there exist constants D_x, D_y such that, when

$$\alpha_0 \geq 2\rho\sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sqrt{C_x}, \quad (39a)$$

$$\beta_0 \geq \max\{1/\mu_g, \tau_g\}, \quad (39b)$$

$$\tau_g \geq \frac{\rho\sigma_{\max}(\mathbf{A}^T \mathbf{A})}{1 - 1/\sqrt{T}}, \quad (39c)$$

$$\tau_f \geq \rho\sigma_{\max}(\mathbf{A}^T \mathbf{A}), \quad (39d)$$

the following inequalities hold

$$\mathbb{E}\|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 \leq \frac{D_x}{\alpha^2}, \quad \forall r, \quad (40a)$$

$$\frac{1}{T} \sum_{r=1}^T \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \leq \frac{D_y}{\beta^2}, \quad \forall k, r, \quad (40b)$$

where the constants are given by

$$C_x = \frac{2(1 + \theta^{-1})}{(1 - (1 + \theta) \left(1 - \frac{\rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}{\tau_f}\right)^2)}, \quad (41a)$$

$$D_x = \frac{4C_x \sigma_f^2}{n} + 4(L_{f,0} + \sigma_f^2) + 2C_x \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A}) B_x, \quad (41b)$$

$$D_y = \frac{\sigma_g^2}{n} + \frac{(D_x \mu_g^{-2} + \mu_g \rho \|\mathbf{y}^1\|^2 \alpha_0) \beta_0^2}{\alpha_0^2} + \frac{\sigma_g^2}{n \mu_g \beta_0}, \quad (41c)$$

and

$$0 < \theta < \frac{1}{(1 - \frac{\rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}{\tau_f})^2} - 1, \quad (42)$$

$$B_x \triangleq \max \left\{ \|\mathbf{x}^1\|^2, \frac{\sigma_f^2(2 + 2/\theta) + n(L_{f,0} + \sigma_f^2)(1 - (1 + \theta)(1 - \frac{\rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}{\tau_f})^2)}{n(2 + 2/\theta)\rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A})} \right\}. \quad (43)$$

Proof. First Part. From (6b) and A3, we know that

$$\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r = \bar{\mathbf{x}}^r - \bar{\mathbf{x}}^{r-1} + \frac{1}{\alpha}(\bar{\mathbf{h}}_f^r - \bar{\mathbf{h}}_f^{r-1}), \quad (44)$$

where

$$\bar{\mathbf{x}} \triangleq \frac{1}{n} \mathbb{1}^T \mathbf{x}, \quad \bar{\mathbf{h}} \triangleq \frac{1}{n} \mathbb{1}^T \mathbf{h}. \quad (45)$$

We then derive, from (5b) and (6b), that

$$\begin{aligned} & \|\mathbf{x}^{r+1} - \mathbf{x}^r - (\mathbb{1} \bar{\mathbf{x}}^{r+1} - \mathbb{1} \bar{\mathbf{x}}^r)\| \\ & \stackrel{(6b),(44)}{\leq} \left\| \left(\mathbf{I} - (1 + \gamma^{-1}) \frac{\rho \mathbf{A}^T \mathbf{A}}{\tau} \right) (\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r) - \left(\mathbf{I} - \frac{\rho \mathbf{A}^T \mathbf{A}}{\tau_f} \right) (\mathbf{x}^{r-1} - \mathbb{1} \bar{\mathbf{x}}^{r-1}) \right\| \\ & \quad + \frac{1}{\alpha} \|\mathbf{h}_f^r - \mathbb{1} \bar{\mathbf{h}}_f^r - (\mathbf{h}_f^{r-1} - \mathbb{1} \bar{\mathbf{h}}_f^{r-1})\| \\ & \leq \left\| \left(\mathbf{I} - \frac{\rho \mathbf{A}^T \mathbf{A}}{\tau_f} \right) (\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r - (\mathbf{x}^{r-1} - \mathbb{1} \bar{\mathbf{x}}^{r-1})) \right\| + \frac{1}{\alpha} \|\rho \mathbf{A}^T \mathbf{A} \mathbf{x}^r\| \\ & \quad + \frac{1}{\alpha} \|\mathbf{h}_f^r - \mathbb{1} \bar{\mathbf{h}}_f^r - (\mathbf{h}_f^{r-1} - \mathbb{1} \bar{\mathbf{h}}_f^{r-1})\| \\ & \stackrel{(a)}{\leq} \left(1 - \frac{\rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}{\tau_f} \right) \|\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r - (\mathbf{x}^{r-1} - \mathbb{1} \bar{\mathbf{x}}^{r-1})\| + \frac{1}{\alpha} \|\rho \mathbf{A}^T \mathbf{A} \mathbf{x}^r\| \\ & \quad + \frac{1}{\alpha} \|\mathbf{h}_f^r - \mathbb{1} \bar{\mathbf{h}}_f^r - (\mathbf{h}_f^{r-1} - \mathbb{1} \bar{\mathbf{h}}_f^{r-1})\|, \end{aligned} \quad (46)$$

where (a) holds because $\mathbb{1}^T (\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r) = 0, \forall r$, and $\tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})$ denotes the minimum nonzero eigenvalue of matrix $\mathbf{A}^T \mathbf{A}$. From the definition of \mathbf{h}_f^r and A3, we have

$$\mathbb{E} \|\mathbf{h}_f^r - \mathbb{1} \bar{\mathbf{h}}_f^r\|^2 = \sum_{i=1}^n \mathbb{E} \|\mathbf{h}_{f,i}^r - \mathbb{E} \bar{\mathbf{h}}_f^r + \mathbb{E} \bar{\mathbf{h}}_f^r - \bar{\mathbf{h}}_f^r\|^2 \leq \frac{2\sigma_f^2}{n}. \quad (47)$$

Next, we use mathematical induction to prove the boundedness of the variable \mathbf{x}^r . When $r = 1$, the size of \mathbf{x}^1 is bounded by a constant, i.e., $\|\mathbf{x}^1\|^2$. Applying Young's inequality, we obtain that, for $\theta > 0$,

$$\begin{aligned} & \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r - (\mathbb{1} \bar{\mathbf{x}}^{r+1} - \mathbb{1} \bar{\mathbf{x}}^r)\|^2 \\ & \leq (1 + \theta) \left(1 - \frac{\rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}{\tau_f} \right)^2 \mathbb{E} \|\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r - (\mathbf{x}^{r-1} - \mathbb{1} \bar{\mathbf{x}}^{r-1})\|^2 \\ & \quad + \left(1 + \frac{1}{\theta} \right) \frac{2}{\alpha^2} \left(\frac{2\sigma_f^2}{n} + \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A}) \|\mathbf{x}^r\|^2 \right). \end{aligned} \quad (48)$$

Define

$$\eta_x \triangleq (1 + \theta) \left(1 - \frac{\rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}{\tau_f} \right)^2, \quad \text{and let } \tau_f \geq \rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A}). \quad (49)$$

When

$$\theta < \frac{1}{\left(1 - \frac{\rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}{\tau_f}\right)^2} - 1, \quad (50)$$

it is obvious that $\eta_x < 1$. Then, we obtain

$$\begin{aligned} & \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r - (\mathbb{1} \bar{\mathbf{x}}^{r+1} - \mathbb{1} \bar{\mathbf{x}}^r)\|^2 \\ & \leq \eta_x \mathbb{E} \|\mathbf{x}^r - \mathbf{x}^{r-1} - (\mathbb{1} \bar{\mathbf{x}}^r - \mathbb{1} \bar{\mathbf{x}}^{r-1})\|^2 + \left(1 + \frac{1}{\theta}\right) \frac{2}{\alpha^2} \left(\frac{2\sigma_f^2}{n} + \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A}) B_1\right) \end{aligned} \quad (51)$$

$$\leq \frac{\left(1 + \frac{1}{\theta}\right) \frac{2}{\alpha^2}}{1 - \eta_x} \left(\frac{2\sigma_f^2}{n} + \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A}) B_1\right) = \frac{C_x \left(\frac{2\sigma_f^2}{n} + \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A}) B_1\right)}{\alpha^2}, \quad (52)$$

where we assume that $\|\mathbf{x}^r\|^2 \leq B_1$ and define $C_x \triangleq \frac{1 + \frac{1}{\theta}}{(1 - \eta_x)}$. Further, from (5b) and (4d), we have

$$\bar{\mathbf{x}}^{r+1} = \bar{\mathbf{x}}^r - \frac{1}{\alpha} \bar{\mathbf{h}}_f^r \quad (53)$$

due to $\mathbb{1}^T \mathbf{A}^T \mathbf{A} = 0$, which yields

$$\|\mathbb{1} \bar{\mathbf{x}}^{r+1} - \mathbb{1} \bar{\mathbf{x}}^r\| = \frac{1}{\alpha} \|\mathbb{1} \bar{\mathbf{h}}_f^r\|, \quad (54)$$

and thus we obtain

$$\|\mathbb{1} \bar{\mathbf{x}}^{r+1} - \mathbb{1} \bar{\mathbf{x}}^r\|^2 \leq \frac{1}{\alpha^2} (2\|\bar{\mathbf{h}}_f^r - \mathbb{E} \bar{\mathbf{h}}_f^r\|^2 + 2\|\mathbb{E} \bar{\mathbf{h}}_f^r\|^2) \leq 2 \frac{\sigma_f^2 + L_{f,0}^2}{n \alpha^2} \quad (55)$$

under A1 and A3.

Therefore, combing (52) and (55) renders

$$\begin{aligned} \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 &= \mathbb{E} \|\mathbf{x}^{r+1} - \mathbb{1} \bar{\mathbf{x}}^{r+1} + \mathbb{1} \bar{\mathbf{x}}^{r+1} - \mathbb{1} \bar{\mathbf{x}}^r + \mathbb{1} \bar{\mathbf{x}}^r - \mathbf{x}^r\|^2 \\ &\leq 2\mathbb{E} \|\mathbf{x}^{r+1} - \mathbb{1} \bar{\mathbf{x}}^{r+1} - (\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r)\|^2 + 2\mathbb{E} \|\mathbb{1} \bar{\mathbf{x}}^{r+1} - \mathbb{1} \bar{\mathbf{x}}^r\|^2 \end{aligned} \quad (56)$$

$$\leq 2 \frac{C_x \left(\frac{2\sigma_f^2}{n} + \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A}) B_1\right) + 2(L_{f,0}^2 + \sigma_f^2)/n}{\alpha^2} \quad (57)$$

$$= 2 \frac{2C_x \sigma_f^2 + 2(L_{f,0} + \sigma_f^2)}{n \alpha^2} + \frac{2C_x \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A}) B_1}{\alpha^2} \sim \mathcal{O}\left(\frac{1}{\alpha^2}\right), \quad (58)$$

and thus we have

$$\|\mathbf{x}^{r+1}\|^2 \leq r \sum_{r=1}^T \|\mathbf{x}^r\|^2 \leq T \sum_{r=1}^T \|\mathbf{x}^r\|^2 \leq 2 \frac{2C_x \sigma_f^2 + 2(L_{f,0} + \sigma_f^2)}{n \alpha_0^2} + \frac{2C_x \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A}) B_1}{\alpha_0^2} \quad (59)$$

where we choose $\alpha = \alpha_0 \sqrt{T}$. To show $\|\mathbf{x}^{r+1}\|^2 \leq B_1$, we only need

$$\frac{2C_x \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A})}{\alpha_0^2} \leq \frac{1}{2} \quad \text{and} \quad \frac{4C_x \sigma_f^2 + 4(L_{f,0} + \sigma_f^2)}{n \alpha_0^2} \leq B_1, \quad (60)$$

which means that

$$\alpha_0 \geq 2\rho \sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sqrt{C_x} \quad \text{and} \quad B_1 \geq \frac{C_x \sigma_f^2 + (L_{f,0} + \sigma_f^2)}{n \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A}) C_x}. \quad (61)$$

Therefore, combining the case where $r = 1$, we conclude that, when

$$\alpha_0 \geq 2\rho \sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sqrt{\frac{2 + \frac{2}{\theta}}{\left(1 - (1 + \theta) \left(1 - \frac{\rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}{\tau}\right)^2\right)}}, \quad (62)$$

then

$$\|\mathbf{x}^{r+1}\|^2 \leq \max \left\{ \|\mathbf{x}^1\|^2, \frac{\sigma_f^2 + \frac{(L_{f,0} + \sigma_f^2)(1 - (1 + \theta)(1 - \frac{\rho \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}{\tau})^2)}{2 + 2/\theta}}{n \rho^2 \sigma_{\max}^2(\mathbf{A}^T \mathbf{A})} \right\} \triangleq B_x. \quad (63)$$

This directly implies

$$\mathbb{E}\|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 \leq \frac{\overbrace{4(C_x\sigma_f^2 + L_{f,0} + \sigma_f^2) + 2C_x\rho^2\sigma_{\max}^2(\mathbf{A}^T\mathbf{A})B_x}^{\triangleq D_x}}{\alpha^2} \sim \mathcal{O}\left(\frac{1}{\alpha^2}\right). \quad (64)$$

Second Part. From (6a), we have

$$\mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] = \mathbb{E}\left[\left(\mathbf{I} - \frac{\rho}{\tau_g}\mathbf{A}^T\mathbf{A}\right)\mathbf{y}_k^r - \left(\mathbf{I} - \frac{\rho}{\tau_g}\mathbf{A}^T\mathbf{A}\right)\mathbf{y}_k^{r-1} - \frac{\rho}{\beta}\mathbf{A}^T\mathbf{A}\mathbf{y}_k^r - \frac{1}{\beta}(\mathbf{h}_{g,k}^r - \mathbf{h}_{g,k}^{r-1})\right]. \quad (65)$$

Multiplying $\mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r]$ on both sides of (65) yields

$$\begin{aligned} \|\mathbb{E}\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 &= \langle \mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}], \mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] \rangle - \frac{\rho}{\tau_g} \langle \mathbf{A}^T\mathbf{A}(\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]), \mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] \rangle \\ &\quad - \frac{\rho}{\beta} \langle \mathbf{A}^T\mathbf{A}\mathbb{E}[\mathbf{y}_k^r], \mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] \rangle - \frac{1}{\beta} \langle g_k(\mathbf{x}^r, \mathbf{y}_k^r) - g_k(\mathbf{x}^{r-1}, \mathbf{y}_k^{r-1}), \mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] \rangle \\ &\leq \frac{\|\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|^2}{2} + \frac{\|\mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r]\|^2}{2} - \frac{\|\mathbb{E}[\mathbf{v}_k^{r+1}]\|^2}{2} \\ &\quad - \frac{\rho}{\tau_g} \left(\frac{\|\mathbf{A}\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|^2}{2} + \frac{\|\mathbf{A}\mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r]\|^2}{2} - \frac{\|\mathbf{A}\mathbb{E}[\mathbf{v}_k^{r+1}]\|^2}{2} \right) \\ &\quad - \frac{\rho}{\beta} \left(\frac{\|\mathbf{A}\mathbb{E}[\mathbf{y}_k^{r+1}]\|^2}{2} - \frac{\|\mathbf{A}\mathbb{E}[\mathbf{y}_k^r]\|^2}{2} - \frac{\|\mathbf{A}\mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r]\|^2}{2} \right) \\ &\quad + \frac{\mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2}{\mu_g\beta} + \frac{\mu_g\|\mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r]\|^2}{4\beta} - \frac{\mu_g}{\beta}\|\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|^2 \\ &\quad + \frac{\sqrt{T}}{2\beta^2}\mathbb{E}\|\mathbf{y}_k^r - \mathbf{y}_k^{r-1}\|^2 + \frac{1}{2\sqrt{T}}\|\mathbb{E}[\mathbf{v}_k^{r+1}]\|^2, \end{aligned}$$

where

$$\begin{aligned} &-\frac{1}{\beta} \langle g(\mathbf{x}^{r-1}, \mathbf{y}_k^r) - g(\mathbf{x}^{r-1}, \mathbf{y}_k^{r-1}), \mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1} + \mathbf{v}_k^{r+1}] \rangle \\ &\stackrel{(a)}{\leq} -\frac{\mu_g}{\beta}\|\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|^2 + \frac{\sqrt{T}}{2\beta^2}\mathbb{E}\|\mathbf{y}_k^r - \mathbf{y}_k^{r-1}\|^2 + \frac{1}{2\sqrt{T}}\|\mathbb{E}\mathbf{v}_k^{r+1}\|^2. \end{aligned} \quad (66)$$

Here (a) follows the strong convexity of function g_k , Young's inequality, and (17).

Note that

$$\frac{\rho}{\tau_g} - \frac{\rho}{\beta} \geq 0, \quad (67)$$

when $\beta \geq \tau_g$, i.e. $\beta_0\sqrt{T} \geq \tau_g$ or $\beta_0 \geq \tau_g$. Moreover, we know that

$$\begin{aligned} &\mathbb{E}_{\zeta^r} [\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 | \mathcal{F}^r] \\ &= \|\mathbb{E}_{\zeta^r}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] | \mathcal{F}^r\|^2 + \mathbb{E}_{\xi^r} [\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r - \mathbb{E}_{\xi^r}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r]\|^2 | \mathcal{F}^r] \end{aligned} \quad (68)$$

$$= \|\mathbb{E}_{\zeta^r}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] | \mathcal{F}^r\|^2 + \mathbb{E}_{\xi^r} \left[\left\| \frac{1}{\alpha} (\mathbb{E}\mathbf{h}_{g,k}^r - \mathbf{h}_{g,k}^r) \right\|^2 | \mathcal{F}^r \right] \quad (69)$$

$$\leq \|\mathbb{E}_{\zeta^r}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] | \mathcal{F}^r\|^2 + \frac{\sigma_g^2}{n\beta^2}. \quad (70)$$

Therefore, we obtain

$$\begin{aligned} &\left(\frac{1}{2} - \frac{\mu_g}{4\beta}\right) \|\mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r]\|^2 \\ &\leq \left(\frac{1}{2} + \frac{\sqrt{T}}{2\beta^2} - \frac{\mu_g}{\beta}\right) \|\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|^2 - \left(\frac{1}{2} - \frac{1}{2\sqrt{T}} - \frac{\rho\sigma_{\max}(\mathbf{A}^T\mathbf{A})}{2\tau_g}\right) \|\mathbb{E}[\mathbf{v}_k^{r+1}]\|^2 \\ &\quad + \frac{\mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2}{\mu_g\beta} - \frac{\rho}{\beta} \left(\frac{\|\mathbf{A}\mathbb{E}[\mathbf{y}_k^{r+1}]\|^2}{2} - \frac{\|\mathbf{A}\mathbb{E}[\mathbf{y}_k^r]\|^2}{2} \right) + \frac{\sigma_g^2\sqrt{T}}{2n\beta^4}. \end{aligned}$$

When $\sqrt{T}/\beta \leq \mu_g$, i.e., $\beta_0 \geq 1/\mu_g$, we have $\frac{\sqrt{T}}{2\beta^2} - \frac{\mu_g}{\beta} \leq -\frac{\mu_g}{2\beta}$; further, we need $\tau_g \geq \frac{\rho\sigma_{\max}(\mathbf{A}^T\mathbf{A})}{1-1/\sqrt{T}}$. Then, we obtain

$$\|\mathbb{E}\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \leq \frac{1 - \frac{2\mu_g}{\beta}}{1 - \frac{\mu_g}{\beta}} \|\mathbb{E}\mathbf{y}_k^r - \mathbf{y}_k^{r-1}\|^2 + \frac{\mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2}{\mu_g\beta} \quad (71)$$

$$- \frac{\rho}{\beta} \left(\frac{\|\mathbf{A}\mathbb{E}\mathbf{y}_k^{r+1}\|^2}{2} - \frac{\|\mathbf{A}\mathbb{E}\mathbf{y}_k^r\|^2}{2} \right) + \frac{\sigma_g^2\sqrt{T}}{2n\beta^4}. \quad (72)$$

Applying the telescoping sum on both sides of (72) yields

$$\frac{1}{T} \sum_{r=1}^T \|\mathbb{E}\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \leq \left(1 - \frac{\mu_g}{\beta - \mu_g}\right) \frac{1}{T} \sum_{r=1}^T \|\mathbb{E}\mathbf{y}_k^r - \mathbf{y}_k^{r-1}\|^2 + \frac{D_x}{\mu_g\beta\alpha^2} + \frac{\rho\|\mathbf{y}_k^1\|^2}{T\beta} + \frac{\sigma_g^2\sqrt{T}}{n\beta^4}, \quad (73)$$

and thus there exists a contraction property for the sum of $\|\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|^2$. Define $S \triangleq \frac{1}{T} \sum_{r=1}^T \|\mathbb{E}\mathbf{y}_k^r - \mathbf{y}_k^{r-1}\|^2$. When

$$\left(1 - \frac{\mu_g}{\beta - \mu_g}\right) S + \frac{D_x}{\mu_g\beta\alpha^2} + \frac{\rho\|\mathbf{y}_k^1\|^2}{T\beta} + \frac{\sigma_g^2\sqrt{T}}{n\beta^4} \leq S, \quad (74)$$

namely

$$S \geq \frac{\beta - \mu_g}{\mu_g} \left(\frac{D_x}{\mu_g\beta\alpha^2} + \frac{\rho\|\mathbf{y}_k^1\|^2}{T\beta} + \frac{\sigma_g^2\sqrt{T}}{n\beta^4} \right), \quad (75)$$

then the sum of $\|\mathbb{E}[\mathbf{y}^r - \mathbf{y}^{r-1}]\|^2$ is upper bounded.

Choosing

$$S = \frac{\beta}{\mu_g} \left(\frac{D_x}{\mu_g\beta\alpha^2} + \frac{\rho\|\mathbf{y}_k^1\|^2}{T\beta} \right) = \frac{D_x}{\mu_g^2\alpha^2} + \frac{\mu_g\rho\|\mathbf{y}_k^1\|^2}{T} + \frac{\sigma_g^2\sqrt{T}}{n\mu_g\beta^3}, \quad (76)$$

and applying (70) yield

$$\frac{1}{T} \sum_{r=1}^T \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \leq \frac{D_x}{\mu_g^2\alpha^2} + \frac{\mu_g\rho\|\mathbf{y}_k^1\|^2\alpha_0}{\alpha^2} + \frac{\sigma_g^2}{n\mu_g\beta_0\beta^2} + \frac{\sigma_g^2}{n\beta^2} \quad (77)$$

$$= \frac{D_x/\mu_g^2 + \mu_g\rho\|\mathbf{y}_k^1\|^2\alpha_0}{\alpha^2} + \left(\frac{1}{\mu_g\beta_0} + 1 \right) \frac{\sigma_g^2}{n\beta^2}. \quad (78)$$

Combining this with (78) renders

$$\frac{1}{T} \sum_{r=1}^T \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \leq \frac{D_y}{\beta^2}, \quad \forall k, \quad (79)$$

where $D_y \triangleq (D_x/\mu_g^2 + \mu_g\rho\|\mathbf{y}_k^1\|^2\alpha_0)\beta_0^2/\alpha_0^2 + (\frac{1}{\mu_g\beta_0} + 1)\sigma_g^2/n$. \square

B.2 Contraction of the LL iterates

Now, we will show the contraction property of the recurrence in the LL optimization process. The proof is adapted from [43, Lemma 3], where the main difference is that we evaluate the iterates in the consensus space. To be more precise, we establish the following result.

Lemma 5.

Under A1, A3, A5. Suppose that sequence $\{\mathbf{x}^r, \mathbf{y}_k^r, \boldsymbol{\lambda}_k^r, \boldsymbol{\omega}_k^r, \forall k, r\}$ is generated by SLAM. Then, when $\beta > 2(\mu_g + L_{g,1})^{-1}$, we have

$$\mathbb{E}\|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \leq \left(1 - \frac{\rho_g}{\beta}\right) \mathbb{E}\|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 + \frac{\sigma_g^2}{n\beta^2}, \quad (80)$$

and there exist constants $\theta', \vartheta > 0$ such that

$$\begin{aligned} \mathbb{E}\|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1})\|^2 &\leq \left(1 + \theta' + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \mathbb{E}\|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \\ &\quad + \left(L_y^2 + \frac{L_y^2}{4\theta'} + \frac{L_{xy}}{4\vartheta}\right) \mathbb{E}\|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 + \left(\frac{L_{xy}}{4\vartheta} + L_y^2\right) \frac{\sigma_f^2}{n\alpha^2}. \end{aligned} \quad (81)$$

Proof. First, we expand the error term in the lower level optimization problem as

$$\begin{aligned} \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1})\|^2 &= \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 + \|\bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1}) - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \\ &\quad + 2 \langle \bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), \bar{\mathbf{y}}_k^*(\mathbf{x}^r) - \bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1}) \rangle. \end{aligned} \quad (82)$$

Then, we will give the upper bound for each term in (82).

$$\begin{aligned} &\mathbb{E} [\|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 | \mathcal{F}^r] \\ \stackrel{(5a)}{=} &\mathbb{E} \left[\left\| \bar{\mathbf{y}}_k^{r+1} - \frac{1}{\beta} \bar{\mathbf{h}}_{g,k}^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r) \right\|^2 | \mathcal{F}^r \right] \end{aligned} \quad (83)$$

$$= \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 - \frac{2}{\beta} \langle \bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), \mathbb{E}[\bar{\mathbf{h}}_{g,k}^r] | \mathcal{F}^r \rangle + \frac{1}{\beta^2} \mathbb{E}[\|\bar{\mathbf{h}}_{g,k}^r\|^2 | \mathcal{F}^r] \quad (84)$$

$$\stackrel{(a)}{\leq} \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 - \frac{2}{\beta} \langle \bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), g_k(\mathbf{x}^r, \bar{\mathbf{y}}_k^r) \rangle + \frac{1}{\beta^2} \|\nabla g_k(\mathbf{x}^r, \bar{\mathbf{y}}_k^r)\|^2 + \frac{\sigma_g^2}{n\beta^2} \quad (85)$$

$$\stackrel{(b)}{\leq} \left(1 - \frac{\rho_g}{\beta}\right) \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 + \frac{\sigma_g^2}{n\beta^2} \quad (86)$$

where in (a) we use A3, (b) follows the μ_g -strong convexity of function $g_k(\mathbf{x}, \mathbf{y}_k)$ [49, Theorem 2.1.11], i.e.,

$$\begin{aligned} - \langle \bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), g_k(\mathbf{x}^r, \bar{\mathbf{y}}_k) \rangle &\leq - \frac{\mu_g L_{g,1}}{\mu_g + L_{g,1}} \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \\ &\quad - \frac{1}{\mu_g + L_{g,1}} \|\nabla g_k(\mathbf{x}^r, \bar{\mathbf{y}}_k) - \nabla g_k(\mathbf{x}^r, \bar{\mathbf{y}}_k^*(\mathbf{x}^r))\|^2 \end{aligned} \quad (87)$$

and we choose $\beta \geq \frac{2}{\mu_g + L_{g,1}}$ with $\rho_g \triangleq \frac{2\mu_g L_{g,1}}{\mu_g + L_{g,1}}$.

Taking the full expectation over \mathcal{F}^r , we have

$$\mathbb{E} \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \leq \left(1 - \frac{\rho_g}{\beta}\right) \mathbb{E} \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 + \frac{\sigma_g^2}{n\beta^2}. \quad (88)$$

Next, we split the cross term into two parts

$$\begin{aligned} \langle \bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), \bar{\mathbf{y}}_k^*(\mathbf{x}^r) - \bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1}) \rangle &= - \langle \bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), \nabla \bar{\mathbf{y}}_k^*(\mathbf{x}^r)(\mathbf{x}^{r+1} - \mathbf{x}^r) \rangle \\ &\quad - \langle \bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), \bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1}) - \bar{\mathbf{y}}_k^*(\mathbf{x}^r) - \nabla \bar{\mathbf{y}}_k^*(\mathbf{x}^r)(\mathbf{x}^{r+1} - \mathbf{x}^r) \rangle. \end{aligned} \quad (89)$$

The first part can be upper bounded by

$$\begin{aligned} &-\mathbb{E} [\langle \bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), \nabla \bar{\mathbf{y}}_k^*(\mathbf{x}^r)(\mathbf{x}^{r+1} - \mathbf{x}^r) \rangle] \\ &= -\mathbb{E} [\langle \bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), \nabla \bar{\mathbf{y}}_k^*(\mathbf{x}^r) \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] | \mathcal{F}^r \rangle] \end{aligned} \quad (90)$$

$$\stackrel{(25)}{\leq} L_y \mathbb{E} \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\| \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\| \quad (91)$$

$$\stackrel{(a)}{\leq} \theta' \mathbb{E} \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 + \frac{L_y^2}{4\theta'} \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2. \quad (92)$$

where in (a) we apply Young's inequality.

The second part can be upper bounded by

$$\begin{aligned} &-\langle \bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r), \bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1}) - \bar{\mathbf{y}}_k^*(\mathbf{x}^r) - \nabla \bar{\mathbf{y}}_k^*(\mathbf{x}^r)(\mathbf{x}^{r+1} - \mathbf{x}^r) \rangle \\ &\leq \mathbb{E} \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\| \|\bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1}) - \bar{\mathbf{y}}_k^*(\mathbf{x}^r) - \nabla \bar{\mathbf{y}}_k^*(\mathbf{x}^r)(\mathbf{x}^{r+1} - \mathbf{x}^r)\| \end{aligned} \quad (93)$$

$$\stackrel{(a)}{\leq} \frac{L_{xy}}{2} \mathbb{E} \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\| \mathbb{E} [\|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 | \mathcal{F}^r] \quad (94)$$

$$\leq \frac{\vartheta L_{xy}}{4} \mathbb{E} [\|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \mathbb{E} [\|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 | \mathcal{F}^r]] + \frac{L_{xy}}{4\vartheta} \mathbb{E} [\|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 | \mathcal{F}^r] \quad (95)$$

$$\stackrel{(b)}{\leq} \frac{\vartheta D_x L_{xy}}{4\alpha^2} \mathbb{E} [\|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2] + \frac{L_{xy}}{4\vartheta} \left(\mathbb{E}[\|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2] + \frac{\sigma_f^2}{n\alpha^2} \right) \quad (96)$$

where (a) follows the Lipschitz continuity of $\nabla \mathbf{y}_k^*(\mathbf{x})$ shown in (26), in (b) we apply Lemma 4, and

$$\begin{aligned} & \mathbb{E}_{\xi^r} [\|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 | \mathcal{F}^{r'}] \\ &= \|\mathbb{E}_{\xi^r} [\mathbf{x}^{r+1} - \mathbf{x}^r] | \mathcal{F}^{r'}\|^2 + \mathbb{E}_{\xi^r} [\|\mathbf{x}^{r+1} - \mathbf{x}^r - \mathbb{E}_{\xi^r} [\mathbf{x}^{r+1} - \mathbf{x}^r] | \mathcal{F}^{r'}\|^2] \end{aligned} \quad (97)$$

$$= \|\mathbb{E}_{\xi^r} [\mathbf{x}^{r+1} - \mathbf{x}^r] | \mathcal{F}^{r'}\|^2 + \mathbb{E}_{\xi^r} \left[\left\| \frac{1}{\alpha} (\mathbb{E} \mathbf{h}_f^r - \mathbf{h}_f^r) \right\|^2 | \mathcal{F}^{r'} \right] \quad (98)$$

$$\leq \|\mathbb{E}_{\xi^r} [\mathbf{x}^{r+1} - \mathbf{x}^r] | \mathcal{F}^{r'}\|^2 + \frac{\sigma_f^2}{n\alpha^2}. \quad (99)$$

Combining (92) and (96), we can have

$$\begin{aligned} \mathbb{E} \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1})\|^2 &\leq \left(1 + \theta' + \frac{\vartheta L_{xy} D_x}{4\alpha^2} \right) \mathbb{E} \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \\ &\quad + \left(L_y^2 + \frac{L_y^2}{4\theta'} + \frac{L_{xy}}{4\vartheta} \right) \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 + \left(\frac{L_{xy}}{4\vartheta} + L_y \right) \frac{\sigma_f^2}{n\alpha^2}. \end{aligned} \quad (100)$$

where we use the gradient Lipschitz continuity, i.e., $\mathbb{E} \|\bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1}) - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \leq L_y^2 \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2$. \square

B.3 Upper Bound of Successive Dual Variables

In what follows, we will show the ascent part after one round of the dual variable update.

Lemma 6.

Under A1-A5, define $\mathbf{D} \triangleq \alpha \mathbf{I} - \rho\gamma \mathbf{A}^T \mathbf{A}$. Suppose that the sequence $\{\mathbf{x}^r, \mathbf{y}_k^r, \boldsymbol{\lambda}^r, \boldsymbol{\omega}_k^r, \forall k\}$ is generated by SLAM. Then, we have

$$\begin{aligned} \frac{\gamma}{\rho} \|\mathbb{E} \boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r\|^2 &\leq \frac{4m^2 L_{f,x}^2 D_x}{n^2 \rho \gamma \alpha^2 \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} + \frac{4 \|\mathbb{E} \mathbf{w}^{r+1}\|_{\mathbf{D}^T \mathbf{D}}^2}{\rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \\ &\quad + \frac{4m L_{f,y}^2 \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2}{n^2 \rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} + \frac{4(b_r + b_{r-1})^2}{\rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}, \end{aligned} \quad (101a)$$

$$\frac{\gamma}{\rho} \|\mathbb{E} \boldsymbol{\omega}_k^{r+1} - \boldsymbol{\omega}_k^r\|^2 \leq \frac{3 \|\mathbb{E} \mathbf{v}_k^{r+1}\|_{\mathbf{D}^T \mathbf{D}}^2}{\rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} + \frac{3L_{g,1}^2 D_x}{n^2 \rho \gamma \alpha^2 \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} + \frac{3L_{g,1}^2 \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2}{n^2 \rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \quad (101b)$$

where $\mathbf{w}^{r+1} \triangleq (\mathbf{x}^{r+1} - \mathbf{x}^r) - (\mathbf{x}^r - \mathbf{x}^{r-1})$ and $\mathbf{v}_k^{r+1} \triangleq (\mathbf{y}_k^{r+1} - \mathbf{y}_k^r) - (\mathbf{y}_k^r - \mathbf{y}_k^{r-1}), \forall k$.

Proof. First Part. First, by utilizing the optimality condition of (5b), we have

$$\mathbf{h}_f^r + \gamma \mathbf{A}^T \boldsymbol{\lambda}^{r+1} + \rho\gamma \mathbf{A}^T \mathbf{A} (\mathbf{x}^r - \mathbf{x}^{r+1}) + \alpha (\mathbf{x}^{r+1} - \mathbf{x}^r) = 0. \quad (102)$$

Subtracting the above equality with the same one from the previous iteration, we obtain

$$\begin{aligned} \mathbf{h}_f^r - \mathbf{h}_f^{r-1} + \gamma \mathbf{A}^T (\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r) + \rho\gamma \mathbf{A}^T \mathbf{A} ((\mathbf{x}^r - \mathbf{x}^{r+1}) - (\mathbf{x}^{r-1} - \mathbf{x}^r)) \\ + \alpha ((\mathbf{x}^{r+1} - \mathbf{x}^r) - (\mathbf{x}^r - \mathbf{x}^{r-1})) = 0. \end{aligned} \quad (103)$$

Let $\mathbf{w}^{r+1} \triangleq (\mathbf{x}^{r+1} - \mathbf{x}^r) - (\mathbf{x}^r - \mathbf{x}^{r-1})$. We can easily write (103) concisely as

$$\mathbf{h}_f^r - \mathbf{h}_f^{r-1} + \gamma \mathbf{A}^T (\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r) + (\alpha - \rho\gamma \mathbf{A}^T \mathbf{A}) \mathbf{w}^{r+1} = 0. \quad (104)$$

According to A3, taking expectation on both sides of the above equation, we have

$$\begin{aligned} \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1}) - \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^r) + b'_r + b'_{r-1} + \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^r) - \bar{\nabla} f(\mathbf{x}^{r-1}, \mathbf{y}_k^r) \\ + \mathbb{E}[\gamma \mathbf{A}^T (\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r)] + \mathbb{E}[(\alpha - \rho\gamma \mathbf{A}^T \mathbf{A}) \mathbf{w}^{r+1}] = 0 \end{aligned} \quad (105)$$

where $0 < b'_r \leq b_r, 0 < b'_{r-1} \leq b_{r-1}$.

Utilizing the fact that $\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r$ lies in the column space of \mathbf{A} , we have $\|\mathbf{A}^T (\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r)\|^2 \geq \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A}) \|\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r\|^2$. After applying the triangle inequality, it is easy to see that the following in-

equality is true

$$\begin{aligned}
& \frac{\gamma}{\rho} \|\mathbb{E}\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r\|^2 \\
& \leq \frac{4}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \|\mathbb{E}[\bar{\nabla}f(\mathbf{x}^r, \mathbf{y}_k^r) - \bar{\nabla}f(\mathbf{x}^{r-1}, \mathbf{y}_k^r)]\|^2 + \frac{4}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \|\mathbb{E}\mathbf{w}^{r+1}\|_{\mathbf{D}^T\mathbf{D}}^2 \\
& \quad + \frac{4}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \|\mathbb{E}[\bar{\nabla}f(\mathbf{x}^r, \mathbf{y}_k^{r+1}) - \bar{\nabla}f(\mathbf{x}^r, \mathbf{y}_k^r)]\|^2 + \frac{4(b'_r + b'_{r-1})^2}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \\
& \stackrel{(a)}{\leq} \frac{4m^2L_{f,x}^2}{n^2\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 + \frac{4}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \|\mathbb{E}\mathbf{w}^{r+1}\|_{\mathbf{D}^T\mathbf{D}}^2 \\
& \quad + \frac{4L_{f,y}^2m}{n^2\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \sum_{k=1}^m \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{4(b'_r + b'_{r-1})^2}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \\
& \stackrel{(b)}{\leq} \frac{4m^2L_{f,x}^2D_x}{n^2\rho\gamma\alpha^2\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} + \frac{4\|\mathbb{E}\mathbf{w}^{r+1}\|_{\mathbf{D}^T\mathbf{D}}^2}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} + \frac{4mL_{f,y}^2\sum_{k=1}^m\mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2}{n^2\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} + \frac{4(b_r + b_{r-1})^2}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})}
\end{aligned} \tag{106}$$

where in (a) we use gradient Lipschitz continuity of the UL loss function w.r.t. \mathbf{x} and \mathbf{y}_k , i.e., (33) and (34), and (b) is true by applying Lemma 4.

Second Part. Utilizing the optimality condition of (5a), we have

$$\mathbf{h}_{g,k}^r + \gamma\mathbf{A}^T\boldsymbol{\omega}_k^{r+1} + \rho\gamma\mathbf{A}^T\mathbf{A}(\mathbf{y}_k^r - \mathbf{y}_k^{r+1}) + \alpha(\mathbf{y}_k^{r+1} - \mathbf{y}_k^r) = 0. \tag{107}$$

Similarly, we have

$$\begin{aligned}
& \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^{r+1}) - \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^r) + \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^r) - \nabla g_k(\mathbf{x}^{r-1}, \mathbf{y}_k^r) \\
& \quad + \gamma\mathbb{E}[\mathbf{A}^T(\boldsymbol{\omega}_k^{r+1} - \boldsymbol{\omega}_k^r)] + \mathbb{E}[(\beta - \rho\gamma\mathbf{A}^T\mathbf{A})\mathbf{v}_k^{r+1}] = 0.
\end{aligned} \tag{108}$$

where we have defined $\mathbf{v}_k^{r+1} \triangleq (\mathbf{y}_k^{r+1} - \mathbf{y}_k^r) - (\mathbf{y}_k^r - \mathbf{y}_k^{r-1})$.

Following (106), we have

$$\begin{aligned}
& \frac{\gamma}{\rho} \|\mathbb{E}[\boldsymbol{\omega}_k^{r+1} - \boldsymbol{\omega}_k^r]\|^2 \\
& \leq \frac{3}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \|\mathbb{E}[\nabla g_k(\mathbf{x}^r, \mathbf{y}_k^r) - \nabla g_k(\mathbf{x}^{r-1}, \mathbf{y}_k^r)]\|^2 + \frac{3}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \|\mathbb{E}\mathbf{v}_k^{r+1}\|_{\mathbf{D}^T\mathbf{D}}^2 \\
& \quad + \frac{3}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \|\mathbb{E}[\nabla g_k(\mathbf{x}^r, \mathbf{y}_k^{r+1}) - \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^r)]\|^2 \\
& \stackrel{(a)}{\leq} \frac{3L_{g,1}^2}{n^2\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 + \frac{3}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \|\mathbb{E}\mathbf{v}_k^{r+1}\|_{\mathbf{D}^T\mathbf{D}}^2 \\
& \quad + \frac{3L_{g,1}^2}{n^2\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2,
\end{aligned} \tag{109}$$

which gives (101b) directly by applying Lemma 4. \square

B.4 Recursive Functions

Now, the ascent part measured by the successive difference of the dual variables is upper bounded w.r.t. $\|\mathbb{E}[\mathbf{w}^{r+1}]\|^2$. Using (5b), we can construct the following recursion that establishes descent w.r.t. $\|\mathbb{E}[\mathbf{w}^{r+1}]\|^2$.

Lemma 7.

Under A1-A4, suppose that the sequence is generated by SLAM. Then, there exists a constant $\vartheta > 0$ such that

$$\begin{aligned}
\mathcal{Q}_w^{r+1} - \mathcal{Q}_w^r & \leq -\frac{1}{2} \|\mathbb{E}[\mathbf{w}^{r+1}]\|_{\mathbf{D}}^2 + \left(\frac{L_{f,x} + 1}{2n} + \frac{L_{f,x}m^2}{n} \right) \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2 \\
& \quad + \frac{mL_{f,y}^2}{L_{f,x}n} \sum_{k=1}^m \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{n(b_r + b_{r-1})^2}{2} + \frac{L_{f,x}m^2\sigma_f^2}{n\alpha^2}
\end{aligned} \tag{110}$$

where

$$\mathcal{Q}_w^r \triangleq \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A}\mathbb{E}[\mathbf{x}^r]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r-1}]\|_{\mathbf{D}}^2 + \frac{L_{f,x}m^2}{n} \mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 \geq 0. \tag{111}$$

Proof. We have the following optimality condition for the \mathbf{x} -update step:

$$\mathbb{E}[\mathbf{h}_f^r + \gamma \mathbf{A}^T \boldsymbol{\lambda}^r + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbf{x}^r + \alpha(\mathbf{x}^{r+1} - \mathbf{x}^r)] = 0, \quad (112a)$$

$$\mathbb{E}[\mathbf{h}_f^{r-1} + \gamma \mathbf{A}^T \boldsymbol{\lambda}^{r-1} + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbf{x}^{r-1} + \alpha(\mathbf{x}^r - \mathbf{x}^{r-1})] = 0. \quad (112b)$$

Hence, we have

$$\bar{\mathbf{h}}_f^r + \gamma \mathbf{A}^T \mathbb{E}[\boldsymbol{\lambda}^r] + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbb{E}[\mathbf{x}^r] + \alpha \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] = 0, \quad (113a)$$

$$\bar{\mathbf{h}}_f^{r-1} + \gamma \mathbf{A}^T \mathbb{E}[\boldsymbol{\lambda}^{r-1}] + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbb{E}[\mathbf{x}^{r-1}] + \alpha \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r-1}] = 0. \quad (113b)$$

Although $\bar{\mathbf{h}}_f^r$ is a biased gradient estimate, we can have

$$\bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1}) + b'_r + \gamma \mathbf{A}^T \mathbb{E}[\boldsymbol{\lambda}^r] + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbb{E}[\mathbf{x}^r] + \alpha \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] = 0, \quad (114a)$$

$$\bar{\nabla} f(\mathbf{x}^{r-1}, \mathbf{y}_k^r) + b'_{r-1} + \gamma \mathbf{A}^T \mathbb{E}[\boldsymbol{\lambda}^{r-1}] + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbb{E}[\mathbf{x}^{r-1}] + \alpha \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r-1}] = 0. \quad (114b)$$

Multiplying $\mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}]$ on both sides of (113a), we get

$$\begin{aligned} \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1}) \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}] + \gamma \mathbf{A}^T \mathbb{E}[\boldsymbol{\lambda}^r] \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}] + b'_r \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}] \\ + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbb{E}[\mathbf{x}^r] \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}] + \alpha \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}] = 0, \end{aligned} \quad (115)$$

and similarly multiplying $\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]$ on both sides of (113b), we can further have

$$\begin{aligned} \bar{\nabla} f(\mathbf{x}^{r-1}, \mathbf{y}_k^r) \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] + \gamma \mathbf{A}^T \mathbb{E}[\boldsymbol{\lambda}^{r-1}] \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] + b'_{r-1} \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] \\ + \rho \gamma \mathbf{A}^T \mathbf{A} \mathbb{E}[\mathbf{x}^{r-1}] \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] + \alpha \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r-1}] \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] = 0. \end{aligned} \quad (116)$$

Plugging (4d) into the above two inequalities, we obtain

$$\begin{aligned} & \frac{\rho}{\sqrt{\gamma}} \langle \mathbf{A}^T \mathbf{A} \mathbb{E}[\mathbf{x}^r], \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] \rangle \\ &= \gamma \langle \mathbf{A}^T \mathbb{E}[\boldsymbol{\lambda}^r - \boldsymbol{\lambda}^{r-1}], \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] \rangle \\ &\leq \langle \bar{\nabla} f(\mathbf{x}^{r-1}, \mathbf{y}_k^r) - \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1}) - \mathbf{D} \mathbb{E}[\mathbf{w}^{r+1}], \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] \rangle + (b'_r - b'_{r-1}) \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}], \end{aligned} \quad (117)$$

which gives

$$\begin{aligned} & \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A} \mathbb{E}[\mathbf{x}^{r+1}]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|_{\mathbf{D}}^2 \\ &\leq \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A} \mathbb{E}[\mathbf{x}^r]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r-1}]\|_{\mathbf{D}}^2 - \frac{1}{2} \|\mathbb{E}[\mathbf{w}^{r+1}]\|_{\mathbf{D}}^2 \\ &\quad + \langle \bar{\nabla} f(\mathbf{x}^{r-1}, \mathbf{y}_k^r) - \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1}), \mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r] \rangle + (b'_r - b'_{r-1}) \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}] \\ &\leq \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A} \mathbb{E}[\mathbf{x}^r]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r-1}]\|_{\mathbf{D}}^2 - \frac{1}{2} \|\mathbb{E}[\mathbf{w}^{r+1}]\|_{\mathbf{D}}^2 \\ &\quad + \frac{n}{2L_{f,x}} \mathbb{E} \|\bar{\nabla} f(\mathbf{x}^{r-1}, \mathbf{y}_k^r) - \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1})\|^2 + \frac{L_{f,x}}{2n} \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2 + (b'_r - b'_{r-1}) \mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}] \\ &\stackrel{(a)}{\leq} \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A} \mathbb{E}[\mathbf{x}^r]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r-1}]\|_{\mathbf{D}}^2 - \frac{1}{2} \|\mathbb{E}[\mathbf{w}^{r+1}]\|_{\mathbf{D}}^2 \\ &\quad + \frac{L_{f,x} m^2}{n} \mathbb{E} \|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 + \frac{m L_{f,y}^2}{n L_{f,x}} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{L_{f,x}}{2n} \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2 \\ &\quad + \frac{n(b'_r - b'_{r-1})^2}{2} + \frac{\|\mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r+1}]\|^2}{2n} \end{aligned} \quad (119)$$

where (a) follows gradient Lipschitz continuity.

Therefore, we have

$$\begin{aligned} & \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A} \mathbb{E}[\mathbf{x}^{r+1}]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|_{\mathbf{D}}^2 + \frac{L_{f,x} m^2}{n} \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 \\ &\leq \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A} \mathbb{E}[\mathbf{x}^r]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{x}^r - \mathbf{x}^{r-1}]\|_{\mathbf{D}}^2 + \frac{L_{f,x} m^2}{n} \mathbb{E} \|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 \\ &\quad - \frac{1}{2} \|\mathbb{E}[\mathbf{w}^{r+1}]\|_{\mathbf{D}}^2 + \left(\frac{L_{f,x} + 1}{2n} + \frac{L_{f,x} m^2}{n} \right) \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2 + \frac{m L_{f,y}^2}{n L_{f,x}} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \\ &\quad + \frac{n(b_r + b_{r-1})^2}{2} + \frac{L_{f,x} m^2 \sigma_f^2}{n \alpha^2} \end{aligned} \quad (120)$$

where we use (99) and (17). \square

Lemma 8.

Under A1-A4, suppose that the sequence is generated by SLAM. Let

$$\mathcal{Q}_{v,k}^r \triangleq \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A}\mathbb{E}[\mathbf{y}_k^r]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|_{\mathbf{D}}^2 + \frac{L_{g,1}}{2n} \mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 + \frac{L_{g,1}}{2n} \mathbb{E}\|\mathbf{y}_k^r - \mathbf{y}_k^{r-1}\|^2 \geq 0, \quad (121)$$

then, the following is true,

$$\mathcal{Q}_{v,k}^{r+1} - \mathcal{Q}_{v,k}^r \leq \frac{L_{g,1}}{2n} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 + \frac{3L_{g,1}}{2n} \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 - \frac{1}{2} \|\mathbb{E}[\mathbf{v}_k^{r+1}]\|_{\mathbf{D}}^2, \forall k. \quad (122)$$

Proof. Following steps from (112) to (117), we can similarly obtain the following series of inequalities for sequence $\{\mathbf{y}_k^r, \forall k\}$.

$$\begin{aligned} & \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A}\mathbb{E}[\mathbf{y}_k^{r+1}]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r]\|_{\mathbf{D}}^2 \\ & \leq \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A}\mathbb{E}[\mathbf{y}_k^r]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|_{\mathbf{D}}^2 - \frac{1}{2} \|\mathbb{E}[\mathbf{v}_k^{r+1}]\|_{\mathbf{D}}^2 \\ & \quad + \langle \nabla g_k(\mathbf{x}^{r-1}, \mathbf{y}_k^{r-1}) - \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^r), \mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] \rangle \\ & \leq \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A}\mathbb{E}[\mathbf{y}_k^r]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|_{\mathbf{D}}^2 - \frac{1}{2} \|\mathbb{E}[\mathbf{v}_k^{r+1}]\|_{\mathbf{D}}^2 \\ & \quad + \langle \nabla g_k(\mathbf{x}^{r-1}, \mathbf{y}_k^{r-1}) - \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^{r-1}) + \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^{r-1}) - \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^r), \mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] \rangle \end{aligned} \quad (123)$$

$$\begin{aligned} & \stackrel{(a)}{\leq} \frac{\rho}{2\sqrt{\gamma}} \|\mathbf{A}\mathbb{E}[\mathbf{y}_k^r]\|^2 + \frac{1}{2} \|\mathbb{E}[\mathbf{y}_k^r - \mathbf{y}_k^{r-1}]\|_{\mathbf{D}}^2 - \frac{1}{2} \|\mathbb{E}[\mathbf{v}_k^{r+1}]\|_{\mathbf{D}}^2 \\ & \quad + \frac{L_{g,1}}{2n} \mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 + \frac{L_{g,1}}{n} \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{L_{g,1}}{2n} \mathbb{E}\|\mathbf{y}_k^r - \mathbf{y}_k^{r-1}\|^2 \end{aligned} \quad (124)$$

where in (a) we use Young's inequality

$$\langle \nabla g_k(\mathbf{x}^{r-1}, \mathbf{y}_k^{r-1}) - \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^{r-1}), \mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] \rangle \leq \frac{L_{g,1}}{2n} \mathbb{E}\|\mathbf{x}^r - \mathbf{x}^{r-1}\|^2 + \frac{L_{g,1}}{2n} \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2, \quad (125)$$

and gradient Lipschitz $g(\cdot, \mathbf{y})$

$$\langle \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^{r-1}) - \nabla g_k(\mathbf{x}^r, \mathbf{y}_k^r), \mathbb{E}[\mathbf{y}_k^{r+1} - \mathbf{y}_k^r] \rangle \leq \frac{L_{g,1}}{2n} \mathbb{E}\|\mathbf{y}_k^r - \mathbf{y}_k^{r-1}\|^2 + \frac{L_{g,1}}{2n} \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2. \quad \square$$

B.5 Proof of Theorem 1

Proof. In this subsection, we will show the convergence rate regarding the stationarity, optimality, and consensus violation w.r.t. both UL and LL optimization variables generated by SLAM as follows.

B.5.1 Consensus Violation of the UL Variables

Using the fact $\mathbb{E}[X^2] = \mathbb{E}[X]^2 + \text{Var}[X]$ (17) gives

$$\begin{aligned} & \mathbb{E} [\|\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r\|^2 | \mathcal{F}^{r'}] - \|\mathbb{E}[\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r] | \mathcal{F}^{r'}\|^2 \\ & = \mathbb{E} [\|\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r - \mathbb{E}[\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r] | \mathcal{F}^{r'}\|^2 | \mathcal{F}^{r'}] = \frac{\rho^2}{\gamma^2} \mathbb{E} [\|\mathbf{A}\mathbf{x}^{r+1} - \mathbb{E}\mathbf{A}\mathbf{x}^{r+1}\|^2 | \mathcal{F}^{r'}] \end{aligned} \quad (126)$$

$$\leq \frac{\rho^2}{\gamma^2 \alpha^2} \sigma_{\max}(\mathbf{A}^T \mathbf{A}) \mathbb{E} \|\mathbf{h}_f^r - \mathbb{E}\mathbf{h}_f^r\|^2 \leq \frac{\rho^2 \sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sigma_f^2}{n \gamma^2 \alpha^2}. \quad (127)$$

From (101a), we have

$$\begin{aligned} & \mathbb{E} \|\mathbf{A}^T \mathbf{x}^{r+1}\|^2 \\ & = \frac{\gamma^2}{\rho^2} \|\mathbb{E}\boldsymbol{\lambda}^{r+1} - \boldsymbol{\lambda}^r\|^2 + \frac{\sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sigma_f^2}{n \alpha^2} \\ & \stackrel{(101a)}{\leq} \frac{\gamma}{\rho} \frac{4 \|\mathbb{E}\mathbf{w}^{r+1}\|_{\mathbf{D}^T \mathbf{D}}^2}{\rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} + \frac{4(m+1)}{n^2 \rho^2 \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \left(\frac{(m+1) D_x L_{f,x}^2}{\alpha^2} + L_{f,y}^2 \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \right) \\ & \quad + \frac{4(b_r + b_{r-1})^2}{\rho^2 \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} + \frac{\sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sigma_f^2}{n \alpha^2}. \end{aligned} \quad (128)$$

Considering the recursion derived in (110), we can obtain that there is a constant C such that

$$\begin{aligned} \mathbb{E}\|\mathbf{A}^T \mathbf{x}^{r+1}\|^2 &\leq \frac{\gamma}{\rho} \left(C\mathcal{Q}_w^r - C\mathcal{Q}_w^{r+1} - \mathbb{E}[\mathbf{w}^{r+1}]^T \left(\frac{C\mathbf{D}}{2} - \frac{4\mathbf{D}^T\mathbf{D}}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \right) \mathbb{E}[\mathbf{w}^{r+1}] \right) \\ &\quad + \left(\frac{\gamma}{\rho} \left(\frac{L_{f,1}+1}{2n} + \frac{L_{f,x}m^2}{n} \right) + \frac{4L_{f,x}^2m^2}{n^2\rho^2\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \right) \frac{D_x}{\alpha^2} \\ &\quad + \frac{mL_{f,y}^2}{n} \left(\frac{\gamma}{\rho} \frac{1}{L_{f,x}} + \frac{4}{n\rho^2\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \right) \sum_{k=1}^m \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \\ &\quad + \frac{\gamma}{\rho} \left(\frac{n(b_r + b_{r-1})^2}{2} + \frac{L_{f,x}m^2\sigma_f^2}{n^2\alpha^2} \right) + \frac{4(b_r + b_{r-1})^2}{\rho^2\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} + \frac{\sigma_{\max}(\mathbf{A}^T\mathbf{A})\sigma_f^2}{n\alpha^2} \end{aligned} \quad (130)$$

To show the descent of the right-hand side of the above inequality, we need

$$\frac{C\mathbf{D}}{2} - \frac{4\mathbf{D}^T\mathbf{D}}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \succ 0, \quad (131)$$

so, it is sufficient to show

$$\frac{C}{2} (\alpha\mathbf{I} - \rho\gamma\mathbf{A}^T\mathbf{A}) - \frac{4}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} (\alpha\mathbf{I} - \rho\gamma\mathbf{A}^T\mathbf{A})^T (\alpha\mathbf{I} - \rho\gamma\mathbf{A}^T\mathbf{A}) \succ 0 \quad (132)$$

$$\Rightarrow \frac{C}{2} (\alpha\mathbf{I} - \rho\gamma\mathbf{A}^T\mathbf{A}) - \frac{4\alpha^2\mathbf{I}}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} + \frac{8\alpha\mathbf{A}^T\mathbf{A}}{\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} - \frac{4\rho\gamma(\mathbf{A}^T\mathbf{A})^2}{\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \succ 0 \quad (133)$$

$$\Rightarrow \frac{C}{2}\alpha\mathbf{I} - \frac{4\alpha^2\mathbf{I}}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} - \frac{4\rho\gamma(\mathbf{A}^T\mathbf{A})^2}{\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \succ 0 \quad \text{and} \quad \frac{8\alpha\mathbf{A}^T\mathbf{A}}{\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} - \frac{C}{2}\rho\mathbf{A}^T\mathbf{A} \succ 0 \quad (134)$$

$$\Rightarrow \frac{C}{2}\alpha\mathbf{I} - \frac{4\alpha^2\mathbf{I}}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} - \frac{4\rho\gamma(\mathbf{A}^T\mathbf{A})^2}{\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \succ 0 \quad \text{and} \quad \frac{8\alpha\mathbf{A}^T\mathbf{A}}{\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} - \frac{C}{2}\rho\mathbf{A}^T\mathbf{A} \succ 0 \quad (135)$$

$$\Rightarrow \frac{C}{2}\alpha - \frac{4\alpha}{\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \left(\frac{\alpha}{\rho\gamma} + \frac{\rho\gamma\sigma_{\max}^2(\mathbf{A}^T\mathbf{A})}{\alpha} \right) \geq 0 \quad \text{and} \quad \alpha \geq \frac{C\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})}{16}, \quad (136)$$

which means

$$\frac{\tau_f}{\rho} + \frac{\rho\sigma_{\max}^2(\mathbf{A}^T\mathbf{A})}{\tau_f} \leq \frac{\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})C}{8} \quad \text{and} \quad \alpha \geq \frac{C\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})}{16}. \quad (137)$$

Applying (131) and the telescoping sum, we have

$$\begin{aligned} \frac{1}{T} \sum_{r=1}^T \mathbb{E}\|\mathbf{A}^T \mathbf{x}^{r+1}\|^2 &\leq \frac{C\gamma\mathcal{Q}_w^1}{T\rho} + \left(\frac{\gamma}{\rho} \left(\frac{L_{f,x}+1}{2n} + \frac{L_{f,x}m^2}{n} \right) + \frac{4L_{f,x}^2m^2}{n^2\rho^2\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \right) \frac{D_x}{\alpha^2} \\ &\quad + \left(\frac{\gamma}{\rho} \frac{mL_{f,y}^2}{nL_{f,x}} + \frac{4L_{f,y}^2m}{n^2\rho^2\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \right) \frac{D_y}{\beta^2} \\ &\quad + \frac{\gamma}{\rho} \left(\frac{n(b_r + b_{r-1})^2}{2} + \frac{L_{f,x}m^2\sigma_f^2}{n\alpha^2} \right) + \frac{4(b_r + b_{r-1})^2}{\rho^2\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} + \frac{\sigma_{\max}(\mathbf{A}^T\mathbf{A})\sigma_f^2}{n\alpha^2} \\ &\sim \mathcal{O}\left(\frac{1}{\sqrt{nT}}\right) \end{aligned} \quad (138)$$

where we choose $\alpha = \alpha_0\sqrt{T/n}$, $\beta = \beta_0\sqrt{T/n}$.

B.5.2 Consensus Violation of the LL Optimization Variables

From Lemma 8, we know that $\mathcal{Q}_{v,k}^r$ is lower bounded by 0. According to (122), we have

$$\mathcal{Q}_{v,k}^{r+1} - \mathcal{Q}_{v,k}^r \leq \frac{L_{g,1}}{n} \left(\frac{D_x}{2\alpha^2} + \frac{3}{2}\mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \right) - \frac{1}{2}\|\mathbb{E}[\mathbf{v}_k^{r+1}]\|_{\mathbf{D}}^2. \quad (139)$$

Multiplying (139) by constant C and adding (101b) together, we have

$$\begin{aligned} \frac{\gamma}{\rho} \|\mathbb{E}\mathbf{w}_k^{r+1} - \mathbf{w}_k^r\|^2 &\leq C\mathcal{Q}_{v,k}^r - C\mathcal{Q}_{v,k}^{r+1} - \mathbb{E}[\mathbf{v}_k^{r+1}]^T \left(\frac{C\mathbf{D}}{2} - \frac{3\mathbf{D}^T\mathbf{D}}{\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \right) \mathbb{E}[\mathbf{v}_k^{r+1}] \\ &\quad + \frac{CL_{g,1}}{n} \left(\frac{D_x}{2\alpha^2} + \frac{3}{2}\mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \right) + \frac{3L_{g,1}^2}{n^2\rho\gamma\tilde{\sigma}_{\min}(\mathbf{A}^T\mathbf{A})} \left(\frac{D_x}{\alpha^2} + \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \right). \end{aligned} \quad (140)$$

Note that

$$\begin{aligned} & \mathbb{E} [\|\boldsymbol{\omega}_k^{r+1} - \boldsymbol{\omega}_k^r\|^2 | \mathcal{F}^r] - \|\mathbb{E}[\boldsymbol{\omega}_k^{r+1} - \boldsymbol{\omega}_k^r] | \mathcal{F}^r\|^2 \\ &= \mathbb{E} [\|\boldsymbol{\omega}_k^{r+1} - \boldsymbol{\omega}_k^r - \mathbb{E}[\boldsymbol{\omega}_k^{r+1} - \boldsymbol{\omega}_k^r] | \mathcal{F}^r\|^2 | \mathcal{F}^r] = \frac{\rho^2}{\gamma^2 \beta^2} \mathbb{E} [\|\mathbf{A}\mathbf{y}_k^{r+1} - \mathbb{E}\mathbf{A}\mathbf{y}_k^{r+1}\|^2 | \mathcal{F}^r] \end{aligned} \quad (141)$$

$$\leq \frac{\rho^2}{\gamma^2 \beta^2} \sigma_{\max}(\mathbf{A}^T \mathbf{A}) \|\mathbf{h}_{g,k}^r - \mathbb{E}\mathbf{h}_{g,k}^r\|^2 \leq \frac{\rho^2 \sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sigma_g^2}{n \gamma^2 \beta^2}. \quad (142)$$

From (4b), we have

$$\begin{aligned} & \mathbb{E} \|\mathbf{A}\mathbf{y}_k^{r+1}\|^2 \\ & \stackrel{(4b)}{\leq} \frac{\gamma^2}{\rho^2} \mathbb{E} \|\boldsymbol{\omega}_k^{r+1} - \boldsymbol{\omega}_k^r\|^2 \end{aligned} \quad (143)$$

$$\stackrel{(142)}{\leq} \frac{\gamma^2}{\rho^2} \mathbb{E} \|\boldsymbol{\omega}_k^{r+1} - \boldsymbol{\omega}_k^r\|^2 + \frac{\sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sigma_g^2}{n \beta^2} \quad (144)$$

$$\begin{aligned} & \stackrel{(140)}{\leq} \frac{\gamma}{\rho} \left(C\mathcal{Q}_{v,k}^r - C\mathcal{Q}_{v,k}^{r+1} - \mathbb{E}[\mathbf{v}_k^{r+1}]^T \left(\frac{C\mathbf{D}}{2} - \frac{3\mathbf{D}^T \mathbf{D}}{\rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \right) \mathbb{E}[\mathbf{v}_k^{r+1}] \right) + \frac{\sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sigma_g^2}{n \beta^2} \\ & \quad + \frac{\gamma}{\rho n} \left(\frac{CL_{g,1}}{2} + \frac{3L_{g,1}^2}{n \rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \right) \frac{D_x}{\alpha^2} + \frac{\gamma}{\rho} \left(\frac{3CL_{g,1}}{2n} + \frac{3L_{g,1}^2}{n^2 \rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \right) \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \\ & \stackrel{(131), (137)}{\leq} \frac{\gamma}{\rho} (C\mathcal{Q}_{v,k}^r - C\mathcal{Q}_{v,k}^{r+1}) + \frac{\gamma}{\rho} \left(\frac{3CL_{g,1}}{2n} + \frac{3L_{g,1}^2}{n^2 \rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \right) \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + N_1 \end{aligned} \quad (145)$$

where

$$N_1 \triangleq \frac{C\gamma L_{g,1} D_x}{2n\rho\alpha^2} + \frac{\sigma_{\max}(\mathbf{A}^T \mathbf{A}) \sigma_g^2}{n\beta^2} + \frac{3L_{g,1}^2 D_x}{n^2 \rho^2 \alpha^2 \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})}. \quad (146)$$

Applying telescoping sum, we can get

$$\begin{aligned} & \frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\mathbf{A}\mathbf{y}_k^{r+1}\|^2 \\ & \leq \frac{C\gamma \mathcal{Q}_{v,k}^1}{T\rho} + N_1 + \frac{\gamma}{\rho} \left(\frac{3CL_{g,1}}{2n} + \frac{3L_{g,1}^2}{n^2 \rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \right) \frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \end{aligned} \quad (147)$$

$$\leq \frac{C\gamma \mathcal{Q}_{v,k}^1}{T\rho} + N_1 + \frac{\gamma}{\rho} \left(\frac{3CL_{g,1}}{2n} + \frac{3L_{g,1}^2}{n^2 \rho \gamma \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \right) \frac{D_y}{\beta^2} \quad (148)$$

$$\sim \mathcal{O} \left(\frac{1}{\sqrt{nT}} \right), \quad (149)$$

when $\gamma \sim \beta \sim \alpha \sim \mathcal{O}(\sqrt{T/n})$.

B.5.3 Stationarity of the UL Optimization Variable

From (5a) and (5b), we have

$$\bar{\mathbf{x}}^{r+1} = \bar{\mathbf{x}}^r - \frac{1}{\alpha} \bar{\mathbf{h}}_f^r, \quad \bar{\mathbf{y}}_k^{r+1} = \bar{\mathbf{y}}_k^r - \frac{1}{\beta} \bar{\mathbf{h}}_{g,k}^r. \quad (150)$$

Note that $\mathbf{y}_k^*(\mathbb{1}\bar{\mathbf{x}}^r) = \mathbb{1}\bar{\mathbf{y}}_k^*(\mathbb{1}\bar{\mathbf{x}}^r)$ due to (2c). From the notations shown in (2), we know that

$$f(\mathbb{1}\bar{\mathbf{x}}, \mathbb{1}\bar{\mathbf{y}}_k^*(\mathbb{1}\bar{\mathbf{x}})) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(\bar{\mathbf{x}}, \mathbf{y}_{i,k}^*(\bar{\mathbf{x}})), \quad (151a)$$

$$g_k(\mathbb{1}\bar{\mathbf{x}}, \mathbb{1}\bar{\mathbf{y}}_k) \triangleq \frac{1}{n} \sum_{i=1}^n g_{i,k}(\bar{\mathbf{x}}, \bar{\mathbf{y}}_k), \quad \forall k \in [m]. \quad (151b)$$

To simply the notations, we ignore $\mathbb{1}$. For example, we just use $f(\bar{\mathbf{x}}, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}))$ and $g_k(\bar{\mathbf{x}}, \bar{\mathbf{y}}_k)$ as the abbreviations of $f(\mathbb{1}\bar{\mathbf{x}}, \mathbb{1}\bar{\mathbf{y}}_k^*(\mathbb{1}\bar{\mathbf{x}}))$ and $g_k(\mathbb{1}\bar{\mathbf{x}}, \mathbb{1}\bar{\mathbf{y}}_k)$ in the following derivation in this proof. Similarly, we have

$$\bar{\nabla} f(\mathbf{x}, \mathbf{y}_k) \triangleq \frac{1}{n} \sum_{i=1}^n \bar{\nabla} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}). \quad (152)$$

Then, we can have the descent at the consensus space as follows

$$\begin{aligned} & \mathbb{E}_{\xi^r} [f(\bar{\mathbf{x}}^{r+1}, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^{r+1})) | \mathcal{F}^{r'}] \\ & \stackrel{(a)}{\leq} f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) + \mathbb{E} [\langle \nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)), \bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r \rangle | \mathcal{F}^{r'}] + \mathbb{E} \left[\frac{L_f}{2} \|\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r\|^2 | \mathcal{F}^{r'} \right] \end{aligned} \quad (153)$$

$$\begin{aligned} & = f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) + \langle \nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)), \mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r] | \mathcal{F}^{r'} \rangle \\ & \quad + \frac{L_f}{2} \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 + \frac{L_f}{2} \mathbb{E} [\|\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r - \mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 | \mathcal{F}^{r'}] \end{aligned} \quad (154)$$

$$\begin{aligned} & \stackrel{(b)}{\leq} f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) - \frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r))\|^2 - \left(\frac{\alpha}{2} - \frac{L_f}{2} \right) \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 \\ & \quad + \frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) - \bar{\mathbf{h}}_f^r\|^2 + \frac{L_f}{2\alpha^2} \mathbb{E} [\|\bar{\mathbf{h}}_f^r - \mathbb{E}\bar{\mathbf{h}}_f^r\|^2 | \mathcal{F}^{r'}] \end{aligned} \quad (155)$$

$$\begin{aligned} & \stackrel{(c)}{\leq} f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) - \frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r))\|^2 - \left(\frac{\alpha}{2} - \frac{L_f}{2} \right) \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 \\ & \quad + \left(\frac{2(m+1)L_f^2}{\alpha n^2} + \frac{10(m+1)}{\alpha n^2} \left(\frac{C_{xy} L_{f,0} L_{g,2}}{\mu_g^2} \right)^2 + \frac{8(m+1)L_{y,c}^2 L_y^2}{\alpha n} \right) \|\mathbf{x}^r - \mathbb{1}\bar{\mathbf{x}}^r\|^2 \\ & \quad + \sum_{k=1}^m \frac{8(m+1)L_{y,c}^2}{n\alpha} \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)\|^2 + \left(\frac{4(mL_{f,y}^2 + (m+1)L_{y,c}^2)}{\alpha n^2} \right) \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 \\ & \quad + \frac{2mL_{f,y}^2}{n^2\alpha} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{2b_r^2}{n\alpha} + \frac{L_f \sigma_f^2}{2n\alpha^2} \end{aligned} \quad (156)$$

where (a) follows the gradient Lipschitz continuity of the UL loss function with constant L_f at the consensus space (which is the same as the centralized case, e.g., [31, Lemma 2.2.]), (b) is true because

$$\begin{aligned} & \langle \nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)), \mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r] | \mathcal{F}^{r'} \rangle \\ & = -\frac{1}{\alpha} \langle \nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)), \mathbb{E}\bar{\mathbf{h}}_f^r \rangle \\ & = -\frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r))\|^2 - \frac{1}{2\alpha} \|\mathbb{E}\bar{\mathbf{h}}_f^r\|^2 + \frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) - \mathbb{E}\bar{\mathbf{h}}_f^r\|^2, \end{aligned}$$

and in (c) we use the following steps: 1) the difference between the UL gradient and its distributed stochastic estimate can be quantified by

$$\begin{aligned} & \|\nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) - \mathbb{E}\bar{\mathbf{h}}_f^r\|^2 \\ & \stackrel{(a.1)}{\leq} 4\|\nabla f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) - \nabla f(\mathbf{x}^r, \mathbf{y}_k^*(\mathbf{x}^r))\|^2 + 4\|\nabla f(\mathbf{x}^r, \mathbf{y}_k^*(\mathbf{x}^r)) - \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^r)\|^2 \\ & \quad + 4\|\bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^r) - \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1})\|^2 + 4\|\bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1}) - \mathbb{E}\bar{\mathbf{h}}_f^r\|^2 \\ & \leq \left(\frac{4(m+1)L_f^2}{n^2} + \frac{10(m+1)}{n^2} \left(\frac{C_{xy} L_{f,0} L_{g,2}}{\mu_g^2} \right)^2 \right) \|\mathbf{x}^r - \mathbb{1}\bar{\mathbf{x}}^r\|^2 \\ & \quad + \sum_{k=1}^m \frac{4(m+1)L_{y,c}^2}{n^2} \|\mathbf{y}_k^*(\mathbf{x}^r) - \mathbf{y}_k^r\|^2 + \frac{4mL_{f,y}^2}{n^2} \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 \\ & \quad + \frac{4mL_{f,y}^2}{n^2} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{4b_r^2}{n} \\ & \stackrel{(a.2)}{\leq} \left(\frac{4(m+1)L_f^2}{n^2} + \frac{10(m+1)}{n^2} \left(\frac{C_{xy} L_{f,0} L_{g,2}}{\mu_g^2} \right)^2 + \frac{16(m+1)L_{y,c}^2 L_y^2}{n} \right) \|\mathbf{x}^r - \mathbb{1}\bar{\mathbf{x}}^r\|^2 \\ & \quad + \sum_{k=1}^m \frac{16(m+1)L_{y,c}^2}{n} \|\bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r) - \bar{\mathbf{y}}_k^r\|^2 + \left(\frac{8(mL_{f,y}^2 + (m+1)L_{y,c}^2)}{n^2} \right) \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 \\ & \quad + \frac{4mL_{f,y}^2}{n^2} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{4b_r^2}{n}, \end{aligned} \quad (157)$$

and in (a.1) the Lipschitz continuity of $\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r))$ w.r.t. \mathbf{x}_i with a constant (defined as L_f here) follows the centralized case directly (which has been shown in [31, Lemma 2.2.]); 2) constant $L_{y,c}$ is computed as

follows

$$\begin{aligned}
& \|\nabla f(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x}^r)) - \bar{\nabla} f(\mathbf{x}, \mathbf{y}_k^r)\|^2 \\
& \leq \frac{m+1}{n^2} \sum_{k=1}^m \left\| \sum_{i=1}^n \nabla_{\mathbf{x}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}) - \nabla_{\mathbf{x}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) \right\|^2 + \frac{m+1}{n^2} \left\| \sum_{i=1}^n \Delta_{i,k}^{(1)} + \Delta_{i,k}^{(2)} + \Delta_{i,k}^{(3)} + \Delta_{i,k}^{(4)} + \Delta_{i,k}^{(5)} \right\|^2 \\
& \leq \frac{m+1}{n^2} \sum_{k=1}^m \underbrace{L_{f,1}^2 + 5 \left(\frac{L_{f,0} L_{g,2}}{\mu_g} \right)^2 + 5 \left(\frac{C_{xy} L_{f,0} L_g}{\mu_g^2} \right)^2 + 5 \left(\frac{C_{xy} L_{f,1}}{\mu_g} \right)^2}_{\triangleq L_{y,c}^2} \|\mathbf{y}_k - \mathbf{y}_k^*(\mathbf{x})\|^2 \\
& \quad + \frac{10(m+1)}{n^2} \left(\frac{C_{xy} L_{f,0} L_{g,2}}{\mu_g^2} \right)^2 \|\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r\|^2, \tag{158}
\end{aligned}$$

and terms $\Delta_{i,k}^{(1)}, \Delta_{i,k}^{(2)}, \Delta_{i,k}^{(3)}, \Delta_{i,k}^{(4)}, \Delta_{i,k}^{(5)}$ are defined as

$$\begin{aligned}
\Delta_{i,k}^{(1)} & \triangleq \left[\nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) - \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) \right] \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}, \mathbf{y}_{i,k}) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}), \\
\Delta_{i,k}^{(2)} & \triangleq \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}) \right]^{-1} - \left[\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_{i,k}(\bar{\mathbf{x}}, \mathbf{y}_{i,k}) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}), \\
\Delta_{i,k}^{(3)} & \triangleq \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\bar{\mathbf{x}}, \mathbf{y}_{i,k}) \right]^{-1} - \left[\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\bar{\mathbf{x}}, \mathbf{y}_k^*(\mathbf{x})) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}), \\
\Delta_{i,k}^{(4)} & \triangleq \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) \left[\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\bar{\mathbf{x}}, \mathbf{y}_k^*(\mathbf{x})) \right]^{-1} - \left[\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}), \\
\Delta_{i,k}^{(5)} & \triangleq \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) \left[\nabla_{\mathbf{y}_k \mathbf{y}_k}^2 g_k(\mathbf{x}, \mathbf{y}_k^*(\mathbf{x})) \right]^{-1} \left[\nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}) - \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \mathbf{y}_{i,k}^*(\mathbf{x})) \right];
\end{aligned}$$

3) and in (a.2) we apply the triangle inequality and the fact that

$$\begin{aligned}
\|\bar{\mathbf{y}}_k^*(\mathbf{x}^r) - \bar{\mathbf{y}}_k^r\|^2 & \leq 2\|\bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r) - \bar{\mathbf{y}}_k^r\|^2 + 2\|\bar{\mathbf{y}}_k^*(\mathbf{x}_k) - \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}_k)\|^2 \tag{159} \\
& \leq 2\|\bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r) - \bar{\mathbf{y}}_k^r\|^2 + 2L_y^2 \|\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r\|^2 \tag{160}
\end{aligned}$$

based on (25).

Further, from Lemma 5, we know that

$$\begin{aligned}
\|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^{r+1})\|^2 & \leq \left(1 - \frac{\rho_g}{\beta}\right) \left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)\|^2 \\
& \quad + \left(L_y^2 + \frac{\alpha}{8(m+1)n} + \frac{L_{xy}}{4\vartheta}\right) \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 \\
& \quad + \left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \frac{\sigma_g^2}{n\beta^2} + \left(\frac{L_{xy}}{4\vartheta} + L_y^2\right) \frac{\sigma_f^2}{n\alpha^2} \tag{161}
\end{aligned}$$

where we choose $\theta' \triangleq 2(m+1)nL_y^2/\alpha$.

Next, let us define potential function at the consensus space as

$$\mathcal{P}_c^r = \mathbb{E}[f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) | \mathcal{F}^r] + \frac{m+1}{n} \sum_{k=1}^m \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)\|^2. \tag{162}$$

Subsequently, we can have

$$\begin{aligned}
& \mathcal{P}_c^{r+1} - \mathcal{P}_c^r \\
& \leq -\frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r))\|^2 - \underbrace{\left(\frac{\alpha}{2} - \left(\frac{L_f}{2} + \frac{m+1}{n} L_y^2 + \frac{\alpha}{8n^2} + \frac{(m+1)L_{xy}}{4n\vartheta} \right) \right)}_{\triangleq C_1} \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 \\
& \quad + \frac{m+1}{n} \left(\underbrace{\frac{8L_{y,c}^2}{\alpha} + \left(1 - \frac{\rho_g}{\beta}\right) \left(1 + \frac{2n(m+1)L_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right)}_{\triangleq C_2} - 1 \right) \sum_{k=1}^m \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)\|^2 \\
& \quad + \left(\frac{2(m+1)L_f^2}{\alpha n^2} + \frac{10(m+1)}{\alpha n^2} \left(\frac{C_{xy} L_{f,0} L_{g,2}}{\mu_g^2} \right)^2 + \frac{8(m+1)L_{y,c}^2 L_y^2}{\alpha n} \right) \|\mathbf{x}^r - \mathbb{1}\bar{\mathbf{x}}^r\|^2 \\
& \quad + \left(\frac{4(mL_{f,y}^2 + (m+1)L_{y,c}^2)}{\alpha n^2} \right) \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 + \frac{2mL_{f,y}^2}{n^2\alpha} \sum_{k=1}^m \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \\
& \quad + \frac{m+1}{n} \underbrace{\left(\left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \frac{\sigma_g^2}{n\beta^2} + \left(\frac{L_{xy}}{4\vartheta} + L_y^2\right) \frac{\sigma_f^2}{n\alpha^2} \right) + \frac{4b_r^2}{n\alpha} + \frac{L_f \sigma_f^2}{2n\alpha^2}}_{\triangleq N_2 \sim \mathcal{O}\left(\frac{1}{n\alpha^2}\right)}.
\end{aligned}$$

1) to show $C_1 > 0$: we need

$$\frac{\alpha}{2} - \left(\frac{L_f}{2} + \frac{m+1}{n} L_y^2 + \frac{\alpha}{8n^2} + \frac{(m+1)L_{xy}}{4n\vartheta} \right) > 0, \quad (163)$$

which requires

$$\alpha > 4 \left(\frac{L_f}{2} + \frac{m+1}{n} L_y^2 + \frac{(m+1)L_{xy}}{4n\vartheta} \right). \quad (164)$$

2) to show $C_2 < 0$: we need

$$\frac{2L_{y,c}^2}{\alpha} + \left(1 - \frac{\rho_g}{\beta}\right) \left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \leq 1, \quad (165)$$

which is

$$\beta < \frac{\rho_g \left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right)}{\frac{1}{\alpha} (8L_{y,c}^2 + 2(m+1)L_y^2) + \frac{\vartheta L_{xy} D_x}{4\alpha^2}} = \frac{\rho_g (\alpha^2 + 2(m+1)n\alpha L_y^2 + \vartheta L_{xy} D_x / 4)}{2\alpha ((m+1)L_y^2 + 4L_{y,c}^2) + \vartheta L_{xy} D_x / 4} \quad (166)$$

$$\leq \frac{\rho_g (\alpha + 2(m+1)nL_y^2 + \frac{\vartheta L_{xy} D_x}{4\alpha})}{2(L_y^2 + 4L_{y,c}^2)}. \quad (167)$$

From (164) and (167), we can have constants $C_1 > 0$ and $C_2 < 0$ such that

$$\begin{aligned}
\mathcal{P}_c^{r+1} - \mathcal{P}_c^r & \leq -\frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r))\|^2 - C_1 \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 + \frac{m+1}{n} C_2 \sum_{k=1}^m \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)\|^2 + N_2 \\
& \quad + \left(\frac{2(m+1)L_f^2}{\alpha n^2} + \frac{10(m+1)}{\alpha n^2} \left(\frac{C_{xy} L_{f,0} L_{g,2}}{\mu_g^2} \right)^2 + \frac{8(m+1)L_{y,c}^2 L_y^2}{\alpha n} \right) \|\mathbf{x}^r - \mathbb{1}\bar{\mathbf{x}}^r\|^2 \\
& \quad + \left(\frac{4(mL_{f,y}^2 + (m+1)L_{y,c}^2)}{\alpha n^2} \right) \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 + \frac{2mL_{f,y}^2}{n^2\alpha} \sum_{k=1}^m \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2. \quad (168)
\end{aligned}$$

Note that the consensus errors have been quantified in (138) and (149). Let $\tilde{\mathbf{x}} \triangleq \mathbf{x} - \mathbb{1}\bar{\mathbf{x}}$. Then, we can have

$$\tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A}) \|\mathbf{x} - \mathbb{1}\bar{\mathbf{x}}\|^2 \leq \|\mathbf{A}(\mathbf{x} - \mathbb{1}\bar{\mathbf{x}})\|^2 = \tilde{\mathbf{x}}^T \mathbf{A}^T \mathbf{A} \tilde{\mathbf{x}} = \sum_{i,j} \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|^2 \quad (169)$$

$$= \sum_{j \in \mathcal{N}_i, \forall i} \|\mathbf{x}_i - \mathbf{x}_j\|^2 = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} = \|\mathbf{A}\mathbf{x}\|^2, \quad (170)$$

Applying the telescoping sum, we have

$$\begin{aligned}
& \frac{1}{T} \sum_{r=1}^T \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r))\|^2 \\
& \leq 2\alpha \frac{\mathcal{P}_c^1 - \mathcal{P}_c}{T} + \frac{4mL_{f,y}^2}{n^2} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \left(\frac{8(mL_{f,y}^2 + (m+1)L_{y,c}^2)}{n^2} \right) \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 \\
& + 2 \left(\frac{2(m+1)L_f^2}{n^2} + \frac{10(m+1)}{n^2} \left(\frac{C_{xy}L_{f,0}L_{g,2}}{\mu_g^2} \right)^2 + \frac{8(m+1)L_{y,c}L_y^2}{n} \right) \|\mathbf{x}^r - \mathbb{1}\bar{\mathbf{x}}^r\|^2 + 2\alpha N_2 \\
& \sim \mathcal{O} \left(\frac{1}{\sqrt{nT}} \right), \tag{171}
\end{aligned}$$

since $N_2 \sim \mathcal{O}(1/(n\alpha^2))$, where $\underline{\mathcal{P}}_c$ denotes the lower bound of \mathcal{P}_c .

B.5.4 Optimality of the LL Optimization Variables

From (165) and condition $\beta \geq 2(\mu_g + L_{g,1})^{-1}$, we know that $-1 < C_2 < 0$. Combining (80) and (81) with $\theta' = 2(m+1)nL_y^2/\alpha$, so we can have

$$\begin{aligned}
& (1 + C_2) \mathbb{E} \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \\
& \leq \mathbb{E} \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 - \mathbb{E} \|\bar{\mathbf{y}}_k^{r+1} - \bar{\mathbf{y}}_k^*(\mathbf{x}^{r+1})\|^2 + \left(L_y^2 + \frac{\alpha}{8(m+1)n} + \frac{L_{xy}}{4\theta} \right) \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2 \\
& + \left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy}D_x}{4\alpha^2} \right) \frac{\sigma_g^2}{n\beta^2} + \left(\frac{L_{xy}}{4\theta} + L_y^2 \right) \frac{\sigma_f^2}{n\alpha^2}. \tag{172}
\end{aligned}$$

Applying the telescoping sum on (172), we have

$$\begin{aligned}
& \frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \\
& \leq \frac{\|\bar{\mathbf{y}}_k^1 - \bar{\mathbf{y}}_k^*(\mathbf{x}^1)\|^2}{T(1+C_2)} + \frac{1}{1+C_2} \left(L_y^2 + \frac{\alpha}{8(m+1)n} + \frac{L_{xy}}{4\theta} \right) \frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 \\
& + \frac{1}{1+C_2} \left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy}D_x}{4\alpha^2} \right) \frac{\sigma_g^2}{n\beta^2} + \frac{1}{1+C_2} \left(\frac{L_{xy}}{4\theta} + L_y^2 \right) \frac{\sigma_f^2}{n\alpha^2} \tag{173} \\
& \leq \frac{\|\bar{\mathbf{y}}_k^1 - \bar{\mathbf{y}}_k^*(\mathbf{x}^1)\|^2}{T(1+C_2)} + \frac{1}{1+C_2} \left(L_y^2 + \frac{\alpha}{8(m+1)n} + \frac{L_{xy}}{4\theta} \right) \frac{D_x}{\alpha^2} \\
& + \frac{1}{1+C_2} \left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy}D_x}{4\alpha^2} \right) \frac{\sigma_g^2}{n\beta^2} + \frac{1}{1+C_2} \left(\frac{L_{xy}}{4\theta} + L_y^2 \right) \frac{\sigma_f^2}{n\alpha^2} \sim \mathcal{O} \left(\frac{1}{\sqrt{nT}} \right), \tag{174}
\end{aligned}$$

since we choose $\alpha \sim \beta \sim \mathcal{O}(\sqrt{T/n})$. \square

B.6 Proof of Corollary 1

In this case, we have

$$\mathbf{x}^{r+1} = \mathbf{x}^r - \frac{\mathbf{h}_f^r}{\alpha}. \tag{175}$$

Then, Lemma 4 holds because

$$\mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 \leq \frac{2}{\alpha^2} (\|\mathbf{h}_f^r - \mathbb{E}\mathbf{h}_f^r\|^2 + \|\mathbb{E}\mathbf{h}_f^r\|^2) \leq n \overbrace{\frac{2\sigma_f^2 + 2L_{f,0}}{\alpha^2}}^{\triangleq D_x}. \tag{176}$$

B.6.1 Consensus Violation of the LL Optimization Variables

From (176), we know that $\mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 \leq nD_x/\alpha^2$. Following (139) to (149), we can still have $\frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\mathbf{A}\mathbf{y}_k^r\|^2 \sim \mathcal{O} \left(\frac{1}{\sqrt{nT}} \right)$ when $\rho \geq n$.

B.6.2 Stationarity of the UL Optimization Variables

Note that there will be n objective functions corresponding to $\{\mathbf{x}_i, \forall n\}$. We will show the convergence of SLAM-L for each one of them. According to the gradient Lipschitz continuity and by following (156), we have

$$\begin{aligned}
& \sum_{i=1}^n \mathbb{E}_{\xi^r} [f_i(\mathbf{x}_i^{r+1}, \mathbf{y}_{i,k}^*(\mathbf{x}_i^{r+1})) | \mathcal{F}^r] \\
& \leq \sum_{i=1}^n f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)) + \mathbb{E} [\langle \nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)), \mathbf{x}_i^{r+1} - \mathbf{x}_i^r \rangle | \mathcal{F}^r] + \mathbb{E} \left[\frac{L_{f,1}}{2} \|\mathbf{x}_i^{r+1} - \mathbf{x}_i^r\|^2 | \mathcal{F}^r \right] \\
& = \sum_{i=1}^n f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)) + \langle \nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)), \mathbb{E}[\mathbf{x}_i^{r+1} - \mathbf{x}_i^r] | \mathcal{F}^r \rangle + \frac{L_{f,1}}{2} \|\mathbb{E}[\mathbf{x}_i^{r+1} - \mathbf{x}_i^r]\|^2 \\
& \quad + \frac{L_{f,1}}{2} \mathbb{E} [\|\mathbf{x}_i^{r+1} - \mathbf{x}_i^r - \mathbb{E}[\mathbf{x}_i^{r+1} - \mathbf{x}_i^r]\|^2 | \mathcal{F}^r] \tag{177}
\end{aligned}$$

$$\begin{aligned}
& \stackrel{(a)}{\leq} \sum_{i=1}^n f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)) - \frac{1}{2\alpha} \|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r))\|^2 - \left(\frac{\alpha}{2} - \frac{L_{f,1}}{2} \right) \|\mathbb{E}[\mathbf{x}_i^{r+1} - \mathbf{x}_i^r]\|^2 \\
& \quad + \frac{1}{2\alpha} \|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)) - [\mathbf{h}_f^r]_i\|^2 + \frac{L_{f,1}}{2\alpha^2} \mathbb{E} [\|[\mathbf{h}_f^r]_i - \mathbb{E}[\mathbf{h}_f^r]_i\|^2 | \mathcal{F}^r] \tag{178}
\end{aligned}$$

$$\begin{aligned}
& \stackrel{(b)}{\leq} \sum_{i=1}^n f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)) - \frac{1}{2\alpha} \|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r))\|^2 - \left(\frac{\alpha}{2} - \frac{L_{f,1}}{2} \right) \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2 \\
& \quad + \frac{2n(m+1)}{\alpha} \sum_{k=1}^m L_{f,l}^2 \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 + \frac{2mL_{f,y}^2}{\alpha} \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 \\
& \quad + \frac{2mL_{f,y}^2}{\alpha} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{2b_r^2}{\alpha} + \frac{nL_{f,1}\sigma_f^2}{2\alpha^2} \tag{179}
\end{aligned}$$

where (a) is true because

$$\begin{aligned}
& \langle \nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)), \mathbb{E}[\mathbf{x}_i^{r+1} - \mathbf{x}_i^r] | \mathcal{F}^r \rangle = -\frac{1}{\alpha} \langle \nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)), \mathbb{E}[\mathbf{h}_f^r]_i \rangle \\
& = -\frac{1}{2\alpha} \|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r))\|^2 - \frac{1}{2\alpha} \|\mathbb{E}[\mathbf{h}_f^r]_i\|^2 + \frac{1}{2\alpha} \|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)) - \mathbb{E}[\mathbf{h}_f^r]_i\|^2, \tag{180}
\end{aligned}$$

and in (b) we apply the following facts,

$$\begin{aligned}
& \sum_{i=1}^n \|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)) - \mathbb{E}[\mathbf{h}_f^r]_i\|^2 \\
& \leq \sum_{i=1}^n 4\|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)) - \bar{\nabla} f_i(\mathbf{x}_i^r, \bar{\mathbf{y}}_k^r)\|^2 + 4\|\bar{\nabla} f_i(\mathbf{x}_i^r, \bar{\mathbf{y}}_k^r) - \bar{\nabla} f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^r)\|^2 \\
& \quad + 4\|\bar{\nabla} f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^r) - \bar{\nabla} f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*)\|^2 + 4\|\bar{\nabla} f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*) - \mathbb{E}[\mathbf{h}_f^r]_i\|^2 \tag{181}
\end{aligned}$$

$$\begin{aligned}
& \leq 4n(m+1) \sum_{k=1}^m L_{y,l}^2 \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 + 4mL_{f,y}^2 \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 \\
& \quad + 4mL_{f,y}^2 \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + 4b_r^2, \tag{182}
\end{aligned}$$

and constant $L_{y,l}$ is computed as follows:

$$\begin{aligned}
& \sum_{i=1}^n \|\bar{\nabla} f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k) - \nabla f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k^*(\mathbf{x}))\|^2 \\
& \leq (m+1) \sum_{k=1}^m \sum_{i=1}^n \|\nabla_{\mathbf{x}} f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k) - \nabla_{\mathbf{x}} f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k^*(\mathbf{x}))\|^2 + (m+1) \left\| \Delta_{i,k}^{(1)} + \Delta_{i,k}^{(2)} + \Delta_{i,k}^{(3)} \right\|^2 \\
& \leq n(m+1) \sum_{k=1}^m L_{f,1}^2 + 3 \underbrace{\left(\left(\frac{L_{f,0}L_{g,2}}{\mu_g} \right)^2 + \left(\frac{L_{f,0}C_{xy}L_g}{\mu_g^2} \right)^2 + \left(\frac{C_{xy}L_{f,1}}{\mu_g} \right)^2 \right)}_{\triangleq L_{y,l}^2} \|\bar{\mathbf{y}}_k - \bar{\mathbf{y}}_k^*(\mathbf{x})\|^2,
\end{aligned}$$

and terms $\Delta_{i,k}'^{(1)}, \Delta_{i,k}'^{(2)}, \Delta_{i,k}'^{(3)}$ are defined as

$$\begin{aligned}\Delta_{i,k}'^{(1)} &\triangleq \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \bar{\mathbf{y}}_k) \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \bar{\mathbf{y}}_k) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k) \\ &\quad - \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \bar{\mathbf{y}}_k^*(\mathbf{x})) \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \bar{\mathbf{y}}_k) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k), \\ \Delta_{i,k}'^{(2)} &\triangleq \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \bar{\mathbf{y}}_k^*(\mathbf{x})) \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \bar{\mathbf{y}}_k) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k) \\ &\quad - \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \bar{\mathbf{y}}_k^*(\mathbf{x})) \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}, \bar{\mathbf{y}}_k^*(\mathbf{x})) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k), \\ \Delta_{i,k}'^{(3)} &\triangleq \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \bar{\mathbf{y}}_k^*(\mathbf{x})) \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}, \bar{\mathbf{y}}_k^*(\mathbf{x})) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k) \\ &\quad - \nabla_{\mathbf{x}_i \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}_i, \bar{\mathbf{y}}_k^*(\mathbf{x})) \left[\nabla_{\mathbf{y}_{i,k} \mathbf{y}_{i,k}}^2 g_{i,k}(\mathbf{x}, \bar{\mathbf{y}}_k^*(\mathbf{x})) \right]^{-1} \nabla_{\mathbf{y}_{i,k}} f_i(\mathbf{x}_i, \bar{\mathbf{y}}_k^*(\mathbf{x})).\end{aligned}$$

Let us define potential function as

$$\mathcal{P}_l^r = \sum_{i=1}^n \mathbb{E} [f_i(\mathbf{x}_i^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) | \mathcal{F}^r] + n(m+1) \sum_{k=1}^m \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2. \quad (183)$$

Then, we can have

$$\begin{aligned}&\mathcal{P}_l^{r+1} - \mathcal{P}_l^r \\ &\leq -\frac{1}{2\alpha} \sum_{i=1}^n \|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r))\|^2 \\ &\quad - \underbrace{\left(\frac{\alpha}{2} - \left(\frac{L_{f,1}}{2} + n(m+1)L_y^2 + \frac{\alpha}{8} + \frac{n(m+1)L_{xy}}{4\vartheta} \right) \right)}_{\triangleq C_1} \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2 \\ &\quad + n(m+1) \underbrace{\left(\frac{2L_{y,l}^2}{\alpha} + \left(1 - \frac{\rho_g}{\beta}\right) \left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) - 1 \right)}_{\triangleq C_2} \sum_{k=1}^m \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \\ &\quad + \frac{2mL_{f,y}^2}{\alpha} \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 + \frac{2mL_{f,y}^2}{\alpha} \sum_{k=1}^m \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \\ &\quad + n(m+1) \underbrace{\left(\left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \frac{\sigma_g^2}{n\beta^2} + \left(\frac{L_{xy}}{4\vartheta} + L_y^2\right) \frac{\sigma_f^2}{n\alpha^2} \right) + \frac{2b_r^2}{\alpha} + \frac{nL_{f,1}\sigma_f^2}{2\alpha^2}}_{\triangleq N_3 \sim \mathcal{O}\left(\frac{n}{\alpha^2}\right)}.\end{aligned}$$

1) to show $C_1 > 0$: we need

$$\frac{\alpha}{2} - \left(\frac{L_{f,1}}{2} + n(m+1)L_y^2 + \frac{\alpha}{8} + \frac{n(m+1)L_{xy}}{4\vartheta} \right) > 0, \quad (184)$$

which requires

$$\alpha > 4 \left(\frac{L_{f,1}}{2} + n(m+1)L_y^2 + \frac{n(m+1)L_{xy}}{4\vartheta} \right). \quad (185)$$

2) to show $C_2 < 0$: we need

$$\frac{2L_{y,l}^2}{\alpha} + \left(1 - \frac{\rho_g}{\beta}\right) \left(1 + \frac{2(m+1)nL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \leq 1, \quad (186)$$

which is

$$\beta < \frac{\rho_g(\alpha + 2(m+1)nL_y^2 + \frac{\vartheta L_{xy} D_x}{4\alpha})}{2(L_y^2 + L_{y,l}^2)}. \quad (187)$$

Then, we can have constants $C_1 > 0$ and $C_2 < 0$ such that

$$\begin{aligned}\mathcal{P}_l^{r+1} - \mathcal{P}_l^r &\leq -\frac{1}{2\alpha} \|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r))\|^2 - C_1 \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2 + n(m+1)C_2 \sum_{k=1}^m \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \\ &\quad + \frac{2mL_{f,y}^2}{\alpha} \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbb{1}\bar{\mathbf{y}}_k^r\|^2 + \frac{2mL_{f,y}^2}{\alpha} \sum_{k=1}^m \mathbb{E}\|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + N_3.\end{aligned} \quad (188)$$

Applying the telescoping sum, we can obtain

$$\begin{aligned} & \frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r))\|^2 \\ & \stackrel{(40b)}{\leq} \frac{2\alpha(\mathcal{P}_l^1 - \underline{\mathcal{P}}_l)}{T} + \frac{4mL_{f,y}^2}{T\tilde{\sigma}_{\min}^2(\mathbf{A}^T\mathbf{A})} \sum_{r=1}^T \sum_{i=1}^m \|\mathbf{A}\mathbf{y}_i^r\|^2 + 4m^2L_{f,y}^2 \frac{D_y}{\beta^2} + \alpha N_3 \sim \mathcal{O}\left(\frac{n}{\sqrt{T}}\right), \end{aligned} \quad (189)$$

where $\underline{\mathcal{P}}_l$ denotes the lower bound of \mathcal{P}_l , which gives the result shown in Corollary 1.

B.6.3 Optimality of the LL Optimization Variables

We know that $-1 < C_2 < 0$ based on (187) and condition $\beta \geq 2(\mu_g + L_{g,1})^{-1}$, and $\mathbb{E}_{\xi^r} \left[\left\| \frac{1}{\alpha} (\mathbb{E}\mathbf{h}_f^r - \mathbf{h}_f^r) \right\|^2 | \mathcal{F}^r \right] \leq n\sigma_f^2/\alpha$ in Lemma 5. Following (172) to (174) and choosing $\theta' = 2(m+1)n^{3/2}L_y^2/\alpha$, we have

$$\frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\mathbf{x}^r)\|^2 \sim \mathcal{O}\left(\frac{1}{\sqrt{nT}}\right), \quad (190)$$

when $\alpha \sim \mathcal{O}(\sqrt{T})$, $\beta \sim \mathcal{O}(\sqrt{T/n})$, and $T \gg n$.

B.7 Proof of Corollary 2

In this case, (5a) reduces to

$$\mathbf{y}_k^{r+1} = \mathbf{y}_k^r - \frac{\mathbf{h}_g^r}{\beta} \quad (191)$$

where \mathbf{h}_g^r is the stochastic estimate of $[g_{1,k}(\mathbf{x}_1^r, \mathbf{y}_{1,k}^r), \dots, g_{n,k}(\mathbf{x}_n^r, \mathbf{y}_{n,k}^r)]^T$.

B.7.1 Consensus Violation of the UL Variables

Following (138), we have $\frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\mathbf{A}^T \mathbf{x}^r\|^2 \leq \mathcal{O}(1/\sqrt{nT})$, where $\alpha \sim \mathcal{O}(\sqrt{T/n})$ and $\beta \sim \mathcal{O}(\sqrt{T/n})$.

B.7.2 Stationarity of the UL Optimization Variables

Based on (191), (80) becomes

$$\mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^*(\mathbf{x}^r)\|^2 \leq \left(1 - \frac{\rho_g}{\beta}\right) \mathbb{E} \|\mathbf{y}_k^r - \mathbf{y}_k^*(\mathbf{x}^r)\|^2 + \frac{n\sigma_g^2}{\beta^2}, \quad (192)$$

and there exist $\theta', \vartheta > 0$ such that

$$\begin{aligned} \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^*(\mathbf{x}^{r+1})\|^2 & \leq \left(1 + \theta' + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^*(\mathbf{x}^r)\|^2 \\ & \quad + \left(L_y^2 + \frac{L_y^2}{4\theta'} + \frac{L_{xy}}{4\vartheta}\right) \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 + \left(\frac{L_{xy}}{4\vartheta} + L_y^2\right) \frac{\sigma_f^2}{n\alpha^2}. \end{aligned} \quad (193)$$

Therefore, (161) becomes

$$\begin{aligned} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^*(\mathbf{x}^{r+1})\|^2 & \leq \left(1 - \frac{\rho_g}{\beta}\right) \left(1 + \frac{2mnL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \|\mathbf{y}_k^r - \mathbf{y}_k^*(\mathbf{x}^r)\|^2 \\ & \quad + \left(L_y^2 + \frac{\alpha}{8mn} + \frac{L_{xy}}{4\vartheta}\right) \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 \\ & \quad + \left(1 + \frac{2mnL_f^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \frac{n\sigma_g^2}{\beta^2} + \left(\frac{L_{xy}}{4\vartheta} + L_y^2\right) \frac{\sigma_f^2}{n\alpha^2} \end{aligned} \quad (194)$$

where $\theta' = 2mnL_y^2/\alpha$.

According to the gradient Lipschitz continuity of the UL objective function, we can have

$$\begin{aligned} & \mathbb{E}_{\xi^r} [f(\bar{\mathbf{x}}^{r+1}, \mathbf{y}_k^*(\bar{\mathbf{x}}^{r+1})) | \mathcal{F}^r] \\ & \stackrel{(a)}{\leq} f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)) + \mathbb{E} [\langle \nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)), \bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r \rangle | \mathcal{F}^r] + \mathbb{E} \left[\frac{L_f}{2} \|\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r\|^2 | \mathcal{F}^r \right] \end{aligned} \quad (195)$$

$$\begin{aligned} & = f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)) + \langle \nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)), \mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r] | \mathcal{F}^r \rangle \\ & \quad + \frac{L_f}{2} \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 + \frac{L_f}{2} \mathbb{E} [\|\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r - \mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 | \mathcal{F}^r] \end{aligned} \quad (196)$$

$$\begin{aligned} & \stackrel{(b)}{\leq} f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)) - \frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r))\|^2 - \left(\frac{\alpha}{2} - \frac{L_f}{2} \right) \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 \\ & \quad + \frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)) - \bar{\mathbf{h}}_f^r\|^2 + \frac{L_f}{2\alpha^2} \mathbb{E} [\|\bar{\mathbf{h}}_f^r - \mathbb{E}\bar{\mathbf{h}}_f^r\|^2 | \mathcal{F}^r] \end{aligned} \quad (197)$$

$$\begin{aligned} & \stackrel{(c)}{\leq} f(\bar{\mathbf{x}}^r, \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)) - \frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r))\|^2 - \left(\frac{\alpha}{2} - \frac{L_f}{2} \right) \|\mathbb{E}[\bar{\mathbf{x}}^r - \bar{\mathbf{x}}]\|^2 + \frac{2L_{f,x}^2 m^2}{n^2 \alpha} \|\mathbf{x}^r - \mathbb{1}\bar{\mathbf{x}}\|^2 \\ & \quad + \frac{2mL_{y,u}^2}{n^2 \alpha} \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbf{y}_k^*(\bar{\mathbf{x}}^r)\|^2 + \frac{2mL_{f,y}^2}{n^2 \alpha} \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{2b_r^2}{n\alpha} + \frac{L_f \sigma_f^2}{2n\alpha^2} \end{aligned}$$

where (a) follows the gradient Lipschitz continuity of the UL loss function with constant L_f at the consensus space, (b) is true because

$$\begin{aligned} & \langle \nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)), \mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r] | \mathcal{F}^r \rangle \\ & = -\frac{1}{\alpha} \langle \nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)), \mathbb{E}\bar{\mathbf{h}}_f^r \rangle \\ & = -\frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r))\|^2 - \frac{1}{2\alpha} \|\mathbb{E}\bar{\mathbf{h}}_f^r\|^2 + \frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)) - \mathbb{E}\bar{\mathbf{h}}_f^r\|^2, \end{aligned}$$

and (c) we can apply the following facts,

$$\begin{aligned} & \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)) - \mathbb{E}\bar{\mathbf{h}}_f^r\|^2 \\ & \leq 4\|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)) - \bar{\nabla} f(\bar{\mathbf{x}}^r, \mathbf{y}_k^r)\|^2 + 4\|\bar{\nabla} f(\bar{\mathbf{x}}^r, \mathbf{y}_k^r) - \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^r)\|^2 \\ & \quad + 4\|\bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^r) - \bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1})\|^2 + 4\|\bar{\nabla} f(\mathbf{x}^r, \mathbf{y}_k^{r+1}) - \mathbb{E}\bar{\mathbf{h}}_f^r\|^2 \end{aligned} \quad (198)$$

$$\begin{aligned} & \leq \frac{4mL_{y,u}^2}{n^2} \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbf{y}_k^*(\bar{\mathbf{x}}^r)\|^2 + \frac{4L_{f,x}^2 m^2}{n^2} \|\mathbf{x}^r - \mathbb{1}\bar{\mathbf{x}}\|^2 \\ & \quad + \frac{4mL_{f,y}^2}{n^2} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + \frac{4b_r^2}{n}, \end{aligned} \quad (199)$$

and the continuity of $\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r))$ is the same as the centralized case (e.g., constant C in [31, Lemma 2.2.]) as variable $\mathbf{y}_{i,k}, \forall i$ are decoupled over the nodes, namely $\|\nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k}^*(\mathbf{x}_i^r)) - \nabla f_i(\mathbf{x}_i^r, \mathbf{y}_{i,k})\| \leq L_{y,u} \|\mathbf{y}_{i,k}^*(\mathbf{x}_i^r) - \mathbf{y}_{i,k}\|, \forall i, k$.

Define potential function as

$$\mathcal{P}_u^r = \mathbb{E} [f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r)) | \mathcal{F}^r] + \frac{m}{n^2} \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbf{y}_k^*(\bar{\mathbf{x}}^r)\|^2. \quad (200)$$

Combining (193), we have Then, we can have

$$\begin{aligned}
& \mathcal{P}_u^{r+1} - \mathcal{P}_u^r \\
& \leq -\frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r))\|^2 - \underbrace{\left(\frac{\alpha}{2} - \left(\frac{L_f}{2} + \frac{m}{n^2} L_y^2 + \frac{\alpha}{8n^3} + \frac{mL_{xy}}{4n^2\vartheta} \right) \right)}_{\triangleq C_1} \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 \\
& \quad + \frac{m}{n^2} \underbrace{\left(\frac{2L_{y,u}^2}{\alpha} + \left(1 - \frac{\rho_g}{\beta}\right) \left(1 + \frac{2nmL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) - 1 \right)}_{\triangleq C_2} \sum_{k=1}^m \|\mathbf{y}_k^r - \mathbf{y}_k^*(\bar{\mathbf{x}}^r)\|^2 \\
& \quad + \frac{2L_{f,x}^2 m^2}{n^2 \alpha} \|\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r\|^2 + \frac{2mL_{f,y}^2}{n^2 \alpha} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 \\
& \quad + \underbrace{\frac{m}{n^2} \left(\left(1 + \frac{2mnL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2}\right) \frac{n\sigma_g^2}{\beta^2} + \left(\frac{L_{xy}}{4\vartheta} + L_y^2\right) \frac{\sigma_f^2}{n\alpha^2} \right) + \frac{2b_r^2}{n\alpha} + \frac{L_f \sigma_f^2}{2n\alpha^2}}_{\triangleq N_4 \sim \mathcal{O}\left(\frac{1}{n\alpha^2}\right)}.
\end{aligned}$$

Following (164) and (167), we can have constants $C_1 > 0$ and $C_2 < 0$ such that

$$\begin{aligned}
\mathcal{P}_u^{r+1} - \mathcal{P}_u^r & \leq -\frac{1}{2\alpha} \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r))\|^2 - C_1 \|\mathbb{E}[\bar{\mathbf{x}}^{r+1} - \bar{\mathbf{x}}^r]\|^2 + \frac{m}{n^2} C_2 \sum_{k=1}^m \|\bar{\mathbf{y}}_k^r - \bar{\mathbf{y}}_k^*(\bar{\mathbf{x}}^r)\|^2 \\
& \quad + \frac{2L_{f,x}^2 m^2}{n^2 \alpha} \|\mathbf{x}^r - \mathbb{1} \bar{\mathbf{x}}^r\|^2 + \frac{2mL_{f,y}^2}{n^2 \alpha} \sum_{k=1}^m \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^r\|^2 + N_4.
\end{aligned} \tag{201}$$

Applying the telescoping sum, we have that when $\alpha \sim \mathcal{O}(\sqrt{T/n})$ and $\beta \sim \mathcal{O}(\sqrt{T/n})$,

$$\begin{aligned}
\frac{1}{T} \sum_{r=1}^T \|\nabla f(\bar{\mathbf{x}}^r, \mathbf{y}_k^*(\bar{\mathbf{x}}^r))\|^2 & \leq 2\alpha \frac{\mathcal{P}_u^1 - \mathcal{P}_u}{T} + \frac{4L_{f,x}^2 m^2}{n^2 \tilde{\sigma}_{\min}(\mathbf{A}^T \mathbf{A})} \frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\mathbf{A} \mathbf{x}_k^r\|^2 \\
& \quad + \frac{4m^2 L_{f,y}^2 D_y}{n^2 \beta^2} + 2\alpha N_4 \sim \mathcal{O}\left(\frac{1}{\sqrt{nT}}\right),
\end{aligned} \tag{202}$$

where \mathcal{P}_u denotes the lower bound of \mathcal{P}_u .

B.7.3 Optimality of the LL Optimization Variables

From (193) and condition $\beta \geq 2(\mu_g + L_{g,1})^{-1}$, we know that $-1 < C_2 < 0$. With $\theta' = 2mnL_y^2/\alpha$, so (194) can be written as

$$\begin{aligned}
& (1 + C_2) \mathbb{E} \|\mathbf{y}_k^r - \mathbf{y}_k^*(\mathbf{x}^r)\|^2 \\
& \leq \mathbb{E} \|\mathbf{y}_k^r - \mathbf{y}_k^*(\mathbf{x}^r)\|^2 - \mathbb{E} \|\mathbf{y}_k^{r+1} - \mathbf{y}_k^*(\mathbf{x}^{r+1})\|^2 + \left(L_y^2 + \frac{\alpha}{8mn} + \frac{L_{xy}}{4\vartheta} \right) \|\mathbb{E}[\mathbf{x}^{r+1} - \mathbf{x}^r]\|^2 \\
& \quad + \left(1 + \frac{2mnL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2} \right) \frac{n\sigma_g^2}{\beta^2} + \frac{L_{xy}}{4\vartheta} \frac{\sigma_f^2}{n\alpha^2}.
\end{aligned} \tag{203}$$

Applying the telescoping sum on (203), we have

$$\begin{aligned}
& \frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\mathbf{y}_k^r - \mathbf{y}_k^*(\mathbf{x}^r)\|^2 \\
& \leq \frac{\|\mathbf{y}_k^1 - \mathbf{y}_k^*(\mathbf{x}^1)\|^2}{T(1 + C_2)} + \frac{1}{1 + C_2} \left(L_y^2 + \frac{\alpha}{8mn} + \frac{L_{xy}}{4\vartheta} \right) \frac{1}{T} \sum_{r=1}^T \mathbb{E} \|\mathbf{x}^{r+1} - \mathbf{x}^r\|^2 \\
& \quad + \frac{1}{1 + C_2} \left(1 + \frac{2mnL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2} \right) \frac{n\sigma_g^2}{\beta^2} + \frac{1}{1 + C_2} \frac{L_{xy}}{4\vartheta} \frac{\sigma_f^2}{n\alpha^2}
\end{aligned} \tag{204}$$

$$\begin{aligned}
& \leq \frac{\|\mathbf{y}_k^1 - \mathbf{y}_k^*(\mathbf{x}^1)\|^2}{T(1 + C_2)} + \frac{1}{1 + C_2} \left(L_y^2 + \frac{\alpha}{8mn} + \frac{L_{xy}}{4\vartheta} \right) \frac{D_x}{\alpha^2} \\
& \quad + \frac{1}{1 + C_2} \left(1 + \frac{2mnL_y^2}{\alpha} + \frac{\vartheta L_{xy} D_x}{4\alpha^2} \right) \frac{n\sigma_g^2}{\beta^2} + \frac{1}{1 + C_2} \frac{L_{xy}}{4\vartheta} \frac{\sigma_f^2}{n\alpha^2} \sim \mathcal{O}\left(\frac{1}{\sqrt{nT}}\right),
\end{aligned} \tag{205}$$

as we choose $\alpha \sim \mathcal{O}(\sqrt{T/n})$, $\beta \sim \mathcal{O}(\sqrt{T/n})$, and $T \gg n$.

C Additional Numerical Results

C.1 SLAM for MARL: SLAM-AC

In this section, we will introduce how to use SLAM for solving MARL problems. Consider a networked Markov Decision Process (nMDP) $(\mathcal{S}, \{\mathcal{A}_i, \forall i\}, \mathcal{P}, \{R_i, \forall i\}, \eta, \mathcal{G}, \gamma)$, where \mathcal{S} denotes the global states shared by all agents, \mathcal{A}_i is the action space of agent i , subsequently $\mathcal{A} = \prod_{i=1}^n \mathcal{A}_i$ is the joint action space of all agents, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the state transition probability of the nMDP, $R_i(s, \mathbf{a}) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}, \forall i$ denote the local rewards, $\eta(s)$ denotes the initial state, \mathcal{G} is the communication graph, and $\gamma \in (0, 1)$ stands for the discount factor. We assume that the states s and actions \mathbf{a} are globally observable. The goal of this problem is to learn a joint policy π_θ parametrized by θ such that the networked reward function is maximized. Let θ_i be the local policy at each agent and the concatenation of all local policies be $\theta = [\theta_1, \dots, \theta_n]^T$ (which will be the UL optimization variables in the DAC formulation). Under this setting, $\mu_\theta(s)$ is the stationary distribution induced by the policy π_θ at each state, and $d_\theta(s) = (1 - \gamma) \sum_{r=0}^{\infty} \gamma^r \mathcal{P}(s^r = s | s^0 \sim \eta(s))$ denotes the discounted visitation measure. Therefore, given the initial state $\eta(s)$, the policy π_θ can generate a trajectory according to the nMDP. Then, the discounted accumulative reward function maximization problem w.r.t. optimizing the policy is $\max_\theta J(\theta)$, and the objective function is

$$J(\theta) \triangleq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\pi_\theta} \left[\sum_{r=0}^{\infty} \gamma^r R_i(s^r, \mathbf{a}^r) \right] \quad (206)$$

$$= \mathbb{E}_{\pi_\theta} \left[\sum_{r=0}^{\infty} \gamma^r \bar{R}(s^r, \mathbf{a}^r) \right] \quad (207)$$

$$= \mathbb{E}_{s \sim \eta(\cdot)} [V_{\pi_\theta}(s)] \quad (208)$$

where $\bar{R}(s, \mathbf{a}) = n^{-1} \sum_{i=1}^n R_i(s, \mathbf{a})$, and the expectation is taken over all the trajectories generated by policy π_θ , and value function $V_{\pi_\theta}(s) \triangleq \mathbb{E} [\sum_{r=0}^{\infty} \gamma^r \bar{R}(s^r, \mathbf{a}^r) | s^0 = s]$. Note that given policy π_θ , value function $V_{\pi_\theta}(s)$ satisfies the Bellman equation [7] and can be written as

$$V_{\pi_\theta}(s) = \mathbb{E}_{\mathbf{a} \sim \pi_\theta(\cdot | s), s' \sim \mathcal{P}(\cdot | s, \mathbf{a})} [\bar{R}(s, \mathbf{a}) + \gamma V_{\pi_\theta}(s')], \quad (209)$$

so the policy gradient [50] of $J(\theta)$ w.r.t. θ can be expressed by

$$\nabla J(\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\theta(\cdot), \mathbf{a} \sim \pi_\theta(\cdot | s), s' \sim \mathcal{P}(\cdot | s, \mathbf{a})} [(\bar{R}(s, \mathbf{a}) + \gamma V_{\pi_\theta}(s')) \nabla_\theta \log \pi_\theta(\mathbf{a} | s)]. \quad (210)$$

Based on the finite sum structure of the rewards over the network and policy gradient theorem, it is motivated that the policy gradient at each node can be estimated locally when a consensus process is allowed and the global action and state are observable. As the value function is not explicitly known at each agent, we apply the critic step to approximate the global valuation function by formulating the LL function estimation process. Following the existing AC works [30, 23, 15], we also adopt the linear function approximation, i.e., $\widehat{V}(s, w) = \varphi(s)^T w$, to estimate the global value function, where $\varphi(s)$ denotes the feature mapping and w is the model parameter. Then, in this decentralized setting, the global value functions are estimated by $\widehat{V}(s, w_i) = \varphi(s)^T w_i, \forall i$ and subsequently the global reward functions can be estimated through $\widehat{R}_i(s, \mathbf{a}, \phi_i) = \psi(s, \mathbf{a})^T \phi_i, \forall i$, where $\{w_i, \phi_i, \forall i\}$ are local model parameters and $\psi(s, \mathbf{a})$ is the given feature mapping. Towards this end, the DAC problem can be formulated as the following DBO problem,

$$\max_{\theta} J(\theta) = \mathbb{E}_{\pi_\theta} \left[\sum_{r=0}^{\infty} \gamma^r \bar{R}(s^r, \mathbf{a}^r) \right] \quad (211a)$$

$$\text{s.t. } \theta_i^p = \theta_j^p, \quad j \in \mathcal{N}_i, \forall i \quad (211b)$$

$$w^*(\theta) = \arg \min_w \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left(R_i(s, \mathbf{a}) + \gamma \widehat{V}(s', w_i) - \widehat{V}(s, w_i) \right)^2, \text{ s.t. } w_i = w_j, \quad j \in \mathcal{N}_i, \forall i \quad (211c)$$

$$\phi^*(\theta) = \arg \min_{\phi} \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left(R_i(s, \mathbf{a}) - \widehat{R}_i(s, \mathbf{a}, \phi_i) \right)^2, \text{ s.t. } \phi_i = \phi_j \quad j \in \mathcal{N}_i, \forall i \quad (211d)$$

where θ_i^p is the shared part of the agent i 's policy parameter, $w = [w_1, \dots, w_n]^T$ and $\phi = [\phi_1, \dots, \phi_n]^T$. The LL problem (211c) is used for approximating the network value function and (211d) for the averaged reward function, and both of them are needed at the UL problem (211a) for computing the policy gradient w.r.t. local policies at each iteration, which can be explicitly approximated by

$$\widehat{\nabla}_{\theta_i} J(\theta) = (1 - \gamma)^{-1} (\psi(s^r, \mathbf{a}^r)^T \phi_i^r + \gamma \varphi(s^{r+1})^T w_i^r - \varphi(s^r)^T w_i^r) \nabla_{\theta_i} \log \pi_{\theta_i}(\mathbf{a}_i^r | s^r). \quad (212)$$

Given the above setting, we implement SLAM for solving DAC problem (211) in a decentralized way. The results are shown in both main text and the following. All the experiments are executed on an Apple MacBook Pro (8 GB Memory, M1 processor).

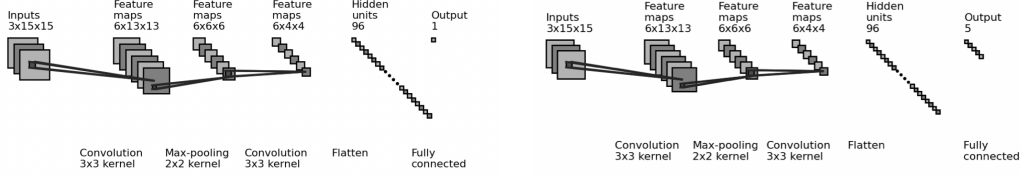


Figure 3: **Neural Network Architecture Diagrams in the Pursuit-Evasion Game.** (Left) The diagram of the **critic** network. (Right) The diagram of the **actor** network.

C.2 Cooperative Navigation Task

In our simulations, the dimension of state s^r is $4 \times n$, since it includes the two-dimensional coordinates of n agents and n landmarks. Both the actor and critic networks maintained at each agent have one hidden layer with 20 neurons followed by the ReLU activation function. Under a common policy $\pi_\theta = \{\pi_{\theta_i}\}_{i=1}^n$, critic networks jointly estimate the global value function $V_\pi(s)$ for all $s \in \mathcal{S}$, where the dimension of the output layer is 1. While the output of actor network corresponds to the probability of choosing each possible action option and the dimension of the output layer is 5 as there are 5 given actions in total. We set the step size for the critic network as 1×10^{-3} and the step size for the actor network as 1×10^{-4} for all the compared algorithms. Moreover, we set $\rho = 1 \times 10^{-1}$ in SLAM-AC. The results have been shown in Figure 1.

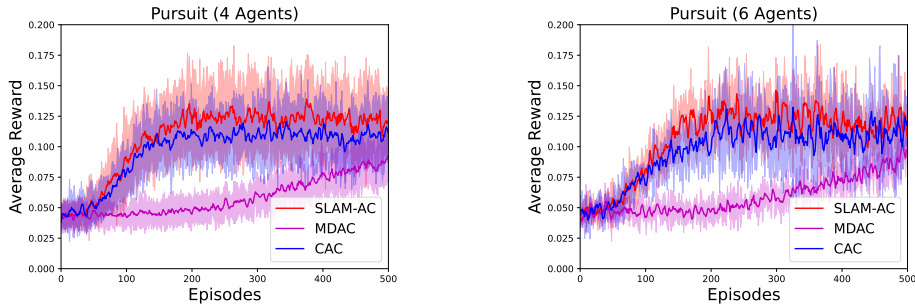


Figure 4: The averaged reward versus the learning process on the pursuit-evasion game. **(With only LL Consensus.)**

C.3 Pursuit-Evasion Game

In this experiment, the “capture” reward for each agent is set to be 5 when a pursuer successfully catches an evader. Additionally, the pursuer will receive a small reward set to be 0.1 when the pursuer encounters an evader at its current location. The environment is divided by a 15×15 grid where there exist obstacles in this two-dimensional (2D) grid such that the agents cannot pass through. Hence, the global state of the pursuit-evasion game consists of three images (binary matrices) with the size of 15×15 . Consequently, the dimension of the global state is $3 \times 15 \times 15$. These three images (binary matrices) respectively present the location of the pursuers, evaders, and obstacles in this 2D grid.

Since the observation of each agent is a 3-channel image, two convolutional neural networks (CNNs) are each agent respectively maintained at each agent, including two convolutional layers, one max-pooling layer, and one fully connected layer for both the actor and critic learners. The ReLU activation function is utilized in each hidden layer of actor network and critic networks. The output of critic network targets approximating the value function $V_\pi(s)$ for all $s \in \mathcal{S}$, where the dimension of the output layer is 1. The output dimension of actor network is 5 which corresponds to the number of possible actions. In each CNN, the raw images (3-channel location matrices), whose dimension is $3 \times 15 \times 15$, are processed by two convolutional layers and one max-pooling layer first and then passed through a fully connected layer as the output layer. The detailed structure diagrams used for the actor and critic networks are shown in Figure 3. For all algorithms, we set the step size for the critic network as 1×10^{-3} and the step size for the actor network as 1×10^{-4} . Moreover, we set $\rho = 5$ in SLAM-AC.

Additional results are shown in Figure 4, where only critic neural networks are used for consensus for estimating the network value function through the communication channel. All the algorithms adopt the same settings. It can be observed from these figures that our proposed algorithm, SLAM-AC, outperforms the state-of-the-art algorithms w.r.t. both convergence speed and achievable average rewards.

C.4 Scalability, robustness, and extendability of SLAM in MARL

In this section, we add more numerical results on MARL and decentralized MAML problems respectively to showcase the superiority of the proposed learning framework and efficiency of SLAM, including scalability, robustness, and extendability.

Based on the problem setup and parameters chosen shown in Section C, we compare the scalability of the three algorithms on a different number of agents and test the performance of SLAM by varying the hyper-parameters.

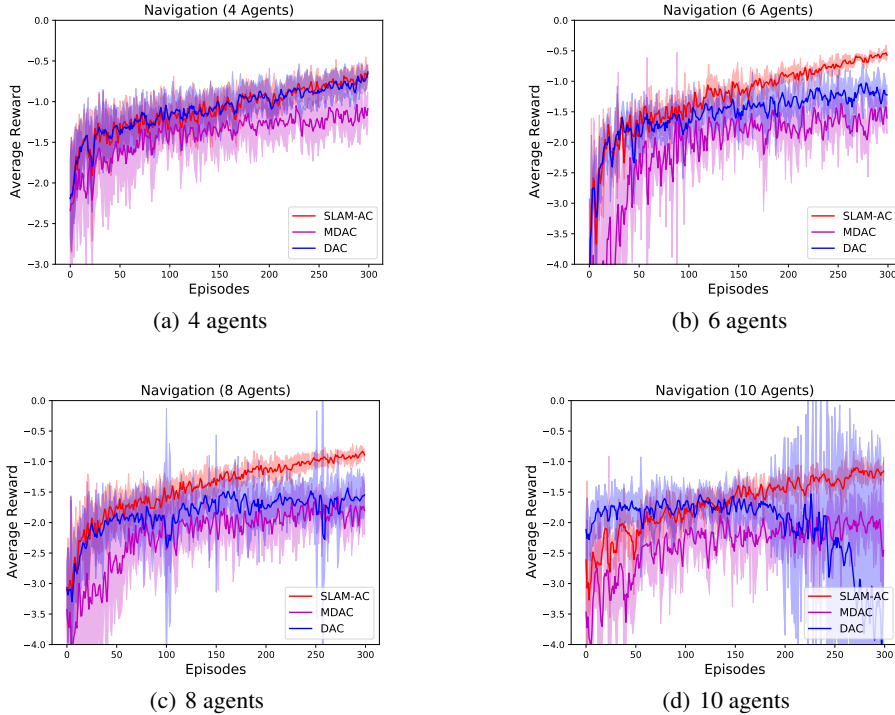


Figure 5: The averaged reward versus the learning process on the cooperative navigation task over different numbers of agents.

C.4.1 Scalability

It is well known that the variance of the multi-agent policy gradient estimate grows as the number of agents increases [18], implying that the difficulty of solving large-scale multi-agent problems rises. We use the same set of step sizes for each algorithm and test their performances by increasing the number of agents from 4 to 10. It can be observed in Figure 5 that SLAM-AC consistently performs well and converges stably with minimum variances compared with the other two existing methods, where the average reward achieved by the classic MARL algorithm, DAC, has higher variances and even diverges when the number of agents is increased to 10, and the reward obtained by MDAC is relatively stable compared to DAC as the multiple consensus steps in the inner loop reduces the variances, but it is still lower than SLAM-AC. Therefore, it is concluded that SLAM-AC scales better than the other two due to the averaged variance over the network.

C.4.2 Effects on Hyper-parameters

Besides the scalability of SLAM, we have also evaluated the numerical performance of SLAM by varying the hyper-parameters, including the actor and critic step sizes and the penalty term. It can be seen in Figure 6 that when the step sizes and penalty term shrink by a half and one-tenth, SLAM-AC converges slower while keeping a similar convergence behavior.

C.4.3 Extendability of SLAM-L with PPO

The proposed SLAM based learning framework is generic, and is amenable to incorporating other optimization techniques at each agent to improve the performance of learning models. In applications of MARL, we further consider two algorithms, denoted as SLAM-PPO and Dec-PPO. Both of them apply proximal policy optimization

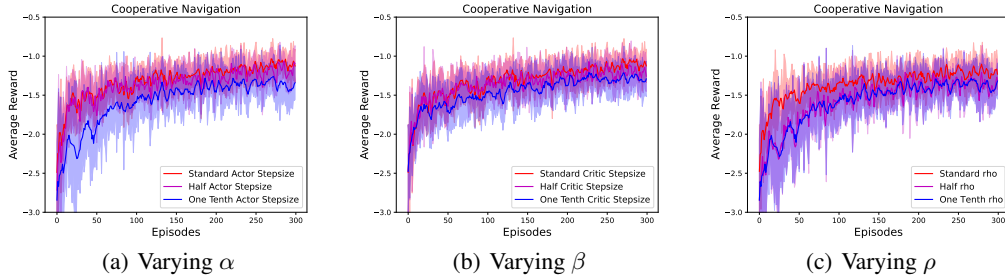


Figure 6: The averaged reward versus the learning process by varying hyper-parameters of SLAM on the cooperative navigation task.

(PPO) instead of the vanilla policy gradient to update their actors (upper-level parameters). In the updates of the critic, the difference is that SLAM-PPO uses SLAM-L to minimize the temporal difference (TD) errors while Dec-PPO directly utilizes the decentralized gradient descent for policy evaluation. We test these two algorithms on the navigation task, where the ϵ clip threshold in PPO is set as 0.2. It can be observed from Figure 7 that SLAM-PPO converges to higher reward values with less variance compared with DEC-PPO, which is similar to the numerical results by employing policy gradient updates for the actor networks.

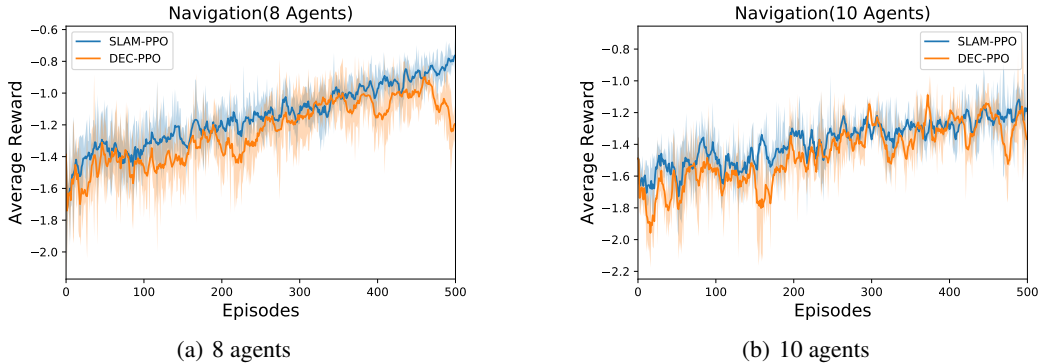


Figure 7: The averaged reward versus the learning process on the cooperative navigation task with PPO for policy improvement.

C.5 Decentralized Meta Learning

We further test the performance of the proposed SLAM with application to a multi-agent MAML problem (9). Following the ANIL structure [6], we partition the weights of the neural network at each agent as two parts denoted by \mathbf{x}_i and \mathbf{y}_i , where the UL optimization problem is to extract the reusable latent space across the connected agents while the LL optimization problem is adopted for adaption of the local model to individual learning tasks.

In this numerical experiment, a two-layer neural network is used, where the numbers of neurons at the hidden and output layers are 32 and 10 respectively and the activation function is sigmoid. A 2-norm regularization term with parameter 0.01 is added to the LL loss function. The communication topology is generated by a random Erdős–Rényi graph. We divide the MNIST dataset as n parts, where each of them only includes 128 data samples for the training. The (standard) initial UL and LL step sizes of SLAM-U are 0.01 and 0.1, and the mini-batch size for gradient estimate is 32.

C.5.1 Linear Speed-up

From the numerical results shown in Figure 8, it can be observed that as the number of agents increases, SLAM-U converges faster in terms of the data samples passed. Setting 93.5% test accuracy as a threshold, we measure the number of data sampled passed w.r.t. the different numbers of agents and report the results in Table 2 and Figure 8(b). If we assume that the computational evaluations of the stochastic gradient estimate are the same at each agent, Figure 8(b) shows that there at least exists a linear speed-up w.r.t. the number of agents in terms of the computational workload, which is consistent with the theoretical analysis.

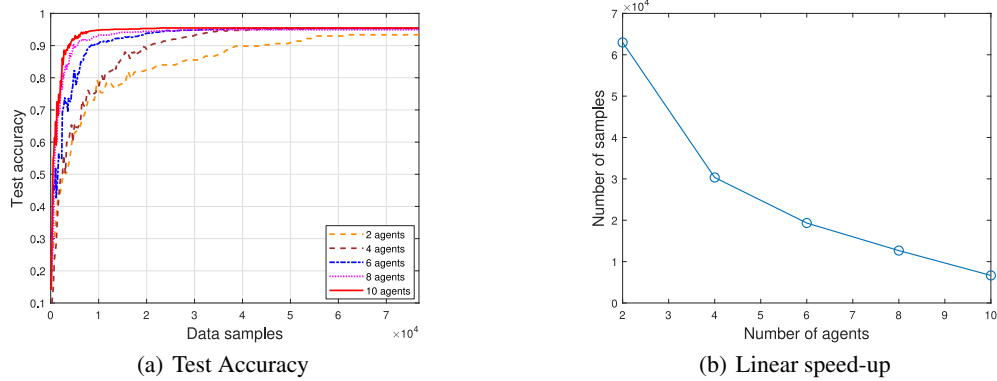


Figure 8: The test accuracy versus the numbers of data samples passed at each agent on the decentralized metal-learning task.

Table 2: Linear speed-up of SLAM-U w.r.t. the number of agents.

number of agents (n)	2	4	6	8	10
required number of samples passed	62976	30336	19328	12672	6656

C.5.2 Effects on Hyper-parameters

Similar to the MARL case, we use different step sizes and penalty parameters to test the robustness of SLAM-U against the changes of these hyper-parameters. From Figure 9, it can be seen that if the step sizes and ρ decrease, SLAM-U converges relatively slower, and the performance of SLAM-U is more sensitive to the UL step size while very robust to the LL step size and ρ .

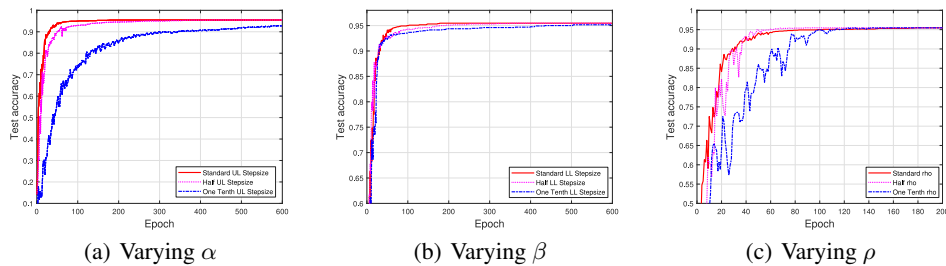


Figure 9: The test accuracy versus the numbers of epochs on the decentralized metal-learning task.