A Regularized Actor-Critic Algorithm for Bi-Level Reinforcement Learning

Sihan Zeng, Sujay Bhatt, Sumitra Ganesh, Alec Koppel

J.P.Morgan AI Research, United States

{sihan.zeng, sujay.bhatt,sumitra.ganesh,alec.koppel}@jpmorgan.com

Abstract

We study a structured bi-level optimization problem where the upper-level objective is a generic smooth function, and the lower-level problem corresponds to policy optimization in a Markov Decision Process (MDP). The decision variable at the upper level parameterizes the reward function of the lower-level MDP, and the upper-level objective is evaluated based on the optimal policy induced by this reward. Such formulations naturally arise in contexts such as reward shaping and reinforcement learning (RL) from human feedback.

Solving this bi-level problem is challenging due to the non-convexity of the lower-level objective and the difficulty of estimating the upper-level hyper-gradient. Existing methods often rely on second-order information, impose strong regularization on the lower-level RL problem, and/or inefficiently use samples through nested-loop procedures. In this work, we propose a single-loop, first-order actor-critic algorithm that optimizes the upper-level objective via a penalty-based reformulation. The algorithm introduces into the lower-level RL objective an entropy regularization with decaying weight, which enables asymptotically unbiased upper-level hyper-gradient estimation without requiring the solution of the exact unregularized lower-level RL problem. Our main contribution is to establish the finite-time and finite-sample convergence of the proposed algorithm to the original, unregularized bi-level optimization problem. We support the theoretical results and numerically validate our method's convergence through simulations in synthetic environments.

1 Introduction

We study bi-level reinforcement learning (RL), a structured bi-level optimization program in which the upper-level decision variable determines the reward function of a lower-level RL problem, and the upper-level objective is evaluated under the lower-level optimal policy. This framework abstracts a wide range of applications where the reward must be tuned to achieve high-level goals while the underlying policy adapts to the reward. Examples include reward shaping [Hu et al., 2020], inverse RL [Zeng et al., 2022b], multi-agent incentive design [Ma et al., 2025], contract design [Zhu et al., 2023], and, notably, reinforcement learning from human feedback (RLHF), one of the central paradigms for fine-tuning large language models (LLMs) [Chakraborty et al., 2024, Ye et al., 2025].

Despite a recent surge of interest in bi-level optimization, bi-level RL remains challenging to solve both in theory and practice. Existing gradient-based approaches to bi-level optimization largely fall into two categories. The first leverages the implicit function theorem to derive the upper-level hyper-gradient [Ghadimi and Wang, 2018] and then applies iterative gradient descent in this direction. However, since the hyper-gradient depends on the Jacobian and Hessian of the lower-level objective, these methods are difficult to apply in bi-level RL, where second-order information either requires oracle access to the transition model or is prohibitively expensive to estimate from trajectory samples.

The second type of methods replaces the lower-level optimality condition with an explicit penalty term [Kwon et al., 2023, Shen and Chen, 2023], enabling an alternative expression of the hypergradient that depends only on first-order information. While this approach bypasses Hessian and Jacobian estimation, existing convergence analyses require strong structural assumptions on the lower-level problem, most commonly strong convexity or the Polyak–Łojasiewicz (PL) condition. In RL, however, these assumptions generally fail: the policy optimization objective may satisfy a weaker "gradient domination" property [Agarwal et al., 2021, Mei et al., 2020], but strong convexity and PL condition do not hold under common policy parameterizations. As a result, existing guarantees in bi-level optimization apply only to regularized or restricted settings of RL, leaving open the question of whether we can provably solve the original, unregularized bi-level RL problem. Moreover, existing penalty-based methods are often implemented in a nested-loop fashion, repeatedly solving the lower-level problem to high precision before each upper-level update to ensure stability. The price of this stability is inefficient use of samples in practice and limited scalability.

Our work addresses this research gap. We advance the second, penalty-based approach in the context of RL and propose a single-loop actor-critic algorithm that provably converges to a stationary point of the bi-level RL objective. The core idea is to enforce the PL condition at the lower level through entropy regularization, and then gradually adjust its weight so that the regularized problem asymptotically recovers the original one. Two important technical innovations enable finite-time and finite-sample analysis for the algorithm. First, we introduce a new technique to tightly characterize the iteration-wise decay of the lower-level optimality error for the regularized RL objective. Combined with the penalty reformulation, this bound allows us to establish convergence of a fully single-loop penalty-based algorithm in the lower-level PL/regularized RL setting, achieving a convergence rate that surpasses the best known rate derived in Kwon et al. [2023] under lower-level strong convexity.

Our goal is to optimize the original, unregularized bi-level RL objective. While a large regularization accelerates the solution of the regularized problem, it also enlarges the discrepancy between the regularized and unregularized optima. Our second innovation addresses this trade-off by dynamically decaying the regularization weight, allowing the algorithm to track the regularized optima as the regularized problem gradually approaches the original one. A key challenge arises here from the time-varying lower-level landscape, which we overcome via a novel multi-time-scale stochastic approximation analysis. Below we detail our innovation and main technical contributions.

Main Contributions

- We establish several fundamental structural properties of bi-level RL, linking the original problem to its entropy-regularized counterpart. In particular, we show that as the regularization weight decays, the optimizer of the entropy-regularized RL objective converges to the unique entropy-maximizing policy within the set of optimizers of the unregularized objective, and we provide a bound on this rate of convergence. Absent in the prior work, this type of structural analysis plays a critical role in justifying the use of the regularized objective as a faithful surrogate for the original formulation, and offers insights of potential independent interest to the broader fields of reinforcement learning and bi-level optimization.
- ullet We present a sample-based, single-loop bi-level RL algorithm and characterize its finite-sample complexity. The algorithm optimizes a regularized bi-level RL objective, while the regularization weight dynamically decays over time. We prove that this algorithm converges to a stationary point of the original bi-level objective with a sample complexity of $\mathcal{O}(\epsilon^{-10})$ through a novel five-time-scale analysis, carefully balancing the regularization weight with the update speed of the dual variable, upper-level decision variable, policy iterates, and value function estimates. To our knowledge, this is the first algorithm that provably solves the unregularized bi-level RL problem, and the first that enjoys finite-time and finite-sample guarantees.
- ullet We show that our proposed algorithm can be instantiated with a constant regularization weight to solve the corresponding regularized bi-level RL problem. In this setting, the algorithm reduces to a single-loop actor-critic algorithm for regularized bi-level RL that relies solely on direct samples from the lower-level MDP. Our finite-sample analysis reveals that the algorithm achieves a complexity of $\mathcal{O}(\epsilon^{-3})$. This rate matches the state-of-the-art complexity of a comparable nested-loop method under lower-level PL condition [Chen et al., 2024], and improves over the $\mathcal{O}(\epsilon^{-3.5})$ complexity derived in Kwon et al. [2023] under strong convexity, which is an even stronger condition.

| | Lower-Level Structure | Single Loop | Only Using First- Order Information | Sample Complexity | Anytime Valid |
|---------------------------|--------------------------|----------------|--|--|------------------|
| Kwon et al. [2023] | Strong convexity | / | ✓ | $\widetilde{\mathcal{O}}(\epsilon^{-3.5})^{\dagger}$ | ✓ |
| Shen and Chen [2023] | PL condition | Х | ✓ | N/A (Iteration complexity derived) | - |
| Xiao and Chen [2024] | PL condition | ✓ | ✓ | N/A (Deterministic setting studied) | - |
| Shen et al. [2024] | Regularized RL | Х | ✓ | N/A (Iteration complexity derived) | - |
| Chakraborty et al. [2024] | Regularized RL | Х | × | N/A (Iteration complexity derived) | - |
| Thoma et al. [2024] | Regularized RL | Х | × | N/A (Iteration complexity derived) | - |
| Xiao et al. [2025] | Strong convexity | Х | ✓ | N/A | × |
| Chen et al. [2024] | PL condition | Х | ✓ | $\widetilde{\mathcal{O}}(\epsilon^{-3})$ | × |
| Yang et al. [2025] | Regularized RL | X | ✓ | $\widetilde{\mathcal{O}}(\epsilon^{-3.5})$ | Х |
| Gaur et al. [2025] | Regularized RL | X | ✓ | $\widetilde{\mathcal{O}}(\epsilon^{-3})$ | Х |
| This work | Regularized RL | ✓ | ✓ | $\widetilde{\mathcal{O}}(\epsilon^{-3})$ | ✓ |
| This work | Original RL | ✓ | ✓ | $\widetilde{\mathcal{O}}(\epsilon^{-10})$ | ✓ |

Table 1: Assumption, structure, and sample complexity (measured by squared gradient norm) of existing algorithms for bi-level optimization and RL. † This is the complexity of the standard F²SA (Fully First-order Stochastic Approximation) algorithm. Kwon et al. [2023] also proposes a momentum-based algorithm (F³SA) that achieves a better rate. For a fair comparison with other non-momentum algorithms in the table, we report the complexity of F²SA. Notably, Yang et al. [2023] further improve the design and analysis of momentum-based algorithms, achieving a complexity of $\mathcal{O}(\epsilon^{-1.5})$, where the key innovation is estimating the lower-level Hessian via finite differences.

1.1 Related Work

Our work relates to the increasing volume of literature on bi-level optimization. Here we discuss the most relevant papers on first-order methods to give context to our contributions, and provide a complete literature comparison with other first- and second-order methods in Table 1.

The penalty reformulations were first introduced in Kwon et al. [2023], Shen and Chen [2023]. The former proposes a fully first-order algorithm and shows that, under lower-level strong convexity, the algorithm converges to a stationary point of the bi-level objective with sample complexity $\widetilde{\mathcal{O}}(\epsilon^{-3.5})$. The latter relaxes the lower-level structure by considering the PL condition. The authors propose a nested-loop algorithm and establish the complexity with respect to the number of outer-loop iteration. However, the complexity with respect to the total number of samples or iterations is unknown.

The penalty reformulation has inspired several works in bi-level RL [Shen et al., 2024, Yang et al., 2025, Gaur et al., 2025], which derive first-order expressions for the hyper-gradient and estimate them directly from environment samples. For technical tractability, these works focus on regularized bi-level RL, where regularization induces the PL condition for the policy optimization objective. Their analyses largely mirror those of generic bi-level optimization under a lower-level PL condition and do not solve the original unregularized problem. In addition, due to the nest-loop structure, these algorithms require a target accuracy to be specified in advance to determine the number of inner-loop iterations, making the convergence guarantees not *anytime valid*: there is no guarantee that the optimality gap decreases monotonically after every iteration, and running the algorithm beyond the prescribed number of iterations does not further reduce the gap. Our paper exactly addresses these limitations.

2 Formulation

Consider an infinite-horizon discounted-reward MDP defined as $\mathcal{M}_x = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r_x, \gamma)$, where $x \in \mathbb{R}^d$ is an exogenous control parameter. Under a fixed x, \mathcal{M}_x is a standard MDP. The state space \mathcal{S} and action space \mathcal{A} are assumed to be finite. The transition kernel is denoted by $\mathcal{P}: \mathcal{S} \times \mathcal{A} \to \Delta(\mathcal{S})$, and we use $\mathcal{P}(s' \mid s, a)$ to represent the probability that the next state is s' when an agent takes action a in state s. The reward function $r_x: \mathcal{S} \times \mathcal{A} \to [0,1]$ is a function of x. The discount factor is denoted by $\gamma \in (0,1)$. Our paper considers the setting where the transition kernel is independent of x, motivated by applications such as RLHF and reward shaping, where the exogenous variable modulates the reward function but not the system dynamics.

An agent learning in this MDP may not directly observe x, and takes actions according to a policy $\pi:\mathbb{S}\to\Delta_{\mathcal{A}}$, which we can represent as a table $\Delta_{\mathcal{A}}^{\mathcal{S}}\in\mathbb{R}^{|\mathcal{S}|\times|\mathcal{A}|}$. Given a control-policy pair (x,π) , we measure its performance in state s by the value function

$$\textstyle V^{x,\pi}(s) \triangleq \mathbb{E}_{a_k \sim \pi(\cdot \mid s_k), s_{k+1} \sim \mathcal{P}(\cdot \mid s_k, a_k)} \left[\sum_{k=0}^{\infty} \gamma^k r_x(s_k, a_k) \mid s_0 = s \right] = \mathbb{E}_{s' \sim d_s^\pi, \; a' \sim \pi(\cdot \mid s')} [r_x(s', a')],$$

where $d_s^{\pi} \in \Delta_{\mathcal{S}}$ is the discounted visitation distribution under initial state s

$$d_s^{\pi}(s') \triangleq (1 - \gamma) \mathbb{E}_{a_k \sim \pi(\cdot \mid s_k), s_{k+1} \sim \mathcal{P}(\cdot \mid s_k, a_k)} \left[\sum_{k=0}^{\infty} \gamma^k \mathbf{1}(s_k = s') \mid s_0 = s \right].$$

Under an initial state distribution $\rho \in \Delta_{\mathcal{S}}$, we define the expected cumulative reward under (x, π)

$$J(x,\pi) \triangleq \mathbb{E}_{s \sim \rho}[V^{x,\pi}(s)] = \mathbb{E}_{s \sim d_{\rho}^{\pi}, a \sim \pi(\cdot | s)}[r_{x}(s,a)], \quad \text{where } d_{\rho}^{\pi} \triangleq \mathbb{E}_{s \sim \rho}[d_{s}^{\pi}].$$

If x were fixed, our goal would be to find a policy that maximizes $J(x,\pi)$. In the mean time, the exogenous controller has its own objective to optimize, anticipating the best response from the policy optimization agent. We denote the controller's objective by $f: \mathbb{R} \times \Delta_{\mathcal{A}}^{\mathcal{S}} \to \mathbb{R}$. Let $\Pi^{\star}(x)$ be the set of optimal policies under control x, which we note may not be singleton, and g be a function that maps $\Pi^{\star}(x)$ to a unique optimal policy within the set (we will shortly introduce g). The controller's optimization problem can then be formulated as the following bi-level program

$$\begin{aligned} \min_{x \in \mathbb{R}^d} \quad & f(x, g(\Pi^\star(x))), & \textbf{Upper-Level} \\ \text{s.t.} \quad & \Pi^\star(x) \triangleq \mathop{\mathrm{argmax}}_\pi J(x, \pi). & \textbf{Lower-Level RL Objective} \end{aligned} \tag{1}$$

Our goal in this paper is to solve (1). This is a challenging problem, as the lower-level objective lacks strong structural properties and may not admit a unique solution. To introduce additional structure and enhance the solvability, we add entropy regularization into the lower-level objective, which leads to solution uniqueness and a strong form of "gradient domination". We stress that the entropy-regularized formulation serves only as an intermediate tool – our ultimate aim remains to solve the original, unregularized problem in (1).

2.1 Entropy Regularization

We discuss the regularized objective and its structural properties. Given (x, π) and regularization weight τ , we define the regularized value function $V_{\tau}^{x,\pi} \in \mathbb{R}^{|\mathcal{S}|}$ and expected cumulative reward J_{τ}

$$V_{\tau}^{x,\pi}(s) \triangleq \mathbb{E}_{a_k \sim \pi(\cdot|s_k), s_{k+1} \sim \mathcal{P}(\cdot|s_k, a_k)} \left[\sum_{k=0}^{\infty} \gamma^k \left(r_x(s_k, a_k) - \tau \log \pi(a_k \mid s_k) \right) \mid s_0 = s \right]$$

$$= \frac{1}{1 - \gamma} \mathbb{E}_{s' \sim d_s^{\pi}, a' \sim \pi(\cdot|s')} [r_x(s', a') + \tau E(\pi, s')], \tag{2}$$

$$J_{\tau}(x,\pi) \triangleq \mathbb{E}_{a_k \sim \pi(\cdot|s_k), s_{k+1} \sim \mathcal{P}(\cdot|s_k, a_k)} \left[\sum_{k=0}^{\infty} \gamma^k \left(r_x(s_k, a_k) - \tau \log \pi(a_k \mid s_k) \right) \mid s_0 \sim \rho \right]$$

$$= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi}, a \sim \pi(\cdot \mid s)} [r_x(s, a) + \tau E(\pi, s)] = \mathbb{E}_{s \sim \rho} [V_{\tau}^{x, \pi}(s)], \tag{3}$$

where $E(\pi,s) = -\sum_a \pi(a\mid s) \log \pi(a\mid s)$ is the entropy function. Under regularization weight $\tau \leq 1$, we have $|V^{x,\pi}_{\tau}(s)| \leq B_V$ for all x,π,s , where $B_V = \frac{1+\log|\mathcal{A}|}{1-\gamma}$.

If the initial state distribution has a full support, an assumption we will shortly introduce and impose throughout the paper, then the optimizer of $J_{\tau}(x,\cdot)$ is unique for any $\tau>0$. We define the operator $\pi_{\tau}^{\star}:\mathbb{R}^{d}\to\Delta_{\mathcal{A}}^{\mathcal{S}}$, which maps a control variable to the optimal policy induced by it

$$\pi_{\tau}^{\star}(x) \triangleq \operatorname{argmax}_{\pi} J_{\tau}(x, \pi), \quad \forall x \in \mathbb{R}^{d}.$$
 (4)

As we use the regularized RL problem to approximate the original one, it is important to understand how $\pi_{\tau}^{\star}(x)$ relates to $\Pi^{\star}(x)$. We make the connection in Lemma 1, under the following assumption on initial state distribution and ergodic Markov chain. The assumption is commonly made in the RL literature to guarantee that the Markov chain of states under any policy has a unique, well-defined stationary distribution [Mei et al., 2020, Wu et al., 2020, Khodadadian et al., 2022].

Assumption 1 (Sufficient Exploration) The initial state distribution ρ is bounded away from zero, i.e. there exists a constant $\rho_{\min} > 0$ such that $\rho(s) \geq \rho_{\min}$ for all $s \in \mathcal{S}$. Additionally, for any π , the Markov chain $\{s_t\}$ generated by P^{π} following $s_{t+1} \sim P^{\pi}(\cdot \mid s_t)$ is ergodic.

Lemma 1 We define $\pi^*(x)$ to be the optimal policy for the original, unregularized MDP with the largest (visitation-weighted) entropy

$$\pi^{\star}(x) \triangleq \operatorname{argmax}_{\pi \in \Pi^{\star}(x)} \mathbb{E}_{s \sim d_{o}^{\pi}}[E(\pi, s)]. \tag{5}$$

Then, under Assumption 1, it holds that $\pi^*(x)$ is unique for all x and is the limit point of $\{\pi_{\tau}^*(x)\}_{\tau}$ $\pi^*(x) = \lim_{\tau \to 0} \pi_{\tau}^*(x)$.

The uniqueness of $\pi^*(x)$ is not obvious and does not directly follow from any known results in convex optimization, since Π^* is not a convex set and the weighted entropy objective in (5) is not concave (see Lemma 3.1 of Hazan et al. [2019] for a proof of non-concavity). Our proof of Lemma 1, presented in detail in Appendix 1, exploits the strict concavity of the unweighted entropy function $E(\cdot,s)$ in the interior of the the simplex, as well as the structure that $\Pi^*(x)$, though non-convex, is a connected set with special structure [Zeng et al., 2023].

Our objective in this work is to solve the optimization problem below, which corresponds to (1) with q mapping a set to the (unique) element of the set maximizing the weighted entropy

$$\min_{x} \Phi(x) \triangleq f(x, \pi^{\star}(x)).$$
 (6)

We also define the regularized version of the objective, serving as a key intermediary in our analysis. Conceptually, the algorithm to be introduced optimizes the regularized objective Φ_{τ} as the regularization attenuation drives Φ_{τ} toward Φ .

$$\min_{x} \Phi_{\tau}(x) \triangleq f(x, \pi_{\tau}^{\star}(x)). \tag{7}$$

3 Algorithm Development

In this section we develop a single-loop first-order algorithm that optimizes Φ_{τ} while gradually decaying τ to zero, thereby recovering the solution to (6). The algorithm operates under stochastic gradient samples of the upper-level objective, as well as state transition and reward samples from the lower-level MDP. We design the algorithm based on a penalty reformulation, a technique recently developed for solving bi-level problems with lower-level strong convexity and PL condition [Shen and Chen, 2023, Kwon et al., 2023]. We begin by presenting an overview of the reformulation in our context.

3.1 Preliminaries - Penalty Reformulation

Our goal is to solve (7) via (stochastic) gradient descent. By the implicit function theorem [Ghadimi and Wang, 2018], $\nabla_x \Phi_{\tau}(x)$ admits a closed-form expression when $\nabla^2_{\pi,\pi} J_{\tau}(x,\pi^*(x))$ is invertible

$$\nabla_x \Phi_{\tau}(x) = \nabla_x f(x, \pi_{\tau}^{\star}(x)) + \nabla_{\pi} f(x, \pi_{\tau}^{\star}(x)) \frac{\partial \pi_{\tau}^{\star}(x)}{\partial x}$$
(8)

$$= \nabla_x f(x, \pi_{\tau}^{\star}(x)) - \nabla_{x,\pi}^2 J_{\tau}(x, \pi_{\tau}^{\star}(x)) \nabla_{\pi,\pi}^2 J_{\tau}(x, \pi_{\tau}^{\star}(x))^{-1} \nabla_{\pi} f(x, \pi_{\tau}^{\star}(x)). \tag{9}$$

Obtaining unbiased samples of $\nabla_x \Phi_{\tau}(x)$ based on (9), however, poses significant challenges, as the expression depends on second-order Jacobian and Hessian terms that cannot be efficiently estimated from state–reward samples. The penalty reformulation is designed to provide an alternative approach of obtaining (asymptotically) unbiased gradient estimates, only requiring first-order information.

Recall the definition of $\pi_{\star}^{\star}(x)$ in (4). We can re-write (7) as follows by introducing a constraint

$$\min_{x,\pi} f(x,\pi)$$
 s.t. $J_{\tau}(x,\pi_{\tau}^{\star}(x)) - J_{\tau}(x,\pi) \le 0.$ (10)

Given a positive constant w, we define

$$\mathcal{L}_{w,\tau}(x,\pi) \triangleq f(x,\pi) + \frac{1}{w} \Big(J_{\tau}(x,\pi_{\tau}^{\star}(x)) - J_{\tau}(x,\pi) \Big), \tag{11}$$

$$\Phi_{w,\tau}(x) \triangleq \min_{\pi} \mathcal{L}_{w,\tau}(x,\pi) = \min_{\pi} f(x,\pi) + \frac{1}{w} \Big(J_{\tau}(x,\pi_{\tau}^{\star}(x)) - J_{\tau}(x,\pi) \Big). \tag{12}$$

We can regard $\mathcal{L}_{w,\tau}$ as the Lagrangian associated with (10), in which 1/w plays the role of the dual variable. To solve (10), it may be tempting to find a minimax saddle point of the Lagrangian using gradient descent ascent. However, as pointed out in Kwon et al. [2023], the solution of (10) is only attained in the limit as the dual variable becomes infinitely large (i.e. w=0). This motivates us to treat w as a parameter governed by a prescribed decay schedule towards zero, rather than as a dual variable updated via gradient ascent. It is known from Kwon et al. [2023][Lemma 3.1] that $\nabla_x \Phi_{w,\tau}(x)$ admits the following expression involving only first-order terms¹

$$\nabla_x \Phi_{w,\tau}(x) = \nabla_x \mathcal{L}_{w,\tau}(x, \pi_{w,\tau}^{\star}(x))$$

¹We follow the convention and use $\nabla_x \mathcal{L}_{w,\tau}(x, \pi_{w,\tau}^\star(x))$ to denote the partial gradient with respect to x evaluated at $(x, \pi_{w,\tau}^\star(x))$, i.e. $\nabla_x \mathcal{L}_{w,\tau}(x, \pi_{w,\tau}^\star(x)) = \nabla_x \mathcal{L}_{w,\tau}(x, \pi) \mid_{\pi=\pi_{w,\tau}^\star(x)}$. The same principle will be used for other functions, such as f and J_τ .

$$= \nabla_x f(x, \pi_\tau^{\star}(x)) + \frac{1}{w} \Big(\nabla_x J_\tau(x, \pi_\tau^{\star}(x)) - \nabla_x J_\tau(x, \pi_{w,\tau}^{\star}(x)) \Big), \tag{13}$$

where we define $\pi_{w,\tau}^{\star}(x) \triangleq \operatorname{argmin}_{\pi} \mathcal{L}_{w,\tau}(x,\pi)$ for all $w,\tau > 0$. Importantly, $\nabla_x \Phi_{w,\tau}(x)$ closely tracks $\nabla_x \Phi_{\tau}(x)$ – the distance between $\nabla_x \Phi_{w,\tau}(x)$ and $\nabla_x \Phi_{\tau}(x)$ scales linearly in w, a result which we establish later in Lemma 11.

3.2 Single-Loop Algorithm with Decaying Penalty and Regularization

We introduce x_k as an estimate of the solution to the bi-level objective (k) is the iteration index) and design an algorithm that iteratively carries out stochastic gradient descent on x_k in an approximate direction of $\nabla_x \Phi_{w_k,\tau_k}(x_k)$, estimated using online samples from the MDP. Here w_k and τ_k are time-varying penalty and regularization weights. As w_k,τ_k decay according to carefully designed schedules, the surrogate objective $\Phi_{w_k,\tau_k}(x_k)$ increasingly approximates $\Phi(x_k)$, allowing us to solve the original bi-level problem.

To estimate $\nabla_x \Phi_{w_k, \tau_k}$ based on (13), we need estimates of $\nabla_x J_{\tau_k}(x_k, \pi_{\tau_k}^\star(x_k))$ and $\nabla_x J_{\tau_k}(x_k, \pi_{w_k, \tau_k}^\star(x_k))$. Note that $\nabla_x J_{\tau}(x, \pi)$ can be expressed in the simple form below

$$\nabla_x J_{\tau}(x,\pi) = \mathbb{E}_{s \sim d_{\alpha}^{\pi}, a \sim \pi(\cdot|s)} [\nabla_x r_x(s,a)]. \tag{14}$$

Given (14), if we had access to an oracle that generates $\pi=\pi^\star_{\tau_k}(x)$ for any x, we would obtain asymptotically unbiased samples of $\nabla_x J_{\tau_k}(x_k,\pi^\star_{\tau_k}(x_k))$ by simply generating a Markovian chain $\{s_k,a_k\}$ under $\pi=\pi^\star_{\tau_k}(x_k)$ and evaluating $\nabla_x r_{x_k}(s_k,a_k)$ along the trajectory. The same can be done to estimate $\nabla_x J_{\tau_k}(x_k,\pi^\star_{w_k,\tau_k}(x_k))$ if $\pi^\star_{w_k,\tau_k}(x_k)$ were available. However, $\pi=\pi^\star_{\tau_k}(x_k)$ and $\pi^\star_{w_k,\tau_k}(x_k)$ are solutions to (augmented) lower-level RL problems and cannot be directly accessed. To overcome the oracle unavailability, we introduce the iterates π_k,π^ℓ_k as approximations of $\pi^\star_{\tau_k}(x_k),\pi^\star_{w_k,\tau_k}(x_k)$, and update them via another layer of stochastic gradient ascent.

Existing bi-level RL methods [Yang et al., 2025, Gaur et al., 2025] typically introduce a nested-loop algorithmic structure when estimating these optimal policies, ensuring that $\pi_k, \pi_k^{\mathcal{L}}$ from the inner loop fully converges to $\pi_{\tau_k}^{\star}(x_k), \pi_{w_k,\tau_k}^{\star}(x_k)$ up to a desired precision. However, such nested-loop algorithms are usually inconvenient to implement in practice and require setting the precision in advance. Our algorithm instead updates $\pi_k, \pi_k^{\mathcal{L}}$ in the same loop as x_k , with a larger step size (i.e. on a faster time scale). Specifically, we maintain policy parameters $\theta_k, \theta_k^{\mathcal{L}}$ that encode $\pi_k, \pi_k^{\mathcal{L}}$ (with notation $\pi_k = \pi_{\theta_k}, \pi_k^{\mathcal{L}} = \pi_{\theta_k^{\mathcal{L}}}$) and iteratively refine them according to

$$\theta_{k+1} = \theta_k + \alpha_k \widetilde{\nabla}_{\theta} J_{\tau_k}(x_k, \pi_{\theta_k}), \quad \theta_{k+1}^{\mathcal{L}} = \theta_k^{\mathcal{L}} + \alpha_k \Big(- \widetilde{\nabla}_{\theta} f(x_k, \pi_{\theta_k^{\mathcal{L}}}) + \frac{1}{w_k} \widetilde{\nabla}_{\theta} J_{\tau_k}(x_k, \pi_{\theta_k^{\mathcal{L}}}) \Big).$$

Here α_k is a step size properly balanced with the decay rates of w_k and τ_k , and $\widetilde{\nabla}_{\theta} f, \widetilde{\nabla}_{\theta} J_{\tau_k}$ denote stochastic samples of the true gradients. Note that $\nabla_{\theta} J_{\tau_k}$ admits the following closed-form expression, and can be estimated in an asymptotically unbiased manner via an actor-critic approach

$$\nabla_{\theta} J_{\tau}(x, \pi_{\theta}) = \mathbb{E}_{d_{\rho}^{\pi_{\theta}}, \pi_{\theta}} \left[\left(r_{x}(s, a) - \tau \log \pi_{\theta}(a \mid s) + \gamma V_{\tau}^{x, \pi_{\theta}}(s') \right) \frac{\nabla_{\theta} \log \pi_{\theta}(a \mid s)}{1 - \gamma} \right]. \tag{15}$$

Actor-critic methods samples stochastic gradients according to (15), replacing the unobservable value function with an estimate which is updated on an even faster time scale via temporal difference learning. Specifically, we introduce two variables \hat{V}_k , $\hat{V}_k^{\mathcal{L}}$ to track $V_{\tau_k}^{x_k,\pi_k}$, $V_{\tau_k}^{x_k,\pi_k^{\mathcal{L}}}$ and present their update rules in (20), where $\Pi_{B_V}: \mathbb{R}^{|\mathcal{S}|} \to \mathbb{R}^{|\mathcal{S}|}$ denotes the element-wise projection of a vector to the interval $[0,B_V]$. The projection operator guarantees the stability of the value function estimates, and the interval contains the true (regularization) value function under regularization weight $\tau_k \leq 1$.

The algorithm can be described at an abstract level as follows. We perform stochastic gradient descent on x_k along the hyper-gradient direction. The hyper-gradient estimation relies on the solutions of lower-level RL problem and the penalty-augmented RL objective, which we obtain via a single-loop actor-critic method. Our actor-critic procedure follows the standard framework, with the key distinction that we incorporate entropy regularization and gradually attenuate its weight over time. Note that despite the resemblance of our actor-critic updates to existing algorithms, the analysis is significantly more challenging in the bi-level setting. In particular, the learning targets for the

lower-level policies are non-stationary, evolving both with the penalty and regularization schedules and with the updates of the upper-level variable. We overcome this challenge with a novel error decomposition scheme that tightly links the sub-optimality gap of the lower-level RL problem to that of the bi-level objective, under the shifting landscape which becomes less structured over time as the penalty and regularization weights decay.

We formally state the updates in Algorithm 1, in which we represent the policies through tabular softmax parameterization², i.e. the parameter $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ encodes the policy π_{θ} according to

$$\pi_{\theta}(a \mid s) = \frac{\exp(\theta(s, a))}{\sum_{a'} \exp(\theta(s, a'))}.$$

Algorithm 1 employs three step size parameters and two penalty/regularization weights, which are all time-decaying sequences: step size ζ_k for upper-level variable update, step size α_k for policy update, step size β_k for value function update, penalty weight w_k , and regularization τ_k . The step sizes are associated with the primal variable (x_k) update, and we need to choose $\zeta_k \ll \alpha_k \ll \beta_k$ to approximate the nested-loop dynamics, where we run a large number of value function updates per policy update and a large number of policy updates per upper-level variable update.

Algorithm 1 Single-Loop Actor-Critic Algorithm for Bi-Level RL

- 1: **Initialize:** control variable x_0 , policy parameters θ_0 and $\theta_0^{\mathcal{L}}$, value function estimates $\hat{V}_0, \hat{V}_0^{\mathcal{L}} \in \mathbb{R}^{|\mathcal{S}|}$
- 2: **for** iteration k = 0, 1, 2, ... **do**
- 3: Trajectory 1:

With probability $1-\gamma$, restart the trajectory by taking $s_{k+1} \sim \rho$. With probability γ , continue following the current trajectory. Take action $a_k \sim \pi_{\theta_k}(\cdot \mid s_k)$, receive rewards $r_{x_k}(s_k, a_k)$, and observe the next state $s_{k+1} \sim \mathcal{P}(\cdot \mid s_k, a_k)$.

- 4: Trajectory 2: With probability $1-\gamma$, restart the trajectory by taking $\bar{s}_{k+1}\sim \rho$. With probability γ , continue following the current trajectory. Take action $\bar{a}_k\sim \pi_{\theta_k^{\mathcal{L}}}(\cdot\mid \bar{s}_k)$, receive rewards $r_{x_k}(\bar{s}_k,\bar{a}_k)$, and observe the next state $\bar{s}_{k+1}\sim \mathcal{P}(\cdot\mid \bar{s}_k,\bar{a}_k)$.
- 5: Observe/Obtain $\xi_k \sim \mu$
- 6: Control variable update:

$$x_{k+1} = x_k - \zeta_k \left(\widetilde{\nabla}_x f(x_k, \pi_{\theta_k^{\mathcal{L}}}, \xi_k) + \frac{1}{w_k} \left(\nabla_x r_{x_k}(s_k, a_k) - \nabla_x r_{x_k}(\bar{s}_k, \bar{a}_k) \right) \right). \tag{16}$$

7: Policy update:

$$\theta_{k+1} = \theta_k + \alpha_k \Big(r_{x_k}(s_k, a_k) + \tau_k E(\pi_{\theta_k}, s_k) + \gamma \hat{V}_k(s_{k+1}) \Big) \nabla_{\theta} \log \pi_{\theta_k}(a_k \mid s_k), \quad (17)$$

$$\theta_{k+1}^{\mathcal{L}} = \theta_k^{\mathcal{L}} + \alpha_k \Big(\Big(r_{x_k}(\bar{s}_k, \bar{a}_k) + \tau_k E(\pi_{\theta_k^{\mathcal{L}}}, \bar{s}_k) + \gamma \hat{V}_k^{\mathcal{L}}(\bar{s}_{k+1}) \Big) \nabla_{\theta} \log \pi_{\theta_k^{\mathcal{L}}}(\bar{a}_k \mid \bar{s}_k) - w_k \widetilde{\nabla}_{\theta} f(x_k, \pi_{\theta_k^{\mathcal{L}}}, \xi_k) \Big), \quad (18)$$

$$\pi_k = \operatorname{softmax}(\theta_k), \quad \pi_k^{\mathcal{L}} = \operatorname{softmax}(\theta_k^{\mathcal{L}}). \quad (19)$$

8: Value function update:

$$\hat{V}_{k+1} = \Pi_{B_V} \Big(\hat{V}_k + \beta_k e_{s_k} \big(r_{x_k}(s_k, a_k) + \tau_k E(\pi_{\theta_k}, s_k) + \gamma \hat{V}_k(s_{k+1}) - \hat{V}_k(s_k) \big) \Big),
\hat{V}_{k+1}^{\mathcal{L}} = \Pi_{B_V} \Big(\hat{V}_k^{\mathcal{L}} + \beta_k e_{\bar{s}_k} \big(r_{x_k}(\bar{s}_k, \bar{a}_k) + \tau_k E(\pi_{\theta_k}^{\mathcal{L}}, \bar{s}_k) + \gamma \hat{V}_k^{\mathcal{L}}(\bar{s}_{k+1}) - \hat{V}_k^{\mathcal{L}}(\bar{s}_k) \big) \Big).$$
(20)

9: end for

²We consider the softmax parameterization for the purpose of mathematical analysis. The algorithm is compatible with any function approximation in practical implementations.

4 Convergence Analysis

In this section, we present the finite-time and finite-sample analysis of the proposed algorithm. We start by introducing the technical assumptions.

Assumption 2 (Hessian Invertibility) The Hessian $\nabla^2_{\theta,\theta} J_{\tau}(x,\pi_{\theta})$ is invertible for any x,θ and $\tau \geq 0$.

Recall from (9) that the Hessian inverse appears in the hyper-gradient. Though we do not explicit work with the Hessian, the assumption importantly guarantees the differentiability of Φ_{τ} . Let $\sigma_{\min}(\cdot)$ denote the smallest *singular value* of a matrix. Assumption 2 implies that there exists a constant $\sigma > 0$ such that

$$\sigma_{\min}\Big(\nabla_{\theta,\theta}^2 J_{\tau}(x,\pi_{\theta})\Big) \ge \underline{\sigma}, \quad \forall x, \theta, \tau,$$
 (21)

Without losing generality, we let $\underline{\sigma} \leq 1$ for the convenience of combining terms in the analysis. Note that (21) should not be confused with the strong convexity of $J_{\tau}(x, \pi_{\theta})$ with respect to θ , which requires $\nabla^2_{\theta,\theta} J_{\tau}(x, \pi_{\theta})$ to be positive definite, i.e. its smallest *eigenvalue* is positive.

Assumption 3 (Lipschitz and Smooth Upper-Level Objective) The function f is differentiable, and we have access to unbiased stochastic gradient operators $\widetilde{\nabla}_x f(x,\pi,\xi)$, $\widetilde{\nabla}_\pi f(x,\pi,\xi)$ and i.i.d. samples ξ from a distribution μ such that

$$\mathbb{E}_{\xi \sim \mu}[\widetilde{\nabla}_x f(x, \pi, \xi)] = \nabla_x f(x, \pi), \quad \mathbb{E}_{\xi \sim \mu}[\widetilde{\nabla}_\pi f(x, \pi, \xi)] = \nabla_\pi f(x, \pi).$$

In addition, there is a bounded constant L_f such that $\forall x, x', \pi, \pi', \xi$

$$\|\widetilde{\nabla}_{x}f(x,\pi,\xi)\| \leq L_{f}, \quad \|\widetilde{\nabla}_{\pi}f(x,\pi,\xi)\| \leq L_{f},$$
$$\|\widetilde{\nabla}_{x}f(x,\pi,\xi) - \widetilde{\nabla}_{x}f(x',\pi',\xi)\| \leq L_{f}(\|x-x'\| + \|\pi-\pi'\|),$$
$$\|\widetilde{\nabla}_{\pi}f(x,\pi,\xi) - \widetilde{\nabla}_{\pi}f(x',\pi',\xi)\| \leq L_{f}(\|x-x'\| + \|\pi-\pi'\|).$$

We also assume that the minimizer of $f(\cdot,\pi)$ exists for any π , i.e. $f(x,\pi)$ never blows up to negative infinity. Without loss of generality, we can shift the function such that $f(x,\pi) \geq 0, \ \forall x,\pi$.

Assumption 4 (Lipschitz Reward, Gradient, and Hessian) There is a constant $L_r < \infty$ such that $\forall s, a, x_1, x_2$

$$|r_{x_1}(s, a) - r_{x_2}(s, a)| \le L_r ||x_1 - x_2||,$$

$$||\nabla_x r_{x_1}(s, a) - \nabla_x r_{x_2}(s, a)|| \le L_r ||x_1 - x_2||,$$

$$||\nabla_{x_r}^2 r_{x_1}(s, a) - \nabla_{x_r}^2 r_{x_2}(s, a)|| \le L_r ||x_1 - x_2||.$$

Assumptions 3 and 4 are standard regularity assumptions in the bi-level RL literature [Chakraborty et al., 2024, Gaur et al., 2025, Yang et al., 2025]. Comparable conditions on upper- and lower-level objectives are also commonly imposed by works on generic bi-level optimization [Kwon et al., 2023, Shen and Chen, 2023].

Assumption 5 (Regularization-Dependent PL Condition) Recall the definition of $\pi_{w,\tau}^{\star}$ after (13). The minimizer $\pi_{w,\tau}^{\star}(x)$ is unique for all x and $w,\tau>0$. In addition, there exists a constant $C_L>0$ such that for all $w,\tau>0$, we have

$$\|\nabla_{\theta} \mathcal{L}_{w,\tau}(x, \pi_{\theta})\|^2 \ge \frac{C_L \tau}{w} \Big(\mathcal{L}_{w,\tau}(x, \pi_{\theta}) - \mathcal{L}_{w,\tau}(x, \pi_{w,\tau}^{\star}(x)) \Big), \quad \forall x, \theta.$$
 (22)

This structural condition plays an important role in our analysis and states that the Lagrangian defined in (12) satisfies the PL condition with respect to the policy parameter θ , with a PL constant that attenuates as the regularization weight becomes smaller. While we directly impose (22) for convenience, the condition can be derived for a proper range of w under assumptions of full coverage initial state distribution (Assumption 1) and exploratory policy (i.e. $\pi_{\theta}(a \mid s)$ is uniformly lower bounded). To see this, note that as $w \to 0$, $L_{w,\tau}(x,\pi)$ approaches a 1/w-scaled (and shifted) version of $J_{\tau}(x,\pi)$, which is known to satisfy the PL condition with $C_L = \mathcal{O}(\rho_{\min}^2(\min_{s,a} \pi_{\theta}(a \mid s))^2)$ (see Mei et al. [2020][Lemma 15]). The scaling explains the dependence of the right-hand side of (22) on 1/w. For sufficiently small w, the contribution of the f term in $L_{w,\tau}$ remains negligible, so the PL condition of J_{τ} continues to dominate and allows (22) to hold.

4.1 Main Results

While our ultimate goal is to solve the objective (6), our analysis relies on jointly bounding the convergence of all variables through a coupled Lyapunov function that combines all residuals shown below. As Φ is non-convex, we may in general only find a first-order stationary point, and we measure the convergence of x_k by $\|\nabla_x \Phi(x_k)\|^2$, the squared gradient norm. To measure the convergence of θ_k and $\theta_k^{\mathcal{L}}$, we consider their distance to the optimal policy in the function space. Finally, the value function estimates \hat{V}_k , $\hat{V}_k^{\mathcal{L}}$ are measured by their ℓ_2 distance to the value functions under the latest upper-level decision variable and policy.

$$\varepsilon_k^{\theta,\mathcal{L}} = w_k \Big(\mathcal{L}_{w_k,\tau_k}(x_k, \pi_{\theta_k^{\mathcal{L}}}) - \mathcal{L}_{w_k,\tau_k}(x_k, \pi_{w_k,\tau_k}^{\star}(x_k)) \Big), \quad \varepsilon_k^{\theta} = J_{\tau_k}(x_k, \pi_{\tau_k}^{\star}(x_k)) - J_{\tau_k}(x_k, \pi_{\theta_k}),$$

$$\varepsilon_k^{V} = \|\hat{V}_k - V_{\tau_k}^{x_k,\pi_{\theta_k}}\|^2, \quad \varepsilon_k^{V,\mathcal{L}} = \|\hat{V}_k^{\mathcal{L}} - V_{\tau_k}^{x_k,\pi_{\theta_k^{\mathcal{L}}}}\|^2.$$

Theorem 1 Consider the iterates of Algorithm 2 under the step sizes and weights

$$\zeta_k = \frac{\zeta_0}{(k+1)^{c_\zeta}}, \quad \alpha_k = \frac{\alpha_0}{(k+1)^{c_\alpha}}, \quad \beta_k = \frac{\beta_0}{(k+1)^{c_\beta}}, \quad w_k = \frac{w_0}{(k+1)^{c_w}}, \quad \tau_k = \frac{\tau_0}{(k+1)^{c_\tau}},$$

with $c_{\zeta}=\frac{9}{10}, c_{\alpha}=\frac{1}{2}, c_{\beta}=\frac{1}{2}, c_{w}=\frac{3}{20}, c_{\tau}=\frac{1}{20}$ and properly selected $\zeta_{0}, \alpha_{0}, \beta_{0}, w_{0}, \tau_{0}$. Under Assumptions I-5, we have for all $k\geq 0$,

$$\min_{t < k} \mathbb{E}[\|\nabla_x \Phi(x_t)\|^2] \le \mathcal{O}\left(\frac{\Phi(x_0) + \varepsilon_0^{\theta} + \varepsilon_0^{\theta, \mathcal{L}} + \varepsilon_0^{V} + \varepsilon_0^{V, \mathcal{L}}}{(k+1)^{1/10}}\right) + \widetilde{\mathcal{O}}\left(\frac{1}{(k+1)^{1/10}}\right).$$

Theorem 1 shows that the best iterate of Algorithm 1 converges to a stationary point of the bi-level RL objective with rate $\widetilde{\mathcal{O}}(k^{-1/10})$. As the algorithm draws two samples in each iteration, this time complexity translates to a sample complexity of the same order. To our knowledge, this is the first time an algorithm has been shown to provably solve the original, unregularized bi-level RL problem. The key technical insight and novelty enabling our analysis is 1) that we recognize the RL objective as one observing a regularization-dependent PL condition with the PL constant diminishing as regularization approaches zero, 2) a multi-time-scale stochastic approximation analysis that balances the decay of step sizes and w_k with that of τ_k , allowing the algorithm convergence to be established under the challenge of a time-varying optimization landscape.

We may also choose to instantiate Algorithm 1 with a constant regularization weight, targeting to optimize the regularized bi-level objective rather than the original one. In this regime, as the lower-level problem satisfies the PL condition with a fixed PL constant, we establish a faster convergence rate in the following theorem.

Theorem 2 Given any fixed regularization weight τ_0 , i.e. $\tau_k = \tau_0$ for all $k \ge 0$, consider the iterates of Algorithm 2 under the step sizes and penalty weight

$$\zeta_k = \frac{\zeta_0}{(k+1)^{c_\zeta}}, \quad \alpha_k = \frac{\alpha_0}{(k+1)^{c_\alpha}}, \quad \beta_k = \frac{\beta_0}{(k+1)^{c_\beta}}, \quad w_k = \frac{w_0}{(k+1)^{c_w}},$$

with $c_{\zeta}=\frac{2}{3},c_{\alpha}=\frac{1}{2},c_{\beta}=\frac{1}{2},c_{w}=\frac{1}{6}$ and properly selected $\zeta_{0},\alpha_{0},\beta_{0},w_{0}$. Under Assumptions 1-5, we have for all $k\geq0$,

$$\min_{t < k} \mathbb{E}[\|\nabla_x \Phi_{\tau_0}(x_t)\|^2] \le \mathcal{O}\left(\frac{\Phi_{\tau_0}(x_0) + \varepsilon_0^{\theta} + \varepsilon_0^{\theta, \mathcal{L}} + \varepsilon_0^{V} + \varepsilon_0^{V, \mathcal{L}}}{(k+1)^{1/3}}\right) + \widetilde{\mathcal{O}}\left(\frac{1}{(k+1)^{1/3}}\right).$$

Theorem 2 establishes a finite-time convergence rate of $\widetilde{\mathcal{O}}(k^{-1/3})$ to a regularized stationary point, and again implies a sample complexity of the same order. Importantly, this rate surpasses that of the F²SA algorithm in Kwon et al. [2023], which is derived under the stronger assumption of lower-level strong convexity. We have achieved the rate improvement under weaker lower-level structure by designing a novel error decomposition scheme that allows us to tightly bound the residuals in the policy iterates based on the PL condition. We also note that the complexity of Algorithm 1 matches the best-known complexity of a nested-loop algorithm developed in Gaur et al. [2025] for solving the regularized bi-level RL problem.

5 Numerical Simulations

We numerically verify the convergence of Algorithm 1 on a small-scale synthetic bi-level RL problem, where the lower-level MDP is defined a 10x10 grid. The reward for the MDP is the negated distance between the current state and a goal position (which encourages the lower-level RL agent to reach the goal in as few steps as possible), whereas the goal is placed by the upper level decision variable. We design the upper-level objective — as a function of upper-level decision variable x and the optimal policy $\pi^*(x)$ — to penalize deviations of the goal position from the center of the grid, while encouraging $\pi^*(x)$ to have direction biases towards moving down and right.

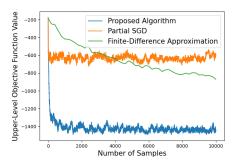


Figure 1: Upper-Level Decision Variable Convergence in Function Value Space

A natural baseline is an alternating (partial) gradient descent—ascent scheme: we maintain iterates (x_k,θ_k) and approximate the full hyper-gradient $\nabla_x\Phi(x_k)$ by its partial component $\nabla_xf(x_k,\pi_{\theta_k})$, while updating π_{θ_k} in the same loop to approximate $\pi^\star(x)$. We refer to this approach as "Partial SGD" and note that it may be stuck at sub-optimal solutions when the partial gradient is misaligned with the true hyper-gradient. This phenomenon is evident in Figure 1, where our proposed algorithm attains a better solution and achieves a smaller objective value than Partial SGD. We also compare against an algorithm that performs iterative gradient descent using the chain-rule expression (8), with $\frac{\partial \pi_\tau^\star(x)}{\partial x}$ estimated via finite differencing. Specifically, for a scalar x, we approximate

$$\frac{\partial \pi_{\tau}^{\star}(x)}{\partial x} \approx \frac{\pi_{\tau}^{\star}(x+\epsilon) - \pi_{\tau}^{\star}(\tau_k)}{\epsilon},$$

where $\pi_{\tau}^{\star}(x+\epsilon)$ and $\pi_{\tau}^{\star}(x)$ are computed through a large number of inner-loop RL iterations. For vector or tensor x, the approximation is carried out entry-wise and then aggregated. We refer to this method as 'Finite-Difference Approximation'" in Figure 1, and note that it is highly sample-inefficient, owing to both the computational overhead and the inaccuracy inherent in finite-difference-based gradient estimation. Further details on the upper-level objective, MDP reward function, step-size parameters, and initialization are provided in Appendix E.

Disclaimer

This paper was prepared for informational purposes ["in part" if the work is collaborative with external partners] by the Artificial Intelligence Research group of JPMorgan Chase & Co. and its affiliates ("JP Morgan") and is not a product of the Research Department of JP Morgan. JP Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

References

Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021.

Souradip Chakraborty, Amrit S Bedi, Alec Koppel, Huazheng Wang, Dinesh Manocha, Mengdi Wang, and Furong Huang. Parl: A unified framework for policy alignment in reinforcement learning from human feedback. In *ICLR*, 2024.

Lesi Chen, Jing Xu, and Jingzhao Zhang. On finding small hyper-gradients in bilevel optimization: Hardness results and improved analysis. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 947–980. PMLR, 2024.

- Mudit Gaur, Amrit Singh Bedi, Raghu Pasupathu, and Vaneet Aggarwal. On the sample complexity bounds in bilevel reinforcement learning. *arXiv preprint arXiv:2503.17644*, 2025.
- Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018.
- Elad Hazan, Sham Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum entropy exploration. In *International Conference on Machine Learning*, pages 2681–2691. PMLR, 2019.
- Yujing Hu, Weixun Wang, Hangtian Jia, Yixiang Wang, Yingfeng Chen, Jianye Hao, Feng Wu, and Changjie Fan. Learning to utilize shaping rewards: A new approach of reward shaping. *Advances in Neural Information Processing Systems*, 33:15931–15941, 2020.
- Sajad Khodadadian, Thinh T Doan, Justin Romberg, and Siva Theja Maguluri. Finite-sample analysis of two-time-scale natural actor–critic algorithm. *IEEE Transactions on Automatic Control*, 68(6): 3273–3284, 2022.
- Jeongyeol Kwon, Dohyun Kwon, Stephen Wright, and Robert D Nowak. A fully first-order method for stochastic bilevel optimization. In *International Conference on Machine Learning*, pages 18083–18113. PMLR, 2023.
- Haoxiang Ma, Shuo Han, Ahmed Hemida, Jie Fu, et al. Adaptive incentive design for markov decision processes with unknown rewards. *Transactions on Machine Learning Research*, 2025.
- Jincheng Mei, Chenjun Xiao, Csaba Szepesvari, and Dale Schuurmans. On the global convergence rates of softmax policy gradient methods. In *International conference on machine learning*, pages 6820–6829. PMLR, 2020.
- Han Shen and Tianyi Chen. On penalty-based bilevel gradient descent method. In *International Conference on Machine Learning*, pages 30992–31015. PMLR, 2023.
- Han Shen, Zhuoran Yang, and Tianyi Chen. Principled penalty-based methods for bilevel reinforcement learning and rlhf. In *International Conference on Machine Learning*, pages 44774–44799. PMLR, 2024.
- Zebang Shen, Alejandro Ribeiro, Hamed Hassani, Hui Qian, and Chao Mi. Hessian aided policy gradient. In *International conference on machine learning*, pages 5729–5738. PMLR, 2019.
- Vinzenz Thoma, Barna Pásztor, Andreas Krause, Giorgia Ramponi, and Yifan Hu. Contextual bilevel reinforcement learning for incentive alignment. Advances in Neural Information Processing Systems, 37:127369–127435, 2024.
- Yue Frank Wu, Weitong Zhang, Pan Xu, and Quanquan Gu. A finite-time analysis of two time-scale actor-critic methods. *Advances in Neural Information Processing Systems*, 33:17617–17628, 2020.
- Quan Xiao and Tianyi Chen. Unlocking global optimality in bilevel optimization: A pilot study. *arXiv preprint arXiv:2408.16087*, 2024.
- Quan Xiao, Hui Yuan, AFM Saif, Gaowen Liu, Ramana Kompella, Mengdi Wang, and Tianyi Chen. A first-order generative bilevel optimization framework for diffusion models. arXiv preprint arXiv:2502.08808, 2025.
- Yan Yang, Bin Gao, and Ya-xiang Yuan. Bilevel reinforcement learning via the development of hyper-gradient without lower-level convexity. In *The 28th International Conference on Artificial Intelligence and Statistics*, 2025.
- Yifan Yang, Peiyao Xiao, and Kaiyi Ji. Achieving $\mathcal{O}(\epsilon^{-1.5})$ complexity in hessian/jacobian-free stochastic bilevel optimization. *Advances in Neural Information Processing Systems*, 36:39491–39503, 2023.
- Ziyu Ye, Rishabh Agarwal, Tianqi Liu, Rishabh Joshi, Sarmishta Velury, Quoc V Le, Qijun Tan, and Yuan Liu. Reward-guided prompt evolving in reinforcement learning for llms. In *Forty-second International Conference on Machine Learning*, 2025.

- Sihan Zeng, Malik Aqeel Anwar, Thinh T Doan, Arijit Raychowdhury, and Justin Romberg. A decentralized policy gradient approach to multi-task reinforcement learning. In *Uncertainty in Artificial Intelligence*, pages 1002–1012. PMLR, 2021.
- Sihan Zeng, Thinh Doan, and Justin Romberg. Regularized gradient descent ascent for two-player zero-sum markov games. *Advances in Neural Information Processing Systems*, 35:34546–34558, 2022a.
- Sihan Zeng, Thinh Doan, and Justin Romberg. Connected superlevel set in (deep) reinforcement learning and its application to minimax theorems. *Advances in Neural Information Processing Systems*, 36:20146–20163, 2023.
- Sihan Zeng, Thinh T Doan, and Justin Romberg. A two-time-scale stochastic optimization framework with applications in control and reinforcement learning. *SIAM Journal on Optimization*, 34(1): 946–976, 2024.
- Siliang Zeng, Chenliang Li, Alfredo Garcia, and Mingyi Hong. Maximum-likelihood inverse reinforcement learning with finite-time guarantees. *Advances in Neural Information Processing Systems*, 35:10122–10135, 2022b.
- Banghua Zhu, Stephen Bates, Zhuoran Yang, Yixin Wang, Jiantao Jiao, and Michael I Jordan. The sample complexity of online contract design. In *Proceedings of the 24th ACM Conference on Economics and Computation*, pages 1188–1188, 2023.
- Shaofeng Zou, Tengyu Xu, and Yingbin Liang. Finite-sample analysis for sarsa with linear function approximation. *Advances in neural information processing systems*, 32, 2019.

Appendix

Contents

| 1 | Intr | oduction | 1 |
|---|------|--|----|
| | 1.1 | Related Work | 3 |
| 2 | For | mulation | 3 |
| | 2.1 | Entropy Regularization | 4 |
| 3 | Algo | orithm Development | 5 |
| | 3.1 | Preliminaries – Penalty Reformulation | 5 |
| | 3.2 | Single-Loop Algorithm with Decaying Penalty and Regularization | 6 |
| 4 | Con | vergence Analysis | 8 |
| | 4.1 | Main Results | 9 |
| 5 | Nun | nerical Simulations | 9 |
| A | Free | quently Used Notations, Equations, and Inequalities | 14 |
| В | Pro | of of Theorems | 17 |
| | B.1 | Proof of Theorem 3 | 20 |
| | B.2 | Proof of Theorem 4 | 23 |
| C | Pro | of of Propositions | 24 |
| | C.1 | Proof of Proposition 1 | 24 |
| | C.2 | Proof of Proposition 3 | 26 |
| | C.3 | Proof of Proposition 4 | 30 |
| | C.4 | Proof of Proposition 5 | 34 |
| | C.5 | Proof of Proposition 2 | 36 |
| D | Pro | of of Supporting Results | 37 |
| | D.1 | Proof of Lemma 1 | 37 |
| | D.2 | Proof of Lemma 2 | 39 |
| | D.3 | Proof of Lemma 3 | 39 |
| | D.4 | Proof of Lemma 4 | 42 |
| | D.5 | Proof of Lemma 5 | 42 |
| | D.6 | Proof of Lemma 6 | 43 |
| | D.7 | Proof of Lemma 7 | 46 |
| | D.8 | Proof of Lemma 8 | 48 |

| C | Simulation Details | 59 |
|---|------------------------|----|
| | D.15 Proof of Lemma 15 | 56 |
| | D.14 Proof of Lemma 14 | 55 |
| | D.13 Proof of Lemma 13 | 55 |
| | D.12 Proof of Lemma 12 | 53 |
| | D.11 Proof of Lemma 11 | 51 |
| | D.10 Proof of Lemma 10 | 49 |
| | D.9 Proof of Lemma 9 | 49 |

A Frequently Used Notations, Equations, and Inequalities

• Besides the value function defined in (2), we also define the regularized Q function and advantage function

$$Q_{\tau}^{x,\pi}(s,a) \triangleq r_x(s,a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s' \mid s,a) V_{\tau}^{x,\pi}(s),$$

$$A_{\tau}^{x,\pi}(s,a) \triangleq Q_{\tau}^{x,\pi}(s,a) - \tau \log \pi(a \mid s) - V_{\tau}^{x,\pi}(s).$$
(23)

- We define the filtration $\mathcal{F}_k \triangleq \{\xi_0, \dots, \xi_k, s_0, \dots, s_k, a_0, \dots, a_k, \bar{s}_0, \dots, \bar{s}_k, \bar{a}_0, \dots, \bar{a}_k\}.$
- The PL condition in Assumption 5 implies quadratic growth, i.e. for any x

$$\mathcal{L}_{w,\tau}(x,\pi_{\theta}) - \mathcal{L}_{w,\tau}(x,\pi_{w,\tau}^{\star}(x)) \ge \frac{C_L \tau}{4w} \|\pi_{\theta} - \pi_{w,\tau}^{\star}(x)\|^2.$$
 (24)

In combination with (22), the inequality implies

$$\|\nabla_{\theta} \mathcal{L}_{w,\tau}(x, \pi_{\theta})\| \ge \frac{C_L \tau}{2w} \|\pi_{\theta} - \pi_{w,\tau}^{\star}(x)\|. \tag{25}$$

• We introduce the following shorthand notations that abstract the update operators in Algorithm 1. For any $x, \pi, \pi^{\mathcal{L}}, \theta, V, s, a, s', \bar{s}, \bar{a}, \xi$, we define

$$D_{w}(x, \pi, \pi^{\mathcal{L}}, s, a, \bar{s}, \bar{a}, \xi) = \widetilde{\nabla}_{x} f(x, \pi^{\mathcal{L}}, \xi) + \frac{1}{w} \Big(\nabla_{x} r_{x}(s, a) - \nabla_{x} r_{x}(\bar{s}, \bar{a}) \Big),$$
(26)

$$F_{w,\tau}(x, \theta, V, s, a, s', \xi) = \Big(r_{x}(s, a) + \tau E(\pi_{\theta}, s) + \gamma V(s') - V(s) \Big) \nabla_{\theta} \log \pi_{\theta}(a \mid s) - w \widetilde{\nabla}_{\theta} f(x, \pi_{\theta}, \xi),$$
(27)

$$G_{\tau}(x, \theta, V, s, a, s') = e_{s} \Big(r_{x}(s, a) + \tau E(\pi_{\theta}, s) + \gamma V(s') - V(s) \Big),$$
(28)

where e_s is the indicator function, i.e. the entry s has a value of one and all other entries are

With (26)-(28), we can rewrite the updates of Algorithm 1

$$x_{k+1} = x_k - \zeta_k D_{w_k}(x_k, \pi_k, \pi_k^{\mathcal{L}}, s_k, a_k, \bar{s}_k, \bar{a}_k, \xi_k),$$

$$\theta_{k+1} = \theta_k + \alpha_k F_{0,\tau_k}(x_k, \theta_k, \hat{V}_k, s_k, a_k, s'_k, \xi_k),$$

$$\theta_{k+1}^{\mathcal{L}} = \theta_k^{\mathcal{L}} + \alpha_k F_{w_k,\tau_k}(x_k, \theta_k^{\mathcal{L}}, \hat{V}_k^{\mathcal{L}}, \bar{s}_k, \bar{a}_k, \bar{s}'_k, \xi_k),$$

$$\hat{V}_{k+1} = \hat{V}_k + \beta_k G_{\tau_k}(x_k, \theta_k, \hat{V}_k, s_k, a_k, s'_k),$$

$$\hat{V}_{k+1}^{\mathcal{L}} = \hat{V}_k^{\mathcal{L}} + \beta_k G_{\tau_k}(x_k, \theta_k^{\mathcal{L}}, \hat{V}_k^{\mathcal{L}}, \bar{s}_k, \bar{a}_k, \bar{s}'_k).$$

We also define operators $\bar{D}, \bar{F}, \bar{F}$ with a proper sense of expectation.

$$\bar{D}_{w}(x,\pi,\pi^{\mathcal{L}}) \triangleq \mathbb{E}_{s \sim d_{\rho}^{\pi},a \sim \pi(\cdot|s),\bar{s} \sim d_{\rho}^{\pi\mathcal{L}},\bar{a} \sim \pi^{\mathcal{L}}(\cdot|\bar{s})),\xi \sim \mu} [D_{w}(x,\pi,\pi^{\mathcal{L}},s,a,\bar{s},\bar{a},\xi)], \quad (29)$$

$$\bar{F}_{w,\tau}(x,\theta,V) \triangleq \mathbb{E}_{s \sim d_{\rho}^{\pi\theta},a \sim \pi\theta(\cdot|s),s' \sim \mathcal{P}(\cdot|s,a),\xi \sim \mu} [F_{w,\tau}(x,\theta,V,s,a,s',\xi)], \quad (30)$$

$$\bar{G}_{\tau}(x,\theta,V) \triangleq \mathbb{E}_{s \sim d_{\rho}^{\pi\theta},a \sim \pi\theta(\cdot|s),s' \sim \mathcal{P}(\cdot|s,a)} [G_{\tau}(x,\theta,V,s,a,s')]. \quad (31)$$

• It can be seen from (13) that the following relationship holds for any x

$$\nabla_x \Phi_{w,\tau}(x) = \bar{D}_w(x, \pi_{\tau}^*(x), \pi_{w,\tau}^*(x)). \tag{32}$$

In addition, we have for any x, θ, τ

$$\bar{F}_{0,\tau}(x,\theta,V_{\tau}^{x,\pi_{\theta}}) = \nabla_{\theta} J_{\tau}(x,\pi_{\theta}), \tag{33}$$

$$\bar{F}_{w,\tau}(x,\theta,V_{\tau}^{x,\pi_{\theta}}) = w\nabla_{\theta}L_{w,\tau}(x,\pi_{\theta}). \tag{34}$$

• We define for $\tau > 0$

$$\ell_{\tau}(x) = J_{\tau}(x, \pi_{\tau}^{\star}(x)). \tag{35}$$

Note that ℓ_{τ} is not to be confused with Φ_{τ} defined in (7), which is the regularized upper-level objective.

• We recall the residuals defined in Section 4.1 and also define the residual in x_k as follows

$$\varepsilon_k^x = \|\nabla_x \Phi_{w_k, \tau_k}(x_k)\|^2, \quad \varepsilon_k^{\theta, \mathcal{L}} = w_k \Big(\mathcal{L}_{w_k, \tau_k}(x_k, \pi_{\theta_k^{\mathcal{L}}}) - \mathcal{L}_{w_k, \tau_k}(x_k, \pi_{w_k, \tau_k}(x_k)) \Big),$$

$$\varepsilon_k^{\theta} = J_{\tau_k}(x_k, \pi_{\tau_k}^{\star}(x_k)) - J_{\tau_k}(x_k, \pi_{\theta_k}), \quad \varepsilon_k^V = \|\hat{V}_k - V_{\tau_k}^{x_k, \pi_{\theta_k}}\|^2, \quad \varepsilon_k^{V, \mathcal{L}} = \|\hat{V}_k^{\mathcal{L}} - V_{\tau_k}^{x_k, \pi_{\theta_k^{\mathcal{L}}}}\|^2.$$
(36)

We also introduce a number of technical lemmas, which will be used in the proofs of the propositions and theorems. We defer the proofs of the lemmas to Appendix D.

Lemma 2 For any $k \ge 0$, we have

$$\tau_k - \tau_{k+1} \le \frac{8\tau_k}{3(k+1)}.$$

Lemma 2 derives a simple bound on the rate of change of the regularization weight τ_k .

Lemma 3 Define $L_V = \max\{\frac{2L_r|\mathcal{S}||\mathcal{A}|}{1-\gamma}, \frac{(12+8\log|\mathcal{A}|)\sqrt{|\mathcal{S}|}}{(1-\gamma)^3}\}$. We have for all $w, \tau \leq 1$ and x, x', θ, θ'

$$|J_{\tau}(x, \pi_{\theta}) - J_{\tau}(x', \pi_{\theta'})| \le L_{V}(\|x - x'\| + \|\theta - \theta'\|), \tag{37}$$

$$\|\nabla_{\theta} J_{\tau}(x, \pi_{\theta}) - \nabla_{\theta} J_{\tau}(x', \pi_{\theta'})\| \le L_{V}(\|x - x'\| + \|\theta - \theta'\|), \tag{38}$$

$$\|\nabla_x J_{\tau}(x,\pi) - \nabla_x J_{\tau}(x',\pi')\| \le L_V(\|x - x'\| + \|\pi - \pi'\|),\tag{39}$$

$$||V_{\tau}^{x,\pi_{\theta}} - V_{\tau}^{x',\pi_{\theta'}}|| \le L_{V}(||x - x'|| + ||\theta - \theta'||), \tag{40}$$

$$\|\nabla_{\theta} V_{\tau}^{x,\pi_{\theta}} - \nabla_{\theta} V_{\tau}^{x',\pi_{\theta'}}\| \le L_{V}(\|x - x'\| + \|\theta - \theta'\|),\tag{41}$$

$$\|\nabla_{\theta,\theta} \mathbb{E}_{s \sim d_{\sigma}^{\pi_{\theta}}} [E(\pi_{\theta}, s)]\| \le L_{V}. \tag{42}$$

In addition, there exists a bounded constant $L_{V,2}^3$ such that for all $\tau \leq 1$ and x, x', θ, θ'

$$\|\nabla_{x,\theta}^{2} J_{\tau}(x,\pi_{\theta}) - \nabla_{x,\theta}^{2} J_{\tau}(x',\pi_{\theta'})\| \leq L_{V,2}(\|x-x'\| + \|\theta-\theta'\|),$$

$$\|\nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta}) - \nabla_{\theta,\theta}^{2} J_{\tau}(x',\pi_{\theta'})\| \leq L_{V,2}(\|x-x'\| + \|\theta-\theta'\|).$$

Lemma 3 shows that the value functions/cumulative returns are Lipschitz, and have Lipschitz continuous gradients and Hessians.

Lemma 4 For any π , π' , we have

$$||d_{\rho}^{\pi} - d_{\rho}^{\pi'}|| \le \frac{\gamma}{1 - \gamma} ||\pi - \pi'||.$$

Lemma 4 shows that the occupancy measure is a Lipschitz function of the policy, a well-known result in the literature. We include the proof in Section D.4 for completeness.

³We skip showing the exact constant here, but note that it depends polynomially on the structural parameters of the problem.

Lemma 5 Recall that $B_V = \frac{1 + \log |A|}{1 - \gamma}$ is the entry-wise upper bound on the magnitude of the value function introduced in the paragraph after (3). We define the constants

$$B_D = 3L_r$$
, $B_F = 2(1+\gamma)B_V + 2\log|\mathcal{A}| + 2 + L_r$, $B_G = (1+\gamma)B_V + \log|\mathcal{A}| + 1$. (43)

Suppose that the regularization parameters satisfy $w \leq \frac{L_r}{L_f}$, $\tau \leq 1$. For all $x, \theta, \theta^{\mathcal{L}}$, $s, a, s', \bar{s}, \bar{a}, \xi$ and $V \in \mathbb{R}^{|\mathcal{S}|}$ satisfying $V(s) \leq B_V$, we have

$$||D_{w}(x, \pi_{\theta}, \pi_{\theta}\varepsilon, s, a, \bar{s}, \bar{a}, \xi)|| \leq \frac{B_{D}}{w},$$

$$||F_{w,\tau}(x, \theta, V, s, a, s', \xi)|| \leq B_{F},$$

$$||G_{\tau}(x, \theta, V, s, a, s')|| \leq B_{G}.$$

Lemma 6 We define the constants

$$L_D = 3L_r + \frac{2B_D}{1 - \gamma}, \quad L_F = 3L_r + 2\log|\mathcal{A}| + \frac{2 + B_F}{1 - \gamma} + 4, \quad L_G = L_r + \frac{B_G}{1 - \gamma} + \log|\mathcal{A}| + 2.$$
(44)

Suppose that the regularization parameters satisfy $w \leq \min\{\frac{L_r}{L_f}, \frac{B_D}{(1-\gamma)L_f}\}, \tau \leq 1$. We have for all $x_1, x_2, \pi_1, \pi_1^{\mathcal{L}}, \pi_2, \pi_2^{\mathcal{L}}$

$$\|\bar{D}_{w}(x_{1}, \pi_{1}, \pi_{1}^{\mathcal{L}}) - \bar{D}_{w}(x_{2}, \pi_{2}, \pi_{2}^{\mathcal{L}})\| \leq \frac{L_{D}}{w} \Big(\|x_{1} - x_{2}\| + \|\pi_{1} - \pi_{2}\| + \|\pi_{1}^{\mathcal{L}} - \pi_{2}^{\mathcal{L}}\| \Big)$$

$$\|\bar{F}_{w,\tau}(x_{1}, \theta_{1}, V_{1}) - \bar{F}_{w,\tau}(x_{2}, \theta_{2}, V_{2})\| \leq L_{F}(\|x_{1} - x_{2}\| + \|\theta_{1} - \theta_{2}\| + \|V_{1} - V_{2}\|),$$

$$\|\bar{G}_{\tau}(x_{1}, \theta_{1}, V_{1}) - \bar{G}_{\tau}(x_{2}, \theta_{2}, V_{2})\| \leq L_{G}(\|x_{1} - x_{2}\| + \|\theta_{1} - \theta_{2}\| + \|V_{1} - V_{2}\|).$$

Lemma 7 Recall the definition of $\pi_{w,\tau}^{\star}$ in Section 3.1. For any $w_1, w_2, \tau_1, \tau_2, x_1, x_2$, we have

$$\|\pi_{w_1,\tau_1}^{\star}(x_1) - \pi_{w_2,\tau_2}^{\star}(x_2)\| \leq \left(\frac{2L_f w_2}{C_L \tau_1} + \frac{2L_V}{C_L \tau_1}\right) \|x_1 - x_2\| + \frac{2L_f |w_1 - w_2|}{C_L \tau_1} + \frac{6|\tau_1 - \tau_2||\mathcal{S}|\log|\mathcal{A}|}{(1 - \gamma)C_L \tau_1}.$$

In addition, for any $w, \tau > 0$, we have

$$\|\pi_{\tau}^{\star}(x) - \pi_{w,\tau}^{\star}(x)\| \le \frac{2L_f w}{C_L \tau}.$$

Lemmas 6 and 7 show that the update operators introduced in (26)-(31) are (approximately) bounded and Lipschitz.

Lemma 8 Define $L_{\star} = \frac{4+8\log|\mathcal{A}|}{\underline{\sigma}(1-\gamma)^4}$. For any $\tau \geq 0$, we have

$$\|\pi_{\tau}^{\star}(x) - \pi^{\star}(x)\| \le L_{\star}\tau.$$

Lemma 8 bounds the distance between the regularized best response $\pi_{\tau}^{\star}(x)$ to $\pi^{\star}(x)$ defined in (5) by a linear function of τ .

Lemma 9 Define $L_L = L_V + L_f + \frac{L_V(C_L + 2L_V)}{C_L}$. We have for all $w, \tau \leq 1$ and x, x', θ, θ'

$$\|\nabla_{\theta} \mathcal{L}_{w,\tau}(x, \pi_{\theta}) - \nabla_{\theta} \mathcal{L}_{w,\tau}(x', \pi_{\theta'})\| \le \frac{L_L}{w} \|x - x'\| + \frac{L_L}{w} \|\theta - \theta'\|, \tag{45}$$

$$\|\nabla_{x}\mathcal{L}_{w,\tau}(x,\pi_{\theta}) - \nabla_{x}\mathcal{L}_{w,\tau}(x',\pi_{\theta'})\| \le \frac{L_{L}}{w\tau}\|x - x'\| + \frac{L_{L}}{w}\|\theta - \theta'\|. \tag{46}$$

Lemma 9 establishes the Lipschitz continuity of the gradients of $\mathcal{L}_{w,\tau}$.

Lemma 10 Recall the definition of ℓ_{τ} in (35). We have for all x_1, x_2

$$\|\nabla \ell_{\tau}(x_1) - \nabla \ell_{\tau}(x_2)\| \le \left(L_V + \frac{2L_V^2}{C_L \tau}\right) \|x_1 - x_2\|,$$

$$\|\nabla_x \Phi_{w,\tau}(x_1) - \nabla_x \Phi_{w,\tau}(x_2)\| \le \left(L_f + \frac{4L_f L_V}{C_L \tau} + \frac{4L_V^2}{C_L w \tau}\right) \|x_1 - x_2\|,$$

$$\|\nabla_x \Phi_{\tau}(x_1) - \nabla_x \Phi_{\tau}(x_2)\| \le \left(1 + \frac{2L_V}{C_L \tau}\right) \left(\frac{2L_f L_V}{C_L \tau} + \frac{2L_f L_V L_{V,2}}{\sigma C_L \tau} + \frac{2L_f L_V^2 L_{V,2}}{\sigma^2 C_L \tau} + \frac{2L_f L_V^2}{\sigma C_L \tau}\right) \|x_1 - x_2\|.$$

Particularly, if $w, \tau \leq 1$, we have

$$\|\nabla \ell_{\tau}(x_1) - \nabla \ell_{\tau}(x_2)\| \le \frac{L_{\Phi}}{\tau} \|x_1 - x_2\|,$$

$$\|\nabla_x \Phi_{w,\tau}(x_1) - \nabla_x \Phi_{w,\tau}(x_2)\| \le \frac{L_{\Phi}}{w\tau} \|x_1 - x_2\|,$$

where $L_{\Phi} = \max\{L_V + \frac{2L_V^2}{C_L}, L_f + \frac{4L_fL_V}{C_L} + \frac{4L_V^2}{C_L}, \frac{4L_fL_V}{C_L\tau} + \frac{4L_fL_VL_{V,2}}{\underline{\sigma}C_L\tau} + \frac{4L_fL_V^2L_{V,2}}{\underline{\sigma}^2C_L\tau} + \frac{4L_fL_V^2}{\underline{\sigma}^2C_L\tau}\}$

In addition, if $\tau \leq \frac{2L_V}{C_L}$, we have

$$\|\nabla_x \Phi_{\tau}(x_1) - \nabla_x \Phi_{\tau}(x_2)\| \le \frac{L_{\Phi}}{\tau} \|x_1 - x_2\|.$$

Lemma 10 establishes the Lipschitz continuity of the gradients of ℓ_{τ} , $\Phi_{w,\tau}$, and Φ_{τ} .

Lemma 11 We have for any $w, \tau > 0$

$$\begin{split} \|\nabla_{x}\Phi_{\tau}(x) - \nabla_{x}\Phi_{w,\tau}(x)\| &\leq \frac{4L_{f}L_{V}w}{C_{L}\underline{\sigma}\tau}(L_{f} + \frac{2L_{f}L_{V,2}}{C_{L}\tau}), \\ \|\nabla_{x}\Phi(x) - \nabla_{x}\Phi_{w,\tau}(x)\| &\leq \frac{4L_{f}L_{V}w}{C_{L}\underline{\sigma}\tau}(L_{f} + \frac{2L_{f}L_{V,2}}{C_{L}\tau}) \\ &+ \frac{L_{\star}L_{f}(L_{V} + 1) + L_{\star}L_{f}L_{V,2} + L_{\star}L_{V}L_{V,2} + L_{f}L_{V,2}(4 + 8\log|\mathcal{A}|)}{(1 - \gamma)^{4}\sigma^{2}}\tau. \end{split}$$

The lemma demonstrate how the magnitude of the difference between $\nabla_x \Phi_{\tau}(x)$ and $\nabla_x \Phi_{w,\tau}(x)$ and that between $\nabla_x \Phi(x)$ and $\nabla_x \Phi_{w,\tau}(x)$ scale with w and τ .

Lemma 12 For any $\tau > 0$, we have

$$\|\nabla_x \Phi_{\tau}(x) - \nabla_x \Phi(x)\| \le \frac{L_{\star} L_f L_V^2 L_{V,2} \tau}{\underline{\sigma}^2}.$$

This lemma shows that distance between $\nabla_x \Phi_{\tau}(x)$ and $\nabla_x \Phi(x)$ is bounded by a function linear in τ .

B Proof of Theorems

We study a slightly simplified variant of Algorithm 1, which is presented as Algorithm 2. The sole distinction between the two is that Algorithm 2 uses i.i.d. samples drawn from the stationary distribution, instead of continuously generated Markovian samples. Stochastic approximation and RL algorithms have been extensively analyzed under Markovian sampling [Zou et al., 2019, Wu et al., 2020], and it is well-established that Markovian samples affect convergence rates only by a logarithmic factor. This simplification enables us to concentrate on the novel aspects we introduce to bi-level RL, without being distracted by standard technical considerations related to Markovian samples.

Theorem 3 (Replicate of Theorem 1 under i.i.d. samples) Consider the iterates of Algorithm 2 under the step sizes and weights

$$\zeta_k = \frac{\zeta_0}{(k+1)^{c_\zeta}}, \quad \alpha_k = \frac{\alpha_0}{(k+1)^{c_\alpha}}, \quad \beta_k = \frac{\beta_0}{(k+1)^{c_\beta}}, \quad w_k = \frac{w_0}{(k+1)^{c_w}}, \quad \tau_k = \frac{\tau_0}{(k+1)^{c_\tau}},$$

Algorithm 2 Actor Critic Algorithm for Bi-Level RL

- 1: **Initialize:** control variable x_0 , policy parameters θ_0 and $\theta_0^{\mathcal{L}}$, value function estimates $\hat{V}_0, \hat{V}_0^{\mathcal{L}} \in \mathbb{R}^{|\mathcal{S}|}$
- 2: **for** iteration k = 0, 1, 2, ... **do**
- 3: Trajectory 1:

Get samples $s_k \sim d_{\rho}^{\pi_{\theta_k}}$, $a_k \sim \pi_{\theta_k}(\cdot \mid s_k)$, $s_k' \sim \mathcal{P}(\cdot \mid s_k, a_k)$. Receive reward $r_{x_k}(s_k, a_k)$.

4: Trajectory 2:

Get samples $\bar{s}_k \sim d_{\rho}^{\pi_{\theta_k^{\mathcal{L}}}}, \bar{a}_k \sim \pi_{\theta_k^{\mathcal{L}}}(\cdot \mid s_k), \bar{s}_k' \sim \mathcal{P}(\cdot \mid \bar{s}_k, a_k)$. Receive reward $r_{x_k}(\bar{s}_k, \bar{a}_k)$.

- 5: Observe/Obtain $\xi_k \sim \mu$
- 6: Control variable update:

$$x_{k+1} = x_k - \zeta_k \left(\widetilde{\nabla}_x f(x_k, \pi_{\theta_k}, \xi_k) + \frac{1}{w_k} \left(\nabla_x r_{x_k}(s_k, a_k) - \nabla_x r_{x_k}(\bar{s}_k, \bar{a}_k) \right) \right). \tag{47}$$

7: Policy update:

$$\theta_{k+1} = \theta_k + \alpha_k \Big(r_{x_k}(s_k, a_k) + \tau_k E(\pi_{\theta_k}, s_k) + \gamma \hat{V}_k(s_{k+1}) \Big) \nabla_{\theta} \log \pi_{\theta_k}(a_k \mid s_k),$$
(48)

$$\theta_{k+1}^{\mathcal{L}} = \theta_k^{\mathcal{L}} + \alpha_k \Big(\big(r_{x_k}(\bar{s}_k, \bar{a}_k) + \tau_k E(\pi_{\theta_k^{\mathcal{L}}}, \bar{s}_k) + \gamma \hat{V}_k^{\mathcal{L}}(\bar{s}_{k+1}) \big) \nabla_{\theta} \log \pi_{\theta_k^{\mathcal{L}}}(\bar{a}_k \mid \bar{s}_k)$$

$$- w_k \widetilde{\nabla}_{\theta} f(x_k, \pi_{\theta_k^{\mathcal{L}}}, \xi_k) \Big),$$
(49)

$$\pi_k = \operatorname{softmax}(\theta_k), \quad \pi_k^{\mathcal{L}} = \operatorname{softmax}(\theta_k^{\mathcal{L}}).$$
(50)

8: Value function update:

$$\hat{V}_{k+1} = \Pi_{B_V} \left(\hat{V}_k + \beta_k e_{s_k} \left(r_{x_k}(s_k, a_k) + \tau_k E(\pi_{\theta_k}, s_k) + \gamma \hat{V}_k(s_{k+1}) - \hat{V}_k(s_k) \right) \right),
\hat{V}_{k+1}^{\mathcal{L}} = \Pi_{B_V} \left(\hat{V}_k^{\mathcal{L}} + \beta_k e_{\bar{s}_k} \left(r_{x_k}(\bar{s}_k, \bar{a}_k) + \tau_k E(\pi_{\theta_k}^{\mathcal{L}}, \bar{s}_k) + \gamma \hat{V}_k^{\mathcal{L}}(\bar{s}_{k+1}) - \hat{V}_k^{\mathcal{L}}(\bar{s}_k) \right) \right).$$
(51)

9: end for

with $c_{\zeta} = \frac{9}{10}$, $c_{\alpha} = \frac{1}{2}$, $c_{\beta} = \frac{1}{2}$, $c_{w} = \frac{3}{20}$, $c_{\tau} = \frac{1}{20}$ and $\zeta_{0}, \alpha_{0}, \beta_{0}, w_{0}, \tau_{0}$ selected such that

$$\zeta_{0} \leq \alpha_{0} \leq \beta_{0} \leq w_{0} \leq \tau_{0} \leq 1,$$

$$\alpha_{0} \leq \min\left\{\frac{4L_{\Phi}}{3L_{V}^{2}}, \frac{4L_{\Phi}}{3L_{L}^{2}}, \frac{3B_{F}}{8}\right\}, \quad \beta_{0} \leq \min\left\{\frac{1-\gamma}{L_{G}^{2}}, \frac{(1-\gamma)L_{V}^{2}}{8|S|\log^{2}|\mathcal{A}|}\right\},$$

$$w_{0} \leq \min\left\{\frac{L_{r}}{L_{f}}, \frac{B_{D}}{(1-\gamma)L_{f}}\right\}, \quad \tau_{0} \leq \min\left\{\frac{L_{f}|\mathcal{S}|}{C_{L}}, \frac{2L_{V}}{C_{L}}, \frac{2L_{V,2}}{C_{L}}\right\},$$

$$\frac{\zeta_{0}}{\alpha_{0}} \leq \min\left\{\frac{C_{L}^{2}\tau_{0}^{2}}{1024(L_{L}^{2}+L_{V}^{2})}, \frac{C_{L}^{2}w_{0}^{2}\tau_{0}^{2}}{512L_{D}^{2}}, \frac{C_{L}^{2}w_{0}\tau_{0}^{2}}{128L_{D}L_{L}}, \frac{(1-\gamma)C_{L}^{2}\tau_{0}^{2}}{6144L_{V}^{2}L_{D}^{2}}, \frac{B_{F}}{B_{D}}\right\},$$

$$\frac{\alpha_{0}}{\beta_{0}} \leq \left\{\frac{1-\gamma}{2\sqrt{6}L_{V}L_{F}}, \frac{1-\gamma}{48L_{V}^{2}}, \frac{8(1-\gamma)(L_{L}^{2}+L_{V}^{2})}{3L_{V}^{2}C_{L}^{2}}, \sqrt{\frac{32B_{G}}{B_{F}^{2}(L_{L}+L_{V})}}, \frac{1-\gamma}{36B_{F}^{2}L_{V}\tau_{0}}\right\},$$

$$\frac{\alpha_{0}}{\beta_{0}^{2}} \leq \frac{4B_{G}}{11B_{F}^{2}L_{V}}, \quad \frac{\alpha_{0}}{\tau_{0}} \leq \frac{1}{L_{V}}.$$
(52)

⁴Note that the step sizes satisfying the conditions always exist and can be found in the order of $\tau_0, w_0, \beta_0, \alpha_0, \zeta_0$ – we first select w_0, τ_0 small enough to satisfy their upper bounds; then we select β_0 ; then we select α_0 with respect to β_0, w_0, τ_0 ; and finally we select ζ_0 with respect to $\alpha_0, \beta_0, w_0, \tau_0$.

Then, under Assumptions 1-5, we have for all $k \geq 0$,

$$\min_{t < k} \mathbb{E}[\|\nabla_x \Phi(x_t)\|^2] \leq \frac{40}{3\zeta_0(k+1)^{1/10}} \left(\Phi_{\tau_0}(x_0) + \varepsilon_0^\theta + \varepsilon_0^{\theta,\mathcal{L}} + \varepsilon_0^V + \varepsilon_0^{V,\mathcal{L}}\right) + \mathcal{O}\left(\frac{\log(k+1)}{(k+1)^{1/10}}\right).$$

Theorem 4 (Replicate of Theorem 2 under i.i.d. samples) Given any fixed regularization weight $\tau_0 \leq \min\{1, \frac{L_f|S|}{C_L}, \frac{2L_V}{C_L}, \frac{2L_{V,2}}{C_L}\}^5$, i.e. $\tau_k = \tau_0$ for all $k \geq 0$, consider the iterates of Algorithm 2 under the step sizes and penalty weight

$$\zeta_k = \frac{\zeta_0}{(k+1)^{c_\zeta}}, \quad \alpha_k = \frac{\alpha_0}{(k+1)^{c_\alpha}}, \quad \beta_k = \frac{\beta_0}{(k+1)^{c_\beta}}, \quad w_k = \frac{w_0}{(k+1)^{c_w}}$$

with $c_{\zeta}=\frac{2}{3}, c_{\alpha}=\frac{1}{2}, c_{\beta}=\frac{1}{2}, c_{w}=\frac{1}{6}$ and $\zeta_{0}, \alpha_{0}, \beta_{0}, w_{0}$ selected such that (52) holds. Then, under Assumptions 1-5, we have for all $k \geq 0$,

$$\min_{t < k} \mathbb{E}[\|\nabla_x \Phi_{\tau_0}(x_t)\|^2] \le \frac{20}{3\zeta_0(k+1)^{1/3}} \left(\Phi_{\tau_0}(x_0) + \varepsilon_0^{\theta} + \varepsilon_0^{\theta, \mathcal{L}} + \varepsilon_0^{V, \mathcal{L}} + \varepsilon_0^{V, \mathcal{L}}\right) + \mathcal{O}\left(\frac{\log(k+1)}{(k+1)^{1/3}}\right).$$

We break down the proofs of the theorems into two parts. First, in the following propositions, we individually establish the iteration-wise convergence of the upper-level decision variable (in Proposition 1 under decaying regularization and Proposition 2 under fixed regularization), policy iterates (in Propositions 3-4), and value function estimates (in Proposition 5). Then, in Sections B.1-B.2, we combine the convergence of these variables and bound their joint convergence through a coupled Lyapunov function. The proofs of the propositions are deferred to Section C.

Proposition 1 Under Assumptions 1-5 and step sizes satisfying (52), the iterates of Algorithm 2 satisfy for all $k \ge 0$

$$\begin{split} \frac{\zeta_k}{4} \mathbb{E}[\|\nabla_x \Phi(x_k)\|^2] &\leq \mathbb{E}[\Phi_{\tau_k}(x_k) - \Phi_{\tau_{k+1}}(x_{k+1})] + \frac{2L_D^2 \zeta_k}{w_k^2} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^{\star}(x_k)\|^2] \\ &\quad + \frac{2L_D^2 \zeta_k}{w_k^2} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^{\star}(x_k)\|^2] + \frac{L_{\star}^2 L_f^2 L_V^4 L_{V, 2}^2 \zeta_k \tau_k^2}{2\underline{\sigma}^4} \\ &\quad + \frac{256L_f^4 L_V^2 L_{V, 2}^2 \zeta_k w_k^2}{C_L^4 \underline{\sigma}^2 \tau_k^4} + \frac{B_D^2 L_\Phi \zeta_k^2}{2\tau_k w_k^2} + \frac{16L_f |\mathcal{S}| \log |\mathcal{A}|}{(1 - \gamma)C_L (k + 1)} \end{split}$$

Proposition 2 Under Assumptions 1-5 and step sizes satisfying (52), the iterates of Algorithm 2 satisfy for all $k \ge 0$

$$\begin{split} \frac{\zeta_k}{2} \mathbb{E}[\|\nabla_x \Phi_{\tau_0}(x_k)\|^2] &\leq \mathbb{E}[\Phi_{\tau_0}(x_k) - \Phi_{\tau_0}(x_{k+1})] + \frac{2L_D^2 \zeta_k}{w_k^2} \mathbb{E}[\|\pi_k - \pi_{\tau_0}^\star(x_k)\|^2] \\ &\qquad \qquad + \frac{2L_D^2 \zeta_k}{w_k^2} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k,\tau_0}^\star(x_k)\|^2] + C_{2,\tau_0} \zeta_k w_k^2 + \frac{B_D^2 L_{\Phi,\tau_0} \zeta_k^2}{2w_k^2}, \end{split}$$
 where $C_{2,\tau_0} = \left(\frac{4L_f L_V}{C_L \sigma_{\tau_0}} (L_f + \frac{2L_f L_{V,2}}{C_L \tau_0})\right)^2$ and $L_{\Phi,\tau_0} = (1 + \frac{2L_V}{C_L \tau_0}) \left(\frac{2L_f L_V}{C_L \tau_0} + \frac{2L_f L_V L_{V,2}}{\underline{\sigma} C_L \tau_0} + \frac{2L_f L_V L_{V,2}}{\underline{\sigma} C_L \tau_0} + \frac{2L_f L_V L_{V,2}}{\underline{\sigma} C_L \tau_0}\right). \end{split}$

Proposition 3 Under Assumptions 1-5 and step sizes satisfying (52), the iterates of Algorithm 2 satisfy for all $k \ge 0$

$$\mathbb{E}[\varepsilon_{k+1}^{\theta} - \varepsilon_{k}^{\theta}] \le -\frac{\alpha_{k}}{8} \mathbb{E}[\|\nabla_{\theta} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}] + 2L_{F}^{2} \alpha_{k} \mathbb{E}[\varepsilon_{k}^{V}] + \frac{32L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} \tau_{k}^{2}} \mathbb{E}[\varepsilon_{k}^{x}]$$

⁵Note that Propositions 3-5 use the Lipschitz continuity conditions established on operator/functions such as $L_{w,\tau}$ and ℓ_{τ} under $\tau \leq \min\{1, \frac{L_f|S|}{C_L}, \frac{2L_V}{C_L}, \frac{2L_{V,2}}{C_L}\}$, a condition imposed so that we can present the associated Lipschitz constants in a more concise form. As the proof of Theorem 4 is based on Propositions 3-5, we state that the result holds for this range of τ_0 . However, the same proof technique applies verbatim for any arbitrary $\tau_0 > 0$; only the values of the Lipschitz constants would change accordingly.

$$-\frac{C_L^2 \alpha_k \tau_k^2}{64} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^{\star}(x_k)\|^2] + \frac{64L_D^2 L_V^2 \zeta_k^2}{C_L^2 \alpha_k w_k^2 \tau_k^2} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^{\star}(x_k)\|^2 + \|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^{\star}(x_k)\|^2]$$

$$+ \frac{2B_D^2 L_\Phi \zeta_k^2}{w_k^2 \tau_k} + \frac{2B_D B_F L_V \zeta_k \alpha_k}{w_k} + \frac{B_F^2 L_V \alpha_k^2}{2} + \frac{16 \log |\mathcal{A}| \tau_k}{3(1 - \gamma)(k + 1)}.$$

Proposition 4 Under Assumptions 1-5 and step sizes satisfying (52), the iterates of Algorithm 2 satisfy for all $k \ge 0$

$$\mathbb{E}[\varepsilon_{k+1}^{\theta,\mathcal{L}} - \varepsilon_k^{\theta,\mathcal{L}}]$$

$$\leq -\frac{w_k^2 \alpha_k}{8} \mathbb{E}[\|\nabla_{\theta} \mathcal{L}_{w_k, \tau_k}(x_k, \pi_{\theta_k^{\mathcal{L}}})\|^2] + 2L_F^2 \alpha_k \mathbb{E}[\varepsilon_k^{V, \mathcal{L}}] + \frac{32L_L^2 \zeta_k^2}{C_L^2 \alpha_k \tau_k^2} \mathbb{E}[\varepsilon_k^x] \\ - \frac{C_L^2 \alpha_k \tau_k^2}{64} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^{\star}(x_k)\|^2] + \frac{64L_D^2 L_L^2 \zeta_k^2}{C_L^2 \alpha_k w_k^2 \tau_k^2} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^{\star}(x_k)\|^2 + \|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^{\star}(x_k)\|^2] \\ + \frac{2B_D^2 L_\Phi \zeta_k^2}{w_k^2 \tau_k} + \frac{2B_D B_F L_L \zeta_k \alpha_k}{w_k} + \frac{B_F^2 L_L \alpha_k^2}{2} + \frac{32 \log |\mathcal{A}| \tau_k}{3(1 - \gamma)(k + 1)}.$$

Proposition 5 (Value Function Convergence) *Under Assumptions 1-5 and step sizes satisfying* (52), the iterates of Algorithm 2 satisfy for all $k \ge 0$

$$\begin{split} \mathbb{E}[\varepsilon_{k+1}^{V}] &\leq \left(1 - \frac{(1 - \gamma)\beta_{k}}{4}\right) \mathbb{E}[\varepsilon_{k}^{V}] + \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{6L_{V}^{2}\alpha_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\|\nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}] + \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\varepsilon_{k}^{x}] \\ &+ \frac{22B_{F}^{2}L_{V}\tau_{0}\alpha_{k}^{2}}{\alpha_{0}} + \frac{64L_{V}^{2}\tau_{k}^{2}}{3(1 - \gamma)\beta_{k}(k + 1)^{2}} + 8B_{G}\beta_{k}^{2}, \\ \mathbb{E}[\varepsilon_{k+1}^{V,\mathcal{L}}] &\leq \left(1 - \frac{(1 - \gamma)\beta_{k}}{4}\right) \mathbb{E}[\varepsilon_{k}^{V,\mathcal{L}}] + \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{6L_{V}^{2}w_{k}^{2}\alpha_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}})\|^{2}] + \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\varepsilon_{k}^{x}] \\ &+ \frac{22B_{F}^{2}L_{V}\tau_{0}\alpha_{k}^{2}}{\alpha_{0}} + \frac{64L_{V}^{2}\tau_{k}^{2}}{3(1 - \gamma)\beta_{k}(k + 1)^{2}} + 8B_{G}\beta_{k}^{2}. \end{split}$$

B.1 Proof of Theorem 3

Combining the bounds in Propositions 1-5, we have for any $k \ge 0$

$$\frac{\zeta_k}{4} \mathbb{E}[\|\nabla_x \Phi(x_k)\|^2]$$

$$\leq \mathbb{E}[\Phi_{\tau_{k}}(x_{k}) - \Phi_{\tau_{k+1}}(x_{k+1})] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{L_{\star}^{2}L_{f}^{2}L_{V}^{4}L_{V,2}^{2}\zeta_{k}\tau_{k}^{2}}{2\underline{\sigma}^{4}} + \frac{256L_{f}^{4}L_{V}^{2}L_{V,2}^{2}\zeta_{k}w_{k}^{2}}{C_{L}^{4}\underline{\sigma}^{2}\tau_{k}^{4}} + \frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{2\tau_{k}w_{k}^{2}} + \frac{16L_{f}|\mathcal{S}|\log|\mathcal{A}|}{(1-\gamma)C_{L}(k+1)} + \mathbb{E}[\varepsilon_{k}^{\theta,\mathcal{L}} - \varepsilon_{k+1}^{\theta,\mathcal{L}}] - \frac{w_{k}^{2}\alpha_{k}}{8} \mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\|^{2}] + 2L_{F}^{2}\alpha_{k}\mathbb{E}[\varepsilon_{k}^{V,\mathcal{L}}] + \frac{32L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}} \mathbb{E}[\varepsilon_{k}^{x}] - \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{2B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{w_{k}^{2}\tau_{k}} + \frac{2B_{D}B_{F}L_{L}\zeta_{k}\alpha_{k}}{w_{k}} + \frac{B_{F}^{2}L_{L}\alpha_{k}^{2}}{2} + \frac{32\log|\mathcal{A}|\tau_{k}}{3(1-\gamma)(k+1)} + \mathbb{E}[\varepsilon_{k}^{\theta} - \varepsilon_{k+1}^{\theta}] - \frac{\alpha_{k}}{8}\mathbb{E}[\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}] + 2L_{F}^{2}\alpha_{k}\mathbb{E}[\varepsilon_{k}^{V}] + \frac{32L_{V}^{V}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}] - \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2}L_{V}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}]$$

$$\begin{split} &+\frac{2B_D^2L_{\Phi}\zeta_k^2}{w_k^2\tau_k} + \frac{2B_DB_FL_V\zeta_k\alpha_k}{w_k} + \frac{B_F^2L_V\alpha_k^2}{2} + \frac{16\log|A|\tau_k}{3(1-\gamma)(k+1)} \\ &- \mathbb{E}[\varepsilon_{k+1}^V] + \left(1 - \frac{(1-\gamma)\beta_k}{4}\right) \mathbb{E}[\varepsilon_k^V] + \frac{12L_V^2L_D^2\zeta_k^2}{(1-\gamma)\beta_k} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^*(x_k)\|^2 + \|\pi_k^{\mathcal{L}} - \pi_{w_k,\tau_k}^*(x_k)\|^2] \\ &+ \frac{6L_V^2\alpha_k^2}{(1-\gamma)\beta_k} \mathbb{E}[\|\nabla_{\theta}J_{\tau_k}(x_k, \pi_{\theta_k})\|^2] + \frac{6L_V^2\zeta_k^2}{(1-\gamma)\beta_k} \mathbb{E}[\varepsilon_k^k] \\ &+ \frac{22B_F^2L_V\tau_0\alpha_k^2}{\alpha_0} + \frac{64L_V^2\tau_k^2}{3(1-\gamma)\beta_k(k+1)^2} + 8B_G\beta_k^2 \\ &- \mathbb{E}[\varepsilon_{k+1}^V] + \left(1 - \frac{(1-\gamma)\beta_k}{4}\right) \mathbb{E}[\varepsilon_k^V.\mathcal{L}] + \frac{12L_V^2L_D^2\zeta_k^2}{(1-\gamma)\beta_k} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^*(x_k)\|^2 + \|\pi_k^{\mathcal{L}} - \pi_{w_k,\tau_k}^*(x_k)\|^2] \\ &+ \frac{6L_V^2\alpha_k^2\alpha_k^2}{(1-\gamma)\beta_k} \mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{w_k,\tau_k}(x_k, \pi_{\theta_k})\|^2] + \frac{6L_V^2\zeta_k^2}{(1-\gamma)\beta_k} \mathbb{E}[\varepsilon_k^x] \\ &+ \frac{22B_F^2L_V\tau_0\alpha_k^2}{\alpha_0} + \frac{64L_V^2\tau_k^2}{3(1-\gamma)\beta_k(k+1)^2} + 8B_G\beta_k^2 \\ \leq \mathbb{E}[\Phi_{\tau_k}(x_k) - \Phi_{\tau_{k+1}}(x_{k+1}) + \varepsilon_k^{\theta,\mathcal{L}} - \varepsilon_{k+1}^{\theta,\mathcal{L}} + \varepsilon_k^{\theta} - \varepsilon_{k+1}^{\theta} + \varepsilon_k^V - \varepsilon_{k+1}^V + \varepsilon_k^{V,\mathcal{L}} - \varepsilon_{k+1}^{V,\mathcal{L}}] \\ &+ \left(\frac{32(L_L^2 + L_V^2)\zeta_k^2}{C_L^2\alpha_k\tau_k^2} + \frac{12L_V^2\zeta_k^2}{(1-\gamma)\beta_k}\right) \mathbb{E}[\varepsilon_k^x] + \left(-\frac{\alpha_k}{8} + \frac{6L_V^2\alpha_k^2}{(1-\gamma)\beta_k}\right) \mathbb{E}[\|\|\nabla_{\theta}\mathcal{L}_{w_k,\tau_k}(x_k, \pi_{\theta_k})\|^2] \\ &+ \left(-\frac{w_k^2\alpha_k}{8} + \frac{6L_V^2w_k^2\alpha_k^2}{(1-\gamma)\beta_k}\right) \mathbb{E}[\|\|\nabla_{\theta}\mathcal{L}_{w_k,\tau_k}(x_k, \pi_{\theta_k})\|^2] \\ &+ \left(-\frac{C_L^2\alpha_k\tau_k^2}{64} + \frac{2L_D^2\zeta_k}{w_k^2} + \frac{128L_D^2L_L^2\zeta_k^2}{C_L^2\alpha_kw_k^2\tau_k^2} + \frac{24L_V^2L_D^2\zeta_k^2}{(1-\gamma)\beta_k}\right) \mathbb{E}[\|\pi_k - \pi_{\tau_k}^*(x_k)\|^2] \\ &+ \left(-\frac{C_L^2\alpha_k\tau_k^2}{64} + 2L_L^2\alpha_k\right) \mathbb{E}[\varepsilon_k^V] + \left(-\frac{(1-\gamma)\beta_k}{4} + 2L_L^2\alpha_k\right) \mathbb{E}[\varepsilon_k^V] \\ &+ \frac{256L_J^4L_V^2L_J^2\zeta_kw_k^2}{44} + \frac{128L_D^2L_L^2\zeta_k^2}{C_L^2\alpha_kw_k^2\tau_k^2} + \frac{24L_V^2L_D^2\zeta_k^2}{2w_k^2\tau_k} + \frac{4B_DB_FL_L\zeta_k\alpha_k}{w_k} + 48B_G\beta_k^2 \\ &+ \frac{128R_V^2\tau_k^2}{3(1-\gamma)\beta_k(k+1)^2} + \frac{16\log|A|\tau_k}{(1-\gamma)(k+1)} + \frac{16L_J|S|\log|A|}{(1-\gamma)C_L(k+1)}, \end{split}$$

where to get the second inequality we combine the terms $\frac{22B_F^2L_V\tau_0\alpha_k^2}{\alpha_0}$ and $8B_G\beta_k^2$ under the condition $\alpha_0 \leq \frac{4B_G\beta_0^2}{11B_F^2L_V}$, and the terms $\frac{B_F^2L_L\alpha_k^2}{2} + \frac{B_F^2L_V\alpha_k^2}{2}$ and $16B_G\beta_k^2$ under the condition $\frac{\alpha_k}{\beta_k} \leq \sqrt{\frac{32B_G}{B_F^2(L_L+L_V)}}$.

Note that the highlighted red coefficients in the inequality above are non-positive and the blue coefficients can be combined under the step size conditions $\zeta_0 \leq \beta_0$, $\tau_k \leq \frac{L_f|\mathcal{S}|}{C_L}$, and

$$\frac{\alpha_0}{\beta_0} \leq \left\{\frac{1-\gamma}{48L_V^2}, \frac{1-\gamma}{8L_F^2}, \frac{8(1-\gamma)(L_L^2 + L_V^2)}{3L_V^2C_L^2}\right\}, \quad \frac{\zeta_0}{\alpha_0} \leq \min\left\{\frac{C_L^2w_0^2\tau_0^2}{512L_D^2}, \frac{C_L^2w_0\tau_0^2}{128L_DL_L}, \frac{(1-\gamma)C_L^2\tau_0^2}{6144L_V^2L_D^2}\right\}$$

This allows us to simplify the inequality and obtain

$$\frac{\zeta_{k}}{4}\mathbb{E}[\|\nabla_{x}\Phi(x_{k})\|^{2}]$$

$$\leq \mathbb{E}[\Phi_{\tau_{k}}(x_{k}) - \Phi_{\tau_{k+1}}(x_{k+1}) + \varepsilon_{k}^{\theta,\mathcal{L}} - \varepsilon_{k+1}^{\theta,\mathcal{L}} + \varepsilon_{k}^{\theta} - \varepsilon_{k+1}^{\theta} + \varepsilon_{k}^{V} - \varepsilon_{k+1}^{V} + \varepsilon_{k}^{V,\mathcal{L}} - \varepsilon_{k+1}^{V,\mathcal{L}}]$$

$$+ \frac{64(L_{L}^{2} + L_{V}^{2})\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}] + \frac{256L_{f}^{4}L_{V}^{2}L_{V,2}^{2}\zeta_{k}w_{k}^{2}}{C_{L}^{4}\underline{\sigma}^{2}\tau_{k}^{4}} + \frac{L_{x}^{2}L_{f}^{2}L_{V}^{4}L_{V,2}^{2}\zeta_{k}\tau_{k}^{2}}{2\underline{\sigma}^{4}} + \frac{9B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{2w_{k}^{2}\tau_{k}}$$

$$+ \frac{4B_{D}B_{F}L_{L}\zeta_{k}\alpha_{k}}{w_{k}} + 48B_{G}\beta_{k}^{2} + \frac{128L_{V}^{2}\tau_{k}^{2}}{3(1-\gamma)\beta_{k}(k+1)^{2}} + \frac{32L_{f}|\mathcal{S}|\log|\mathcal{A}|}{(1-\gamma)C_{L}(k+1)}.$$
(53)

Recall the definition $\varepsilon_k^x = \|\nabla_x \Phi_{w_k, \tau_k}(x_k)\|^2$ in (36). We can relate the second term on the right hand side of (53) to $\|\nabla_x \Phi(x_k)\|^2$ using Lemma 11

$$\frac{64(L_L^2 + L_V^2)\zeta_k^2}{C_L^2\alpha_k\tau_k^2} \varepsilon_k^x \le \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2\alpha_k\tau_k^2} \|\nabla_x \Phi(x_k)\|^2 + \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2\alpha_k\tau_k^2} \|\nabla_x \Phi(x_k) - \nabla_x \Phi_{w_k,\tau_k}(x_k)\|^2 \\
\le \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2\alpha_k\tau_k^2} \|\nabla_x \Phi(x_k)\|^2 + \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2\alpha_k\tau_k^2} \left(\frac{4L_f L_V w_k}{C_L \underline{\sigma} \tau_k} (L_f + \frac{2L_f L_{V,2}}{C_L \tau_k})\right) \\
+ \frac{L_{\star} L_f (L_V + 1) + L_{\star} L_f L_{V,2} + L_{\star} L_V L_{V,2} + L_f L_{V,2} (4 + 8\log|\mathcal{A}|)}{(1 - \gamma)^4 \underline{\sigma}^2} \tau_k\right)^2 \\
\le \frac{\zeta_k}{8} \|\nabla_x \Phi(x_k)\|^2 + \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2\alpha_k\tau_k^2} \left(\frac{4L_f L_V w_k}{C_L \underline{\sigma} \tau_k} (L_f + \frac{2L_f L_{V,2}}{C_L \tau_k})\right) \\
+ \frac{L_{\star} L_f (L_V + 1) + L_{\star} L_f L_{V,2} + L_{\star} L_V L_{V,2} + L_f L_{V,2} (4 + 8\log|\mathcal{A}|)}{(1 - \gamma)^4 \underline{\sigma}^2} \tau_k\right)^2 \\
= \frac{\zeta_k}{8} \|\nabla_x \Phi(x_k)\|^2 + \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2\alpha_k\tau_k^2} \left(\frac{4L_f L_V w_k}{C_L \underline{\sigma} \tau_k} (L_f + \frac{2L_f L_{V,2}}{C_L \tau_k}) + C_1 \tau_k\right)^2, \tag{54}$$

where we derive the third inequality from the step size condition $\frac{\zeta_0}{\alpha_0} \leq \frac{C_L^2 \tau_0^2}{1024(L_L^2 + L_V^2)}$, and define in the last equation $C_1 = \frac{L_\star L_f(L_V+1) + L_\star L_f L_{V,2} + L_\star L_V L_{V,2} + L_f L_{V,2}(4+8\log|\mathcal{A}|)}{(1-\gamma)^4\underline{\sigma}^2}$.

Combining (53) and (54) and substituting in the step size decay rates,

$$\begin{split} \frac{\zeta_k}{8} \mathbb{E}[\|\nabla_x \Phi(x_k)\|^2] &\leq \mathbb{E}[\Phi_{\tau_k}(x_k) - \Phi_{\tau_{k+1}}(x_{k+1}) + \varepsilon_k^{\theta, \mathcal{L}} - \varepsilon_{k+1}^{\theta, \mathcal{L}} + \varepsilon_k^{\theta} - \varepsilon_{k+1}^{\theta} + \varepsilon_k^{V} - \varepsilon_{k+1}^{V} + \varepsilon_k^{V, \mathcal{L}} - \varepsilon_{k+1}^{V, \mathcal{L}}] \\ &\quad + \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2\alpha_k\tau_k^2} \left(\frac{4L_fL_Vw_k}{C_L\underline{\sigma}\tau_k}(L_f + \frac{2L_fL_{V,2}}{C_L\tau_k}) + C_1\tau_k\right)^2 \\ &\quad + \frac{256L_f^4L_V^2L_{V,2}^2\zeta_kw_k^2}{C_L^4\underline{\sigma}^2\tau_k^4} + \frac{L_\star^2L_f^2L_V^4L_{V,2}^2\zeta_k\tau_k^2}{2\underline{\sigma}^4} + \frac{9B_D^2L_\Phi\zeta_k^2}{2w_k^2\tau_k} \\ &\quad + \frac{4B_DB_FL_L\zeta_k\alpha_k}{w_k} + 48B_G\beta_k^2 + \frac{128L_V^2\tau_k^2}{3(1-\gamma)\beta_k(k+1)^2} + \frac{32L_f|\mathcal{S}|\log|\mathcal{A}|}{(1-\gamma)C_L(k+1)} \\ &\quad (55) \\ &\leq \mathbb{E}[\Phi_{\tau_k}(x_k) - \Phi_{\tau_{k+1}}(x_{k+1}) + \varepsilon_k^{\theta, \mathcal{L}} - \varepsilon_{k+1}^{\theta, \mathcal{L}} + \varepsilon_k^{\theta} - \varepsilon_{k+1}^{\theta} + \varepsilon_k^{V} - \varepsilon_{k+1}^{V, \mathcal{L}} + \varepsilon_k^{V, \mathcal{L}} - \varepsilon_{k+1}^{V, \mathcal{L}}] \\ &\quad + \mathcal{O}\left(\frac{1}{(k+1)}\right). \end{split}$$

Re-arranging the terms and summing over iterations,

$$\sum_{t=0}^{k-1} \frac{\zeta_0}{8(t+1)^{9/10}} \mathbb{E}[\|\nabla_x \Phi(x_t)\|^2] \le \Phi_{\tau_0}(x_0) + \varepsilon_0^{\theta} + \varepsilon_0^{\theta, \mathcal{L}} + \varepsilon_0^{V} + \varepsilon_0^{V, \mathcal{L}} + \mathcal{O}\left(\sum_{t=0}^{k-1} \frac{1}{(t+1)}\right). \tag{56}$$

The following inequalities on the summation of step sizes are standard results in the literature (the ones in our paper are specifically adapted from Lemma 3 of an earlier version of Zeng et al. [2024])

$$\sum_{t=0}^{k-1} \frac{1}{t+1} \leqslant \frac{\log(k+1)}{\log(2)},$$

$$\sum_{t=0}^{k} \frac{1}{(t+1)^u} \ge \frac{(1 - \frac{1}{2^{1-u}})(k+1)^{1-u}}{1-u}, \quad \forall u \in (0,1)$$

and with u = 9/10,

$$\sum_{t=0}^k \frac{1}{(t+1)^{9/10}} \ge \frac{0.06(k+1)^{1/10}}{1/10} = \frac{3(k+1)^{1/10}}{5}.$$

This allows us to further simplify (56)

$$\min_{t < k} \mathbb{E}[\|\nabla_x \Phi(x_t)\|^2] \le \frac{1}{\sum_{t=0}^{k-1} \frac{\zeta_0}{8(t+1)^{9/10}}} \sum_{t=0}^{k-1} \frac{\zeta_0}{8(t+1)^{9/10}} \mathbb{E}[\|\nabla_x \Phi(x_t)\|^2]
\le \frac{40}{3\zeta_0(k+1)^{1/10}} \left(\Phi_{\tau_0}(x_0) + \varepsilon_0^{\theta} + \varepsilon_0^{\theta, \mathcal{L}} + \varepsilon_0^{V, \mathcal{L}} + \mathcal{O}(\log(k+1))\right).$$

B.2 Proof of Theorem 4

We combine Propositions 3-2. Note that $\tau_k = \tau_0$. We can follow a line of analysis identical to what leads to (53) in the proof of Theorem 3 and show the following inequality

$$\frac{\zeta_{k}}{2}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{0}}(x_{k})\|^{2}] \leq \mathbb{E}[\Phi_{\tau_{0}}(x_{k}) - \Phi_{\tau_{0}}(x_{k+1}) + \varepsilon_{k}^{\theta,\mathcal{L}} - \varepsilon_{k+1}^{\theta,\mathcal{L}} + \varepsilon_{k}^{\theta} - \varepsilon_{k+1}^{\theta} + \varepsilon_{k}^{V} - \varepsilon_{k+1}^{V} + \varepsilon_{k}^{V,\mathcal{L}} - \varepsilon_{k+1}^{V,\mathcal{L}}] \\
+ \frac{64(L_{L}^{2} + L_{V}^{2})\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{0}^{2}} \mathbb{E}[\varepsilon_{k}^{x}] + C_{2,\tau_{0}}\zeta_{k}w_{k}^{2} + \frac{4B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{w_{k}^{2}\tau_{0}} + \frac{B_{D}^{2}L_{\Phi,\tau_{0}}\zeta_{k}^{2}}{2w_{k}^{2}} \\
+ \frac{4B_{D}B_{F}L_{L}\zeta_{k}\alpha_{k}}{w_{k}} + 48B_{G}\beta_{k}^{2} + \frac{128L_{V}^{2}\tau_{0}^{2}}{3(1-\gamma)\beta_{k}(k+1)^{2}} + \frac{48\log|\mathcal{A}|\tau_{0}}{3(1-\gamma)(k+1)}. \tag{57}$$

The step size conditions required to show (57) are

$$\frac{\alpha_0}{\beta_0} \le \left\{ \frac{1-\gamma}{48L_V^2}, \frac{1-\gamma}{8L_F^2}, \frac{8(1-\gamma)(L_L^2 + L_V^2)}{3L_V^2 C_L^2}, \sqrt{\frac{32B_G}{B_F^2 (L_L + L_V)}} \right\}, \quad \frac{\alpha_0}{\beta_0^2} \le \frac{4B_G}{11B_F^2 L_V},$$

$$\zeta_0 \le \beta_0, \quad \frac{\zeta_0}{\alpha_0} \le \min \left\{ \frac{C_L^2 w_0^2 \tau_0^2}{512L_D^2}, \frac{C_L^2 w_0 \tau_0^2}{128L_D L_L}, \frac{(1-\gamma)C_L^2 \tau_0^2}{6144L_V^2 L_D^2} \right\}.$$

Recall the definition $\varepsilon_k^x = \|\nabla_x \Phi_{w_k, \tau_k}(x_k)\|^2 = \|\nabla_x \Phi_{w_k, \tau_0}(x_k)\|^2$ in (36). We can relate the second term on the right hand side of (57) to $\|\nabla_x \Phi(x_k)\|^2$ using Lemma 11

$$\frac{64(L_L^2 + L_V^2)\zeta_k^2}{C_L^2 \alpha_k \tau_0^2} \varepsilon_k^x \leq \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2 \alpha_k \tau_0^2} \|\nabla_x \Phi_{\tau_0}(x_k)\|^2 + \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2 \alpha_k \tau_0^2} \|\nabla_x \Phi_{\tau_0}(x_k) - \nabla_x \Phi_{w_k, \tau_0}(x_k)\|^2 \\
\leq \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2 \alpha_k \tau_0^2} \|\nabla_x \Phi_{\tau_0}(x_k)\|^2 + \frac{128(L_L^2 + L_V^2)\zeta_k^2}{C_L^2 \alpha_k \tau_0^2} \left(\frac{4L_f L_V w_k}{C_L \sigma \tau_0} (L_f + \frac{2L_f L_{V,2}}{C_L \tau_0})\right)^2 \\
\leq \frac{\zeta_k}{4} \|\nabla_x \Phi_{\tau_0}(x_k)\|^2 + \frac{128(L_L^2 + L_V^2)\zeta_k^2 w_k^2}{C_L^2 \alpha_k \tau_0^2} \left(\frac{4L_f L_V}{C_L \sigma \tau_0} (L_f + \frac{2L_f L_{V,2}}{C_L \tau_0})\right)^2 \\
= \frac{\zeta_k}{4} \|\nabla_x \Phi_{\tau_0}(x_k)\|^2 + \frac{128(L_L^2 + L_V^2)C_{2,\tau_0}\zeta_k^2 w_k^2}{C_I^2 \alpha_k \tau_0^2}, \tag{58}$$

where the third inequality is due to the step size condition $\frac{\zeta_0}{\alpha_0} \leq \frac{C_L^2 \tau_0^2}{512(L_\tau^2 + L_\tau^2)}$.

Combining (57) and (58) and plugging in the step size decay rates,

$$\frac{\zeta_{k}}{4} \mathbb{E}[\|\nabla_{x} \Phi_{\tau_{0}}(x_{k})\|^{2}] \leq \mathbb{E}[\Phi_{\tau_{0}}(x_{k}) - \Phi_{\tau_{0}}(x_{k+1}) + \varepsilon_{k}^{\theta, \mathcal{L}} - \varepsilon_{k+1}^{\theta, \mathcal{L}} + \varepsilon_{k}^{\theta} - \varepsilon_{k+1}^{\theta} + \varepsilon_{k}^{V} - \varepsilon_{k+1}^{V} + \varepsilon_{k}^{V, \mathcal{L}} - \varepsilon_{k+1}^{V, \mathcal{L}}] \\
+ \frac{128(L_{L}^{2} + L_{V}^{2})C_{2,\tau_{0}}\zeta_{k}^{2}w_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{0}^{2}} + 2C_{2,\tau_{0}}\zeta_{k}w_{k}^{2} + \frac{4B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{w_{k}^{2}\tau_{0}} + \frac{B_{D}^{2}L_{\Phi,\tau_{0}}\zeta_{k}^{2}}{2w_{k}^{2}} \\
+ \frac{4B_{D}B_{F}L_{L}\zeta_{k}\alpha_{k}}{w_{k}} + 48B_{G}\beta_{k}^{2} + \frac{128L_{V}^{2}\tau_{0}^{2}}{3(1-\gamma)\beta_{k}(k+1)^{2}} + \frac{48\log|\mathcal{A}|\tau_{0}}{3(1-\gamma)(k+1)} \\
\leq \mathbb{E}[\Phi_{\tau_{0}}(x_{k}) - \Phi_{\tau_{0}}(x_{k+1}) + \varepsilon_{k}^{\theta, \mathcal{L}} - \varepsilon_{k+1}^{\theta, \mathcal{L}} + \varepsilon_{k}^{\theta} - \varepsilon_{k+1}^{\theta} + \varepsilon_{k}^{V} - \varepsilon_{k+1}^{V, \mathcal{L}} + \varepsilon_{k}^{V, \mathcal{L}} - \varepsilon_{k+1}^{V, \mathcal{L}}] \\
+ \mathcal{O}\left(\frac{\zeta_{k}^{2}w_{k}^{2}}{\alpha_{k}} + \zeta_{k}w_{k}^{2} + \frac{\zeta_{k}^{2}}{w_{k}^{2}} + \frac{\zeta_{k}\alpha_{k}}{w_{k}} + \beta_{k}^{2} + \frac{1}{\beta_{k}(k+1)^{2}} + \frac{1}{k+1}\right)$$

$$\leq \mathbb{E}[\Phi_{\tau_0}(x_k) - \Phi_{\tau_0}(x_{k+1}) + \varepsilon_k^{\theta,\mathcal{L}} - \varepsilon_{k+1}^{\theta,\mathcal{L}} + \varepsilon_k^{\theta} - \varepsilon_{k+1}^{\theta} + \varepsilon_k^{V} - \varepsilon_{k+1}^{V} + \varepsilon_k^{V,\mathcal{L}} - \varepsilon_{k+1}^{V,\mathcal{L}}] + \mathcal{O}\left(\frac{1}{k+1}\right).$$

Re-arranging the terms and summing over iterations,

$$\sum_{t=0}^{k-1} \frac{\zeta_0}{4(t+1)^{2/3}} \mathbb{E}[\|\nabla_x \Phi_{\tau_0}(x_t)\|^2] \le \Phi_{\tau_0}(x_0) + \varepsilon_0^{\theta} + \varepsilon_0^{\theta, \mathcal{L}} + \varepsilon_0^{V} + \varepsilon_0^{V, \mathcal{L}} + \mathcal{O}\left(\sum_{t=0}^{k-1} \frac{1}{(t+1)}\right).$$
(59)

The following inequalities on the summation of step sizes are standard results in the literature (the ones in our paper are specifically adapted from Lemma 3 of an earlier version of Zeng et al. [2024])

$$\sum_{t=0}^{k-1} \frac{1}{t+1} \leqslant \frac{\log(k+1)}{\log(2)},$$

$$\sum_{t=0}^{k} \frac{1}{(t+1)^u} \ge \frac{(1-\frac{1}{2^{1-u}})(k+1)^{1-u}}{1-u}, \quad \forall u \in (0,1)$$

and with u = 2/3,

$$\sum_{t=0}^{k} \frac{1}{(t+1)^{2/3}} \ge \frac{0.2(k+1)^{1/3}}{1/3} = \frac{3(k+1)^{1/3}}{5}.$$

This allows us to further simplify (59)

$$\begin{split} \min_{t < k} \mathbb{E}[\|\nabla_x \Phi_{\tau_0}(x_t)\|^2] &\leq \frac{1}{\sum_{t=0}^{k-1} \frac{\zeta_0}{4(t+1)^{2/3}}} \sum_{t=0}^{k-1} \frac{\zeta_0}{4(t+1)^{2/3}} \mathbb{E}[\|\nabla_x \Phi_{\tau_0}(x_t)\|^2] \\ &\leq \frac{20}{3\zeta_0(k+1)^{1/3}} \left(\Phi_{\tau_0}(x_0) + \varepsilon_0^{\theta} + \varepsilon_0^{\theta, \mathcal{L}} + \varepsilon_0^{V, \mathcal{L}} + \mathcal{O}(\log(k+1)) \right). \end{split}$$

C Proof of Propositions

C.1 Proof of Proposition 1

We know from Lemma 10 that under the step size condition $\tau_k \leq \frac{2L_V}{C_L}$, the objective Φ_{τ_k} has $\frac{L_{\Phi}}{\tau_k}$ -Lipschitz gradients. This implies

$$\begin{split} &\Phi_{\tau_{k}}(x_{k+1}) - \Phi_{\tau_{k}}(x_{k}) \\ &\leq \langle \nabla_{x}\Phi_{\tau_{k}}(x_{k}), x_{k+1} - x_{k} \rangle + \frac{L_{\Phi}}{2\tau_{k}} \|x_{k+1} - x_{k}\|^{2} \\ &= -\zeta_{k} \langle \nabla_{x}\Phi_{\tau_{k}}(x_{k}), D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \rangle + \frac{L_{\Phi}\zeta_{k}^{2}}{2\tau_{k}} \|D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \|^{2} \\ &= -\zeta_{k} \|\nabla_{x}\Phi_{\tau_{k}}(x_{k})\|^{2} + \zeta_{k} \langle \nabla_{x}\Phi_{\tau_{k}}(x_{k}), \nabla_{x}\Phi_{\tau_{k}}(x_{k}) - D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \rangle \\ &+ \frac{L_{\Phi}\zeta_{k}^{2}}{2\tau_{k}} \|D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \|^{2}. \end{split}$$

By the law of total expectation,

$$\mathbb{E}[\Phi_{\tau_{k}}(x_{k+1}) - \Phi_{\tau_{k}}(x_{k})]$$

$$\leq -\zeta_{k}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{k}}(x_{k})\|^{2}] + \frac{L_{\Phi}\zeta_{k}^{2}}{2\tau_{k}}\mathbb{E}[\|D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k})\|^{2}]$$

$$+ \zeta_{k} \mathbb{E}[\langle \nabla_{x} \Phi_{\tau_{k}}(x_{k}), \nabla_{x} \Phi_{\tau_{k}}(x_{k}) - \mathbb{E}[D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \mid \mathcal{F}_{k-1}] \rangle]$$

$$= -\zeta_{k} \mathbb{E}[\|\nabla_{x} \Phi_{\tau_{k}}(x_{k})\|^{2}] + \frac{L_{\Phi} \zeta_{k}^{2}}{2\tau_{k}} \mathbb{E}[\|D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k})\|^{2}]$$

$$+ \zeta_{k} \mathbb{E}[\langle \nabla_{x} \Phi_{\tau_{k}}(x_{k}), \nabla_{x} \Phi_{\tau_{k}}(x_{k}) - \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}) \rangle]$$

$$\leq -\zeta_{k} \mathbb{E}[\|\nabla_{x} \Phi_{\tau_{k}}(x_{k})\|^{2}] + \frac{B_{D}^{2} L_{\Phi} \zeta_{k}^{2}}{2\tau_{k} w_{k}^{2}} + \frac{\zeta_{k}}{2} \mathbb{E}[\|\nabla_{x} \Phi_{\tau_{k}}(x_{k})\|^{2}] + \frac{\zeta_{k}}{2} \mathbb{E}[\|\nabla_{x} \Phi_{\tau_{k}}(x_{k}) - \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}})\|^{2}]$$

$$= -\frac{\zeta_{k}}{2} \mathbb{E}[\|\nabla_{x} \Phi_{\tau_{k}}(x_{k})\|^{2}] + \frac{B_{D}^{2} L_{\Phi} \zeta_{k}^{2}}{2\tau_{k} w_{k}^{2}}$$

$$+ \frac{\zeta_{k}}{2} \mathbb{E}\left[\|\left(\nabla_{x} \Phi_{\tau_{k}}(x_{k}) - \nabla_{x} \Phi_{w_{k}, \tau_{k}}(x_{k})\right) + \left(\bar{D}_{w_{k}}(x_{k}, \pi_{\tau_{k}}^{\star}(x_{k}), \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})) - \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}})\right)\|^{2}\right]$$

$$\leq -\frac{\zeta_{k}}{2} \mathbb{E}[\|\nabla_{x} \Phi_{\tau_{k}}(x_{k})\|^{2}] + \zeta_{k} \mathbb{E}[\|\bar{D}_{w_{k}}(x_{k}, \pi_{\tau_{k}}^{\star}(x_{k}), \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})) - \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}})\|^{2}]$$

$$+ \zeta_{k} \mathbb{E}[\|\nabla_{x} \Phi_{\tau_{k}}(x_{k}) - \nabla_{x} \Phi_{w_{k}, \tau_{k}}(x_{k})\|^{2}] + \frac{B_{D}^{2} L_{\Phi} \zeta_{k}^{2}}{2\tau_{k} w_{k}^{2}},$$

$$(60)$$

where the second inequality applies Lemma 5 and the last equation follows from (32).

To bound the second term on the right hand side of (60), we apply Lemma 6

$$\|\bar{D}_{w_{k}}(x_{k}, \pi_{\tau_{k}}^{\star}(x_{k}), \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})) - \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}})\|^{2}$$

$$\leq \frac{L_{D}^{2}}{w_{k}^{2}} \Big(\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\| + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})\| \Big)^{2}$$

$$\leq \frac{2L_{D}^{2}}{w_{k}^{2}} \|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \frac{2L_{D}^{2}}{w_{k}^{2}} \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})\|^{2}.$$

$$(61)$$

For the third term of (60), we have from Lemma 11

$$\|\nabla_x \Phi_{\tau_k}(x_k) - \nabla_x \Phi_{w_k, \tau_k}(x_k)\|^2 \le \left(\frac{4L_f L_V w_k}{C_L \underline{\sigma} \tau_k} (L_f + \frac{2L_f L_{V,2}}{C_L \tau_k})\right)^2.$$
 (62)

Substituting (61) and (62) into (60),

$$\mathbb{E}[\Phi_{\tau_{k}}(x_{k+1}) - \Phi_{\tau_{k}}(x_{k})] \\
\leq -\frac{\zeta_{k}}{2} \mathbb{E}[\|\nabla_{x}\Phi_{\tau_{k}}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\
+ \zeta_{k} \left(\frac{4L_{f}L_{V}w_{k}}{C_{L}\underline{\sigma}\tau_{k}}(L_{f} + \frac{2L_{f}L_{V,2}}{C_{L}\tau_{k}})\right)^{2} + \frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{2\tau_{k}w_{k}^{2}}.$$
(63)

Under the choice of step size $\tau_k \leq \frac{2L_{V,2}}{C_L}$, we can further simplify (63)

$$\mathbb{E}[\Phi_{\tau_{k}}(x_{k+1}) - \Phi_{\tau_{k}}(x_{k})]$$

$$\leq -\frac{\zeta_{k}}{2}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{k}}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}]$$

$$+ \zeta_{k}(\frac{4L_{f}L_{V}w_{k}}{C_{L}\underline{\sigma}\tau_{k}} \cdot \frac{4L_{f}L_{V,2}}{C_{L}\tau_{k}})^{2} + \frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{2\tau_{k}w_{k}^{2}}$$

$$= -\frac{\zeta_{k}}{2}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{k}}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}]$$

$$+ \frac{256L_{f}^{4}L_{V}^{2}L_{V,2}^{2}\zeta_{k}w_{k}^{2}}{C_{f}^{4}\sigma^{2}\tau_{k}^{4}} + \frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{2\tau_{k}w_{k}^{2}}.$$
(64)

The next step is to bridge the gap between $\Phi_{\tau_k}(x_{k+1})$ and $\Phi_{\tau_{k+1}}(x_{k+1})$. By the definition of Φ_{τ} in (7)

$$\Phi_{\tau_{k+1}}(x_{k+1}) - \Phi_{\tau_k}(x_{k+1}) = f(x_{k+1}, \pi^\star_{\tau_{k+1}}(x_{k+1})) - f(x_{k+1}, \pi^\star_{\tau_k}(x_{k+1}))$$

$$\leq L_f \| \pi_{\tau_{k+1}}^*(x_{k+1}) - \pi_{\tau_k}^*(x_{k+1}) \|$$

$$\leq L_f \cdot \frac{6(\tau_k - \tau_{k+1})|\mathcal{S}|\log|\mathcal{A}|}{(1 - \gamma)C_L \tau_k}$$

$$\leq \frac{16L_f |\mathcal{S}|\log|\mathcal{A}|}{(1 - \gamma)C_L (k+1)}, \tag{65}$$

where the first inequality is due to Assumption 3, and the second inequality is due to Lemma 7, and the last inequality applies Lemma 2.

Finally, we bridge the gap between $\|\nabla_x \Phi_{\tau_k}(x_k)\|^2$ and $\|\nabla_x \Phi(x_k)\|^2$ by invoking Lemma 12

$$\|\nabla_x \Phi_{\tau_k}(x) - \nabla_x \Phi(x)\| \le \frac{L_{\star} L_f L_V^2 L_{V,2} \tau_k}{\sigma^2}.$$
 (66)

Combining (64)-(66),

$$\begin{split} &\mathbb{E}[\Phi_{\tau_{k+1}}(x_{k+1}) - \Phi_{\tau_{k}}(x_{k})] \\ &= \mathbb{E}[\Phi_{\tau_{k}}(x_{k+1}) - \Phi_{\tau_{k}}(x_{k})] + \mathbb{E}[\Phi_{\tau_{k+1}}(x_{k+1}) - \Phi_{\tau_{k}}(x_{k+1})] \\ &\leq -\frac{\zeta_{k}}{2}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{k}}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &\quad + \frac{256L_{f}^{4}L_{V}^{2}L_{V,2}^{2}\zeta_{k}w_{k}^{2}}{C_{L}^{4}\sigma^{2}\tau_{k}^{4}} + \frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{2\tau_{k}w_{k}^{2}} + \frac{16L_{f}|\mathcal{S}|\log|\mathcal{A}|}{(1-\gamma)C_{L}(k+1)} \\ &\leq -\frac{\zeta_{k}}{4}\mathbb{E}[\|\nabla_{x}\Phi(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{2L_{D}^{2}\zeta_{k}}{w_{k}^{2}}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &\quad + \frac{L_{\star}^{2}L_{f}^{2}L_{V}^{4}L_{V,2}^{2}\zeta_{k}\tau_{k}^{2}}{2\sigma^{4}} + \frac{256L_{f}^{4}L_{V}^{2}L_{V,2}^{2}\zeta_{k}w_{k}^{2}}{C_{f}^{4}\sigma^{2}\tau_{k}^{4}} + \frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{2\tau_{k}w_{k}^{2}} + \frac{16L_{f}|\mathcal{S}|\log|\mathcal{A}|}{(1-\gamma)C_{L}(k+1)}, \end{split}$$

where the last inequality follows from the simple fact that $-\frac{a^2}{2} \le -\frac{b^2}{4} + \frac{(a-b)^2}{2}$ for any scalar a,b.

C.2 Proof of Proposition 3

The proof depends on an intermediate result that bounds an important cross term. We state it in the lemma below and defer its proof to Section D.13.

Lemma 13 Under the assumptions and step sizes of Proposition 3, we have for all $k \ge 0$

$$\begin{split} & \mathbb{E}[-\langle \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}) - \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{\tau_{k}}^{\star}(x_{k})}), x_{k+1} - x_{k} \rangle \rangle] \\ & \leq \frac{C_{L}^{2} \alpha_{k} \tau_{k}^{2}}{64} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] \\ & \quad + \frac{64 L_{D}^{2} L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{32 L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E}[\varepsilon_{k}^{x}]. \end{split}$$

We now proceed to the proof of Proposition 3. We consider the following decomposition and bound each term on the right hand side individually.

$$-J_{\tau_{k+1}}(x_{k+1}, \pi_{\theta_{k+1}}) + J_{\tau_k}(x_k, \pi_{\theta_k})$$

$$= \left(-J_{\tau_k}(x_{k+1}, \pi_{\theta_{k+1}}) + J_{\tau_k}(x_{k+1}, \pi_{\theta_k})\right) + \left(-J_{\tau_k}(x_{k+1}, \pi_{\theta_k}) + J_{\tau_k}(x_k, \pi_{\theta_k})\right)$$

$$+ \left(-J_{\tau_{k+1}}(x_{k+1}, \pi_{\theta_{k+1}}) + J_{\tau_k}(x_{k+1}, \pi_{\theta_{k+1}})\right). \tag{67}$$

Bound the First Term of (67). As J_{τ} has L_V -Lipschitz gradients,

$$-J_{\tau_k}(x_{k+1},\pi_{\theta_{k+1}})+J_{\tau_k}(x_{k+1},\pi_{\theta_k})$$

$$\leq \langle -\nabla_{\theta} J_{\tau_{k}}(x_{k+1}, \pi_{\theta_{k}}), \theta_{k+1} - \theta_{k} \rangle + \frac{L_{V}}{2} \|\theta_{k+1} - \theta_{k}\|^{2}
= -\alpha_{k} \langle \nabla_{\theta} J_{\tau_{k}}(x_{k+1}, \pi_{\theta_{k}}), F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s'_{k}) \rangle + \frac{L_{V} \alpha_{k}^{2}}{2} \|F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s'_{k})\|^{2}
= -\alpha_{k} \langle \nabla_{\theta} J_{\tau_{k}}(x_{k+1}, \pi_{\theta_{k}}), F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s'_{k}) - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) \rangle
- \alpha_{k} \langle \nabla_{\theta} J_{\tau_{k}}(x_{k+1}, \pi_{\theta_{k}}), \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}) \rangle
- \alpha_{k} \langle \nabla_{\theta} J_{\tau_{k}}(x_{k+1}, \pi_{\theta_{k}}), \nabla_{\theta} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}) \rangle + \frac{L_{V} \alpha_{k}^{2}}{2} \|F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s'_{k})\|^{2}, \tag{68}$$

where the final equation follows from $\nabla_{\theta} J_{\tau_k}(x_k, \pi_{\theta_k}) = \bar{F}_{0,\tau_k}(x_k, \theta_k, V_{\tau_k}^{x_k, \pi_{\theta_k}})$.

To bound the first term of (68),

$$-\alpha_{k}\mathbb{E}[\langle \nabla_{\theta}J_{\tau_{k}}(x_{k+1}, \pi_{\theta_{k}}), F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s'_{k}) - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\rangle]$$

$$= -\alpha_{k}\mathbb{E}[\langle \nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}), \mathbb{E}[F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s'_{k}) \mid \mathcal{F}_{k-1}] - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\rangle]$$

$$+ \alpha_{k}\mathbb{E}[\langle \nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}) - \nabla_{\theta}J_{\tau_{k}}(x_{k+1}, \pi_{\theta_{k}}), F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s'_{k}) - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\rangle]$$

$$= \alpha_{k}\mathbb{E}[\langle \nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}) - \nabla_{\theta}J_{\tau_{k}}(x_{k+1}, \pi_{\theta_{k}}), F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s'_{k}) - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\rangle]$$

$$\leq \alpha_{k} \cdot L_{V}\mathbb{E}[||x_{k+1} - x_{k}||] \cdot 2B_{F}$$

$$\leq 2B_{F}L_{V}\alpha_{k} \cdot \frac{B_{D}\zeta_{k}}{w_{k}}$$

$$= \frac{2B_{D}B_{F}L_{V}\zeta_{k}\alpha_{k}}{w_{k}}, \qquad (69)$$

where the second inequality follows from Lemma 5.

For the second term of (68),

$$-\alpha_{k}\langle\nabla_{\theta}J_{\tau_{k}}(x_{k+1},\pi_{\theta_{k}}),\bar{F}_{0,\tau_{k}}(x_{k},\theta_{k},\hat{V}_{k}) - \bar{F}_{0,\tau_{k}}(x_{k},\theta_{k},V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}})\rangle$$

$$\leq \frac{\alpha_{k}}{8}\|\nabla_{\theta}J_{\tau_{k}}(x_{k+1},\pi_{\theta_{k}})\|^{2} + 2\alpha_{k}\|\bar{F}_{0,\tau_{k}}(x_{k},\theta_{k},\hat{V}_{k}) - \bar{F}_{0,\tau_{k}}(x_{k},\theta_{k},V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}})\|^{2}$$

$$\leq \frac{\alpha_{k}}{4}\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2} + \frac{\alpha_{k}}{4}\|\nabla_{\theta}J_{\tau_{k}}(x_{k+1},\pi_{\theta_{k}}) - \nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}$$

$$+ 2\alpha_{k}\|\bar{F}_{0,\tau_{k}}(x_{k},\theta_{k},\hat{V}_{k}) - \bar{F}_{0,\tau_{k}}(x_{k},\theta_{k},V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}})\|^{2}$$

$$\leq \frac{\alpha_{k}}{4}\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2} + \frac{L_{V}^{2}\alpha_{k}}{4}\left(\|x_{k+1} - x_{k}\|^{2}\right) + 2L_{F}^{2}\alpha_{k}\|V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}} - \hat{V}_{k}\|^{2}$$

$$\leq \frac{\alpha_{k}}{4}\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2} + \frac{L_{V}^{2}\alpha_{k}}{4} \cdot \left(\frac{B_{D}\zeta_{k}}{w_{k}}\right)^{2} + 2L_{F}^{2}\alpha_{k}\varepsilon_{k}^{V}$$

$$\leq \frac{\alpha_{k}}{4}\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2} + \frac{B_{D}^{2}L_{V}^{2}\zeta_{k}^{2}\alpha_{k}}{4w_{k}^{2}} + 2L_{F}^{2}\alpha_{k}\varepsilon_{k}^{V},$$

$$(70)$$

where the third inequality is a result of the Lipschitz continuity of J_{τ_k} and \bar{F}_{0,τ_k}

For the third term of (68),

$$-\alpha_{k}\langle\nabla_{\theta}J_{\tau_{k}}(x_{k+1},\pi_{\theta_{k}}),\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\rangle$$

$$=-\alpha_{k}\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}+\alpha_{k}\langle\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})-\nabla_{\theta}J_{\tau_{k}}(x_{k+1},\pi_{\theta_{k}}),\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\rangle$$

$$\leq -\frac{\alpha_{k}}{2}\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}+\frac{\alpha_{k}}{2}\|\nabla_{\theta}J_{\tau_{k}}(x_{k+1},\pi_{\theta_{k}})-\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}$$

$$\leq -\frac{\alpha_{k}}{2}\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}+\frac{L_{V}^{2}\alpha_{k}}{2}\left(\|x_{k+1}-x_{k}\|^{2}\right)$$

$$\leq -\frac{\alpha_{k}}{2}\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}+\frac{L_{V}^{2}\alpha_{k}}{2}\cdot\left(\frac{B_{D}\zeta_{k}}{w_{k}}\right)^{2}$$

$$\leq -\frac{\alpha_{k}}{2}\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}+\frac{B_{D}^{2}L_{V}^{2}\zeta_{k}^{2}\alpha_{k}}{2w_{k}^{2}},$$

$$(71)$$

where the second inequality follows from the L_V -smoothness continuity of J_τ from (38) of Lemma 3. Substituting (69)-(71) into (68),

$$\mathbb{E}\left[-J_{\tau_{k+1}}(x_{k+1}, \pi_{\theta_{k+1}}) + J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\right] \\
\leq \frac{2B_{D}B_{F}L_{V}\zeta_{k}\alpha_{k}}{w_{k}} + \frac{\alpha_{k}}{4}\mathbb{E}\left[\|\nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}\right] + \frac{B_{D}^{2}L_{V}^{2}\zeta_{k}^{2}\alpha_{k}}{4w_{k}^{2}} + 2L_{F}^{2}\alpha_{k}\mathbb{E}\left[\varepsilon_{k}^{V}\right] \\
- \frac{\alpha_{k}}{2}\mathbb{E}\left[\|\nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}\right] + \frac{B_{D}^{2}L_{V}^{2}\zeta_{k}^{2}\alpha_{k}}{2w_{k}^{2}} + \frac{L_{V}\alpha_{k}^{2}}{2}\mathbb{E}\left[\|F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}')\|^{2}\right] \\
\leq -\frac{\alpha_{k}}{4}\mathbb{E}\left[\|\nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}\right] + 2L_{F}^{2}\alpha_{k}\mathbb{E}\left[\varepsilon_{k}^{V}\right] + \frac{3B_{D}^{2}L_{V}^{2}\zeta_{k}^{2}\alpha_{k}}{4w_{k}^{2}} + \frac{2B_{D}B_{F}L_{V}\zeta_{k}\alpha_{k}}{w_{k}} + \frac{B_{F}^{2}L_{V}\alpha_{k}^{2}}{2} \\
\leq -\frac{\alpha_{k}}{8}\mathbb{E}\left[\|\nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}\right] + 2L_{F}^{2}\alpha_{k}\mathbb{E}\left[\varepsilon_{k}^{V}\right] + \frac{3B_{D}^{2}L_{V}^{2}\zeta_{k}^{2}\alpha_{k}}{4w_{k}^{2}} + \frac{2B_{D}B_{F}L_{V}\zeta_{k}\alpha_{k}}{w_{k}} + \frac{B_{F}^{2}L_{V}\alpha_{k}^{2}}{2} \\
-\frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{32}\mathbb{E}\left[\|\pi_{\theta_{k}} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}\right], \tag{72}$$

where the last inequality plugs in the relationship

$$\|\nabla_{\theta} J_{\tau_k}(x_k, \pi_{\theta_k})\| \ge \frac{C_L \tau_k}{2} \|\pi_{\theta_k} - \pi_{\tau_k}^{\star}(x_k)\|.$$
 (73)

Note that $J_{\tau}(x,\pi) = \lim_{w\to 0} w \mathcal{L}_{w,\tau}(x,\pi)$, which implies that (73) follows from (25).

Bound the Second Term of (67). We use $\theta_{\tau}^{\star}(x)$ to denote a softmax parameter that encodes the policy $\pi_{\tau}^{\star}(x)$. Again, as J_{τ} has L_V -Lipschitz gradients,

$$\begin{split} &-J_{\tau_{k}}(x_{k+1},\pi_{\theta_{k}}) + J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) \\ &\leq -\langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}), x_{k+1} - x_{k} \rangle + \frac{L_{V}}{2} \|x_{k+1} - x_{k}\|^{2} \\ &= -\langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) - \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}), x_{k+1} - x_{k} \rangle + \frac{L_{V}}{2} \|x_{k+1} - x_{k}\|^{2} \\ &- \langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}), x_{k+1} - x_{k} \rangle \\ &= -\langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) - \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}), x_{k+1} - x_{k} \rangle + \frac{L_{V}}{2} \|x_{k+1} - x_{k}\|^{2} \\ &- \langle \nabla_{x}\ell_{\tau_{k}}(x_{k}), x_{k+1} - x_{k} \rangle \\ &\leq -\langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) - \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}), x_{k+1} - x_{k} \rangle + \frac{L_{V}}{2} \|x_{k+1} - x_{k}\|^{2} \\ &+ \left(-\ell_{\tau_{k}}(x_{k+1}) + \ell_{\tau_{k}}(x_{k}) \right) + \frac{L_{\Phi}}{2\tau_{k}} \|x_{k+1} - x_{k} \rangle \\ &+ \left(-J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) - \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}), x_{k+1} - x_{k} \rangle \\ &+ \left(-J_{\tau_{k}}(x_{k+1},\pi_{\tau_{k}}^{\star}(x_{k+1})) + J_{\tau_{k}}(x_{k},\pi_{\tau_{k}}^{\star}(x_{k})) \right) + \frac{L_{\Phi}}{\tau_{k}} \|x_{k+1} - x_{k}\|^{2}, \end{split}$$

where the last inequality is due to $L_V \leq L_\Phi$ and the step size condition $\tau_k \leq 1$, and the second equation uses the relationship $\nabla_x J_{\tau_k}(x_k, \pi_{\theta^\star_{\tau_k}(x_k)}) = \nabla_x \ell_{\tau_k}(x_k)$, which is due to $\nabla_\theta J_{\tau_k}(x_k, \pi_{\theta^\star_{\tau_k}(x_k)}) = 0$ by the first-order optimality condition. The second inequality is due to the fact that ℓ_τ is $\frac{L_\Phi}{\tau}$ -smooth when $\tau \leq 1$ (established in Lemma 10) and that for an L-smooth function f, we have

$$f(y) - f(x) \le \langle \nabla f(x), y - x \rangle + \frac{L}{2} ||x - y||^2.$$

Taking the expectation and plugging in the result from Lemma 13,

$$\mathbb{E}[-J_{\tau_k}(x_{k+1}, \pi_{\theta_k}) + J_{\tau_k}(x_k, \pi_{\theta_k})]$$

$$\leq \frac{C_L^2 \alpha_k \tau_k^2}{64} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^{\star}(x_k)\|^2]$$

$$+\frac{64L_{D}^{2}L_{V}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k}-\pi_{\tau_{k}}^{\star}(x_{k})\|^{2}+\|\pi_{k}^{\mathcal{L}}-\pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}]+\frac{32L_{V}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}]$$

$$+\mathbb{E}[-J_{\tau_{k}}(x_{k+1},\pi_{\tau_{k}}^{\star}(x_{k+1}))+J_{\tau_{k}}(x_{k},\pi_{\tau_{k}}^{\star}(x_{k}))]+\frac{L_{\Phi}}{\tau_{k}}\mathbb{E}[\|x_{k+1}-x_{k}\|^{2}]$$

$$\leq \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64}\mathbb{E}[\|\pi_{k}-\pi_{\tau_{k}}^{\star}(x_{k})\|^{2}]+\frac{64L_{D}^{2}L_{V}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k}-\pi_{\tau_{k}}^{\star}(x_{k})\|^{2}+\|\pi_{k}^{\mathcal{L}}-\pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}]$$

$$+\mathbb{E}[-J_{\tau_{k}}(x_{k+1},\pi_{\tau_{k}}^{\star}(x_{k+1}))+J_{\tau_{k}}(x_{k},\pi_{\tau_{k}}^{\star}(x_{k}))]+\frac{32L_{V}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}]+\frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{w_{k}^{2}\tau_{k}}, \tag{74}$$

where the last inequality follows from $\|x_{k+1} - x_k\| \leq \frac{B_D \zeta_k}{w_k}$.

Bound the Third Term of (67).

$$-J_{\tau_{k+1}}(x_{k+1}, \pi_{\theta_{k+1}}) + J_{\tau_{k}}(x_{k+1}, \pi_{\theta_{k+1}}) = \frac{\tau_{k} - \tau_{k+1}}{1 - \gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi_{k+1}^{\mathcal{L}}}} [E(\pi_{k+1}^{\mathcal{L}}, s)]$$

$$\leq \frac{\log |\mathcal{A}|(\tau_{k} - \tau_{k+1})}{(1 - \gamma)}$$

$$\leq \frac{8 \log |\mathcal{A}| \tau_{k}}{3(1 - \gamma)(k + 1)}, \tag{75}$$

where the second inequality follows from Lemma 2.

Collecting the bounds in (72)-(75) and plugging them into (67), we get

$$\begin{split} &\mathbb{E}[-J_{\tau_{k+1}}(x_{k+1},\pi_{\theta_{k+1}}) + J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})] \\ &\leq -\frac{\alpha_{k}}{8}\mathbb{E}[\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}] + 2L_{F}^{2}\alpha_{k}\mathbb{E}[\varepsilon_{k}^{V}] + \frac{3B_{D}^{2}L_{V}^{2}\zeta_{k}^{2}\alpha_{k}}{4w_{k}^{2}} + \frac{2B_{D}B_{F}L_{V}\zeta_{k}\alpha_{k}}{w_{k}} + \frac{B_{F}^{2}L_{V}\alpha_{k}^{2}}{2} \\ &\quad - \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{32}\mathbb{E}[\|\pi_{\theta_{k}} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &\quad + \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2}L_{V}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &\quad + \mathbb{E}[-J_{\tau_{k}}(x_{k+1}, \pi_{\tau_{k}}^{\star}(x_{k+1})) + J_{\tau_{k}}(x_{k}, \pi_{\tau_{k}}^{\star}(x_{k}))] + \frac{32L_{V}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}] + \frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{w_{k}^{2}\tau_{k}} \\ &\quad + \frac{8\log|\mathcal{A}|\tau_{k}}{3(1-\gamma)(k+1)} \\ \leq -\frac{\alpha_{k}}{8}\mathbb{E}[\|\nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}] + 2L_{F}^{2}\alpha_{k}\mathbb{E}[\varepsilon_{k}^{V}] + \frac{32L_{V}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}] \\ &\quad - \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2}L_{V}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &\quad + \frac{2B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{64} + \frac{2B_{D}B_{F}L_{V}\zeta_{k}\alpha_{k}}{w_{k}} + \frac{B_{F}^{2}L_{V}\alpha_{k}^{2}}{2} + \frac{8\log|\mathcal{A}|\tau_{k}}{3(1-\gamma)(k+1)} \\ &\quad + \mathbb{E}[-J_{\tau_{k}}(x_{k+1}, \pi_{\tau_{k}}^{\star}(x_{k+1})) + J_{\tau_{k}}(x_{k}, \pi_{\tau_{k}}^{\star}(x_{k}))], \end{split}$$

where in the last inequality we have combined the terms $\frac{3B_D^2L_V^2\zeta_k^2\alpha_k}{4w_k^2}$ and $\frac{B_D^2L_\Phi\zeta_k^2}{w_k^2\tau_k}$ under the step size conditions $\tau_k \leq 1$ and $\alpha_k \leq \frac{4L_\Phi}{3L_V^2}$.

Recall the definition of ε_k^{θ} in (36). We can re-arrange the terms in the inequality above and obtain

$$\begin{split} & \mathbb{E}[\varepsilon_{k+1}^{\theta,\mathcal{L}} - \varepsilon_{k}^{\theta,\mathcal{L}}] \\ & = \mathbb{E}[-J_{\tau_{k+1}}(x_{k+1}, \pi_{\theta_{k+1}}) + J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})] \\ & + \mathbb{E}[J_{\tau_{k}}(x_{k+1}, \pi_{\tau_{k}}^{\star}(x_{k+1})) - J_{\tau_{k}}(x_{k}, \pi_{\tau_{k}}^{\star}(x_{k}))] \\ & + \mathbb{E}[J_{\tau_{k+1}}(x_{k+1}, \pi_{\tau_{k}}^{\star}(x_{k+1})) - J_{\tau_{k}}(x_{k+1}, \pi_{\tau_{k}}^{\star}(x_{k+1}))] \end{split}$$

$$\leq -\frac{\alpha_{k}}{8} \mathbb{E}[\|\nabla_{\theta} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}] + 2L_{F}^{2} \alpha_{k} \mathbb{E}[\varepsilon_{k}^{V}] + \frac{32L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} \tau_{k}^{2}} \mathbb{E}[\varepsilon_{k}^{x}] \\
- \frac{C_{L}^{2} \alpha_{k} \tau_{k}^{2}}{64} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2} L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})\|^{2}] \\
+ \frac{2B_{D}^{2} L_{\Phi} \zeta_{k}^{2}}{w_{L}^{2} \tau_{k}} + \frac{2B_{D} B_{F} L_{V} \zeta_{k} \alpha_{k}}{w_{k}} + \frac{B_{F}^{2} L_{V} \alpha_{k}^{2}}{2} + \frac{16 \log |\mathcal{A}| \tau_{k}}{3(1 - \gamma)(k + 1)},$$

where the bound on $J_{\tau_{k+1}}(x_{k+1}, \pi_{\tau_k}^{\star}(x_{k+1})) - J_{\tau_k}(x_{k+1}, \pi_{\tau_k}^{\star}(x_{k+1})) \leq \frac{8 \log |\mathcal{A}| \tau_k}{3(1-\gamma)(k+1)}$ can be obtained in a manner similar to (75).

C.3 Proof of Proposition 4

The proof depends on an intermediate result that bounds an important cross term. We state it in the lemma below and defer its proof to Section D.14.

Lemma 14 Under the assumptions and step sizes of Proposition 4, we have for all $k \ge 0$

$$\begin{split} & \mathbb{E}[\langle \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}(x_{k},\pi_{\theta_{w_{k},\tau_{k}}^{\star}}(x_{k})),x_{k+1} - x_{k} \rangle] \\ & \leq \frac{C_{L}^{2} \alpha_{k} \tau_{k}^{2}}{64} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ & \quad + \frac{64 L_{D}^{2} L_{L}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{32 L_{L}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E}[\varepsilon_{k}^{x}]. \end{split}$$

We now proceed to the proof of the proposition. We define the re-weighted functions

$$\mathcal{L}_{w,\tau}^{\text{reweight}}(x,\pi) \triangleq w \mathcal{L}_{w,\tau}(x,\pi) = w f(x,\pi) + (J_{\tau}(x,\pi_{\tau}^{\star}(x)) - J_{\tau}(x,\pi)), \quad \Phi_{w,\tau}^{\text{reweight}}(x) \triangleq w \Phi_{w,\tau}(x).$$

We consider the following decomposition and bound each term on the right hand side individually.

$$\mathcal{L}_{w_{k+1},\tau_{k+1}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}}) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})
= \left(\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}}) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k}^{\mathcal{L}}})\right) + \left(\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k}^{\mathcal{L}}}) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\right)
+ \left(\mathcal{L}_{w_{k+1},\tau_{k+1}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}}) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}})\right).$$
(76)

Bound the First Term of (76). As $\mathcal{L}_{w,\tau}^{\mathrm{reweight}}$ has L_L -Lipschitz gradients with respect to θ (shown in Lemma 9) under the condition $w, \tau \leq 1$,

$$\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k+1}^{\mathcal{L}}}) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}})
\leq \langle \nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}), \theta_{k+1}^{\mathcal{L}} - \theta_{k}^{\mathcal{L}} \rangle + \frac{L_{L}}{2} \|\theta_{k+1}^{\mathcal{L}} - \theta_{k}^{\mathcal{L}}\|^{2}
= \alpha_{k} \langle \nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}), F_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}, \bar{s}_{k}, \bar{a}_{k}, \bar{s}'_{k}, \xi_{k}) \rangle
+ \frac{L_{L}\alpha_{k}^{2}}{2} \|F_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}, \bar{s}_{k}, \bar{a}_{k}, \bar{s}'_{k}, \xi_{k})\|^{2}
= \alpha_{k} \langle \nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}), F_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}, \bar{s}_{k}, \bar{a}_{k}, \bar{s}'_{k}, \xi_{k}) - \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}) \rangle
+ \alpha_{k} \langle \nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}), \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}) - \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, V_{\tau_{k}}^{\tau_{k}}) \rangle
- \alpha_{k} \langle \nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}), \nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) \rangle
+ \frac{L_{L}\alpha_{k}^{2}}{2} \|F_{w_{k},\tau_{k}}(\theta_{k}, \omega_{k}, \hat{\mu}_{k}, \hat{V}_{f,k}, s_{k}, a_{k}, b_{k}, s'_{k}, \xi_{k})\|^{2},$$
(77)

where the final equation follows from

$$\nabla_{\theta} \mathcal{L}_{w_k, \tau_k}^{\text{reweight}}(x_k, \pi_{\theta_k^{\mathcal{L}}}) = w_k \nabla_{\theta} \mathcal{L}_{w_k, \tau_k}(x_k, \pi_{\theta_k^{\mathcal{L}}}) = -\bar{F}_{w_k, \tau_k}(x_k, \theta_k^{\mathcal{L}}, V_{\tau_k}^{x_k, \pi_{\theta_k^{\mathcal{L}}}}).$$

To bound the first term of (77),

$$\alpha_{k}\mathbb{E}[\langle \nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}), F_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}, \bar{s}_{k}, \bar{a}_{k}, \bar{s}_{k}', \xi_{k}) - \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}) \rangle]$$

$$= \alpha_{k}\mathbb{E}[\langle \nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}), \mathbb{E}[F_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}, \bar{s}_{k}, \bar{a}_{k}, \bar{s}_{k}', \xi_{k}) \mid \mathcal{F}_{k-1}] - \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}) \rangle]$$

$$- \alpha_{k}\mathbb{E}[\langle \nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}), F_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}, \bar{s}_{k}, \bar{a}_{k}, \bar{s}_{k}', \xi_{k}) - \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}) \rangle]$$

$$= -\alpha_{k}\mathbb{E}[\langle \nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}), F_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}, \bar{s}_{k}, \bar{a}_{k}, \bar{s}_{k}', \xi_{k}) - \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}) \rangle]$$

$$\leq \alpha_{k} \cdot L_{L}\mathbb{E}[\|x_{k+1} - x_{k}\|] \cdot 2B_{F}$$

$$\leq 2B_{F}L_{L}\alpha_{k} \cdot \frac{B_{D}\zeta_{k}}{w_{k}}$$

$$= \frac{2B_{D}B_{F}L_{L}\zeta_{k}\alpha_{k}}{w_{k}}, \qquad (78)$$

where the first inequality is due to (45) of Lemma 9, and the second inequality follows from Lemma 5. For the second term of (77),

$$\alpha_{k} \langle \nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}), \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}) - \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}}) \rangle \\
\leq \frac{\alpha_{k}}{8} \|\nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}})\|^{2} + 2\alpha_{k} \|\bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}) - \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}}) \|^{2} \\
\leq \frac{\alpha_{k}}{4} \|\nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}})\|^{2} + \frac{\alpha_{k}}{4} \|\nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}) \|^{2} \\
+ 2\alpha_{k} \|\bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, \hat{V}_{k}^{\mathcal{L}}) - \bar{F}_{w_{k},\tau_{k}}(x_{k}, \theta_{k}^{\mathcal{L}}, V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}}) \|^{2} \\
\leq \frac{\alpha_{k}}{4} \|\nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) \|^{2} + \frac{L_{L}^{2} \alpha_{k}}{4} \left(\|x_{k+1} - x_{k}\|^{2} \right) + 2L_{F}^{2} \alpha_{k} \|V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}} - \hat{V}_{k}^{\mathcal{L}} \|^{2} \\
\leq \frac{\alpha_{k}}{4} \|\nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) \|^{2} + \frac{L_{L}^{2} \alpha_{k}}{4} \cdot \left(\frac{B_{D} \zeta_{k}}{w_{k}} \right)^{2} + 2L_{F}^{2} \alpha_{k} \varepsilon_{k}^{V,\mathcal{L}} \\
\leq \frac{\alpha_{k}}{4} \|\nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) \|^{2} + \frac{B_{D}^{2} L_{L}^{2} \zeta_{k}^{2} \alpha_{k}}{4w_{k}^{2}} + 2L_{F}^{2} \alpha_{k} \varepsilon_{k}^{V,\mathcal{L}}, \tag{79}$$

where the third inequality is again a result of the Lipschitz continuity of $\nabla_{\theta} \mathcal{L}_{w_k, \tau_k}^{\text{reweight}}$ from Lemma 3 and the Lipschitz continuity of \bar{F}_{w_k, τ_k} .

For the third term of (77),

$$-\alpha_{k}\langle\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k}^{\mathcal{L}}}),\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\rangle$$

$$=-\alpha_{k}\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\|^{2} + \alpha_{k}\langle\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k}^{\mathcal{L}}}),\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\rangle$$

$$\leq -\frac{\alpha_{k}}{2}\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\|^{2} + \frac{L_{L}^{2}\alpha_{k}}{2}(\|x_{k+1}-x_{k}\|^{2})$$

$$\leq -\frac{\alpha_{k}}{2}\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\|^{2} + \frac{L_{L}^{2}\alpha_{k}}{2}(\|x_{k+1}-x_{k}\|^{2})$$

$$\leq -\frac{\alpha_{k}}{2}\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\|^{2} + \frac{L_{L}^{2}\alpha_{k}}{2}\cdot\left(\frac{B_{D}\zeta_{k}}{w_{k}}\right)^{2}$$

$$\leq -\frac{\alpha_{k}}{2}\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\|^{2} + \frac{B_{L}^{2}L_{L}^{2}\zeta_{k}^{2}\alpha_{k}}{2},$$

$$\leq -\frac{\alpha_{k}}{2}\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\|^{2} + \frac{B_{L}^{2}L_{L}^{2}\zeta_$$

where the second inequality again follows from the L_L -smoothness of $\mathcal{L}_{w_k,\tau_k}^{\text{reweight}}$ shown in (45) of Lemma 9.

Substituting (78)-(80) into (77),

$$\begin{split} & \mathbb{E}[\mathcal{L}_{w_k,\tau_k}^{\text{reweight}}(x_{k+1},\pi_{\theta_k^{\mathcal{L}}}) - \mathcal{L}_{w_k,\tau_k}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}})] \\ & \leq \frac{2B_DB_FL_L\zeta_k\alpha_k}{w_k} + \frac{\alpha_k}{4}\mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{w_k,\tau_k}^{\text{reweight}}(x_k,\pi_{\theta_k^{\mathcal{L}}})\|^2] + \frac{B_D^2L_L^2\zeta_k^2\alpha_k}{4w_k^2} + 2L_F^2\alpha_k\mathbb{E}[\varepsilon_k^{V,\mathcal{L}}] \\ & - \frac{\alpha_k}{2}\mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{w_k,\tau_k}^{\text{reweight}}(x_k,\pi_{\theta_k^{\mathcal{L}}})\|^2] + \frac{B_D^2L_L^2\zeta_k^2\alpha_k}{2w_t^2} + \frac{L_L\alpha_k^2}{2}\mathbb{E}[\|F_{w_k,\tau_k}(\theta_k,\omega_k,\hat{\mu}_k,\hat{V}_{f,k},s_k,a_k,b_k,s_k',\xi_k)\|^2] \end{split}$$

$$\leq -\frac{\alpha_{k}}{4} \mathbb{E}[\|\nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}})\|^{2}] + 2L_{F}^{2} \alpha_{k} \mathbb{E}[\varepsilon_{k}^{V,\mathcal{L}}] + \frac{3B_{D}^{2} L_{L}^{2} \zeta_{k}^{2} \alpha_{k}}{4w_{k}^{2}} + \frac{2B_{D} B_{F} L_{L} \zeta_{k} \alpha_{k}}{w_{k}} + \frac{B_{F}^{2} L_{L} \alpha_{k}^{2}}{2} \\
\leq -\frac{\alpha_{k}}{8} \mathbb{E}[\|\nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}})\|^{2}] + 2L_{F}^{2} \alpha_{k} \mathbb{E}[\varepsilon_{k}^{V,\mathcal{L}}] + \frac{3B_{D}^{2} L_{L}^{2} \zeta_{k}^{2} \alpha_{k}}{4w_{k}^{2}} + \frac{2B_{D} B_{F} L_{L} \zeta_{k} \alpha_{k}}{w_{k}} + \frac{B_{F}^{2} L_{L} \alpha_{k}^{2}}{2} \\
-\frac{C_{L}^{2} \alpha_{k} \tau_{k}^{2}}{32} \mathbb{E}[\|\pi_{\theta_{k}^{\mathcal{L}}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}], \tag{81}$$

where the last inequality follows from (25), which states that

$$\|\nabla_{\theta} \mathcal{L}_{w_k, \tau_k}^{\text{reweight}}(x_k, \pi_{\theta_k^{\mathcal{L}}})\| \ge \frac{C_L \tau_k}{2} \|\pi_{\theta_k^{\mathcal{L}}} - \pi_{w_k, \tau_k}^{\star}(x_k)\|.$$

Bound the Second Term of (76). We use $\theta_{w,\tau}^\star(x)$ to denote a softmax parameter that encodes the policy $\pi_{w,\tau}^\star(x)$. We know from Lemma 9 that $\mathcal{L}_{w,\tau}^{\mathrm{reweight}}$ is $\frac{L_L}{\tau}$ -smooth with respect to x under $w,\tau \leq 1$,

$$\begin{split} &\mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_{k+1},\pi_{\theta_k^{\mathcal{L}}}) - \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) \\ &\leq \langle \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}), x_{k+1} - x_k \rangle + \frac{L_L}{2\tau_k} \|x_{k+1} - x_k\|^2 \\ &= \langle \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_{w_k,\tau_k}}(x_k)), x_{k+1} - x_k \rangle + \frac{L_L}{2\tau_k} \|x_{k+1} - x_k\|^2 \\ &+ \langle \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_{w_k,\tau_k}}(x_k)), x_{k+1} - x_k \rangle \\ &= \langle \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_{w_k,\tau_k}}(x_k)), x_{k+1} - x_k \rangle + \frac{L_L}{2\tau_k} \|x_{k+1} - x_k\|^2 \\ &+ \langle \nabla_x \Phi^{\text{reweight}}_{w_k,\tau_k}(x_k), x_{k+1} - x_k \rangle \\ &\leq \langle \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_{w_k,\tau_k}}(x_k)), x_{k+1} - x_k \rangle + \frac{L_L}{2\tau_k} \|x_{k+1} - x_k\|^2 \\ &+ \left(\Phi^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k) \right) + \frac{L_\Phi}{2\tau_k} \|x_{k+1} - x_k\|^2 \\ &\leq \langle \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_{w_k,\tau_k}}(x_k)), x_{k+1} - x_k \rangle \\ &+ \left(\mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_{w_k,\tau_k}}(x_k)), x_{k+1} - x_k \rangle \\ &+ \left(\mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_{w_k,\tau_k}}(x_k)), x_{k+1} - x_k \right) \\ &+ \left(\mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_{w_k,\tau_k}}(x_k)), x_{k+1} - x_k \right) \\ &+ \left(\mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_{w_k,\tau_k}}(x_k)), x_{k+1} - x_k \right) \\ &+ \left(\mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) \right) \\ &+ \left(\mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) - \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k,\pi_{\theta_k^{\mathcal{L}}}) \right) \\ &+ \mathcal{L}^{\text{reweight}}_{w_k,\tau_k}(x_k$$

where the last inequality is due to $L_L \leq L_{\Phi}$, and the second equation uses the relationship $\nabla_x \mathcal{L}^{\text{reweight}}_{w_k, \tau_k}(x_k, \pi_{\theta^\star_{w_k, \tau_k}(x_k)}) = \nabla_x \Phi^{\text{reweight}}_{w_k, \tau_k}(x_k)$, which is due to $\nabla_\theta \mathcal{L}_{w_k, \tau_k}(x_k, \pi_{\theta^\star_{w_k, \tau_k}(x_k)}) = 0$ by the first-order optimality condition. The second inequality is due to the fact that $\Phi_{w,\tau}$ is $\frac{L_{\Phi}}{w\tau}$ -smooth when $w, \tau \leq 1$ (established in Lemma 10) and that for an L-smooth function f, we have

$$-f(y) + f(x) \le \langle -\nabla f(x), y - x \rangle + \frac{L}{2} ||x - y||^2.$$

Taking the expectation and plugging in the result from Lemma 14,

$$\mathbb{E}[\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k}^{\mathcal{L}}}) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}})] \\
\leq \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\
+ \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{32L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}} \mathbb{E}[\varepsilon_{k}^{x}] \\
+ \mathbb{E}[\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k+1})) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k}))] + \frac{L_{\Phi}}{\tau_{k}} \mathbb{E}[\|x_{k+1} - x_{k}\|^{2}] \\
\leq \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\
+ \mathbb{E}[\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k+1})) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k}))] + \frac{32L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}} \mathbb{E}[\varepsilon_{k}^{x}] + \frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{w_{k}^{2}\tau_{k}^{2}}, \\
(82)$$

where the last inequality follows from $||x_{k+1} - x_k|| \le \frac{B_D \zeta_k}{w_k}$.

Bound the Third Term of (76). By the definition of $\mathcal{L}_{w,\tau}$ in (11),

$$\mathcal{L}_{w_{k+1},\tau_{k+1}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}}) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}}) \\
= (w_{k+1} - w_{k})f(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}}) \\
+ \left(J_{\tau_{k+1}}(x_{k+1},\pi_{\tau_{k+1}}^{\star}(x_{k+1})) - J_{\tau_{k}}(x_{k+1},\pi_{\tau_{k}}^{\star}(x_{k+1}))\right) - \left(J_{\tau_{k+1}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}}) - J_{\tau_{k}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}})\right) \\
\leq (w_{k+1} - w_{k})f(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}}) + \frac{\log|\mathcal{A}|(\tau_{k} - \tau_{k+1})}{1 - \gamma} + \frac{\tau_{k} - \tau_{k+1}}{1 - \gamma} \mathbb{E}_{s \sim d_{\rho}^{\tau_{k+1}^{\mathcal{L}}}}[E(\pi_{k+1}^{\mathcal{L}},s)] \\
\leq 0 + \frac{\log|\mathcal{A}|(\tau_{k} - \tau_{k+1})}{1 - \gamma} + \frac{\tau_{k} - \tau_{k+1}}{1 - \gamma} \mathbb{E}_{s \sim d_{\rho}^{\tau_{k+1}^{\mathcal{L}}}}[E(\pi_{k+1}^{\mathcal{L}},s)] \\
\leq \frac{16\log|\mathcal{A}|\tau_{k}}{3(1 - \gamma)(k + 1)}, \tag{83}$$

where the first inequality follows from Zeng et al. [2022a][Lemma 3], the second inequality is due to the fact that f is non-negative from Assumption 3, and the third inequality follows from Lemma 2 and the relationship $E(\pi,s) \leq \log |\mathcal{A}|$ for any policy π .

Collecting the bounds in (81)-(83) and plugging them into (76), we get

$$\begin{split} &\mathbb{E}[\mathcal{L}_{w_{k+1},\tau_{k+1}}^{\text{eweight}}(x_{k+1},\pi_{\theta_{k+1}^{\mathcal{L}}}) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})] \\ &\leq -\frac{\alpha_{k}}{8}\mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}})\|^{2}] + 2L_{F}^{2}\alpha_{k}\mathbb{E}[\varepsilon_{k}^{V,\mathcal{L}}] + \frac{3B_{D}^{2}L_{L}^{2}\zeta_{k}^{2}\alpha_{k}}{4w_{k}^{2}} + \frac{2B_{D}B_{F}L_{L}\zeta_{k}\alpha_{k}}{w_{k}} + \frac{B_{F}^{2}L_{L}\alpha_{k}^{2}}{2} \\ &- \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{32}\mathbb{E}[\|\pi_{\theta_{k}^{\mathcal{L}}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \mathbb{E}[\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k+1})) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k}))] + \frac{32L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}] + \frac{B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{w_{k}^{2}\tau_{k}} \\ &+ \frac{16\log|\mathcal{A}|\tau_{k}}{3(1-\gamma)(k+1)} \\ \leq -\frac{\alpha_{k}}{8}\mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}})\|^{2}] + 2L_{F}^{2}\alpha_{k}\mathbb{E}[\varepsilon_{k}^{V,\mathcal{L}}] + \frac{32L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}] \\ &- \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{8}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{2B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{64} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{2B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{64} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{2B_{D}^{2}L_{\Phi}\zeta_{k}^{2}}{64} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{2B_{D}^{2}L_{L}^{2}\zeta_{k}^{$$

where in the last inequality we have combined the terms $\frac{3B_D^2L_L^2\zeta_k^2\alpha_k}{4w_k^2}$ and $\frac{B_D^2L_\Phi\zeta_k^2}{w_k^2\tau_k}$ under the step size conditions $\tau_k \leq 1$ and $\alpha_k \leq \frac{4L_\Phi}{3L^2}$.

Recall the definition of $\varepsilon_k^{\theta,\mathcal{L}}$ in (36). We can re-arrange the terms in the inequality above and obtain

$$\begin{split} & \mathbb{E}[\varepsilon_{k+1}^{\theta,\mathcal{L}} - \varepsilon_{k}^{\theta,\mathcal{L}}] \\ &= \mathbb{E}[\mathcal{L}_{w_{k+1},\tau_{k+1}}^{\text{reweight}}(x_{k+1}, \pi_{\theta_{k+1}^{\mathcal{L}}}) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}})] \\ &- \mathbb{E}[\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k+1})) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k}))] \\ &- \mathbb{E}[\mathcal{L}_{w_{k+1},\tau_{k+1}}^{\text{reweight}}(x_{k+1}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k+1})) - \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k+1}, \pi_{w_{k},\tau_{k}}^{\star}(x_{k+1}))] \end{split}$$

$$\leq -\frac{\alpha_{k}}{8} \mathbb{E}[\|\nabla_{\theta} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}})\|^{2}] + 2L_{F}^{2} \alpha_{k} \mathbb{E}[\varepsilon_{k}^{V,\mathcal{L}}] + \frac{32L_{L}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} \tau_{k}^{2}} \mathbb{E}[\varepsilon_{k}^{x}] \\
- \frac{C_{L}^{2} \alpha_{k} \tau_{k}^{2}}{64} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{64L_{D}^{2} L_{L}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\
+ \frac{2B_{D}^{2} L_{\Phi} \zeta_{k}^{2}}{w_{k}^{2} \tau_{k}} + \frac{2B_{D} B_{F} L_{L} \zeta_{k} \alpha_{k}}{w_{k}} + \frac{B_{F}^{2} L_{L} \alpha_{k}^{2}}{2} + \frac{32 \log |\mathcal{A}| \tau_{k}}{3(1 - \gamma)(k + 1)},$$

where the bound on $\mathcal{L}^{\mathrm{reweight}}_{w_{k+1},\tau_{k+1}}(x_{k+1},\pi^{\star}_{w_k,\tau_k}(x_{k+1})) - \mathcal{L}^{\mathrm{reweight}}_{w_k,\tau_k}(x_{k+1},\pi^{\star}_{w_k,\tau_k}(x_{k+1})) \leq \frac{16\log|\mathcal{A}|\tau_k}{3(1-\gamma)(k+1)}$ can be obtained in a manner similar to (83).

C.4 Proof of Proposition 5

We first establish the convergence of ε_k^V . To this end, we introduce the following technical lemma, which bounds an important cross term.

Lemma 15 Under the assumptions and step sizes of Proposition 5, we have for all $k \ge 0$

$$\begin{split} & \mathbb{E}[\langle \hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} + \beta_{k} \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}), V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} - V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}} \rangle] \\ & \leq \frac{(1 - \gamma)\beta_{k}}{2} \mathbb{E}[\|\hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} + \beta_{k} \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\|^{2}] + \frac{6L_{V}^{2} \zeta_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\varepsilon_{k}^{x}] \\ & + \frac{12L_{V}^{2} L_{D}^{2} \zeta_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{12L_{V}^{2} L_{D}^{2} \zeta_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})\|^{2}] \\ & + \frac{6L_{V}^{2} \alpha_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\|\nabla_{\theta} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}] + \frac{6L_{V}^{2} L_{F}^{2} \alpha_{k}^{2}}{(1 - \gamma)\beta_{k}} \mathbb{E}[\varepsilon_{k}^{V}] + \frac{6B_{F}^{2} L_{V} \tau_{0} \alpha_{k}^{2}}{\alpha_{0}} + \frac{32L_{V}^{2} \tau_{k}^{2}}{3(1 - \gamma)\beta_{k}(k + 1)^{2}}. \end{split}$$

We defer the proof of the lemma to Appendix D.15.

By the update rule in (51), we have

$$\varepsilon_{k+1}^{V} = \|\hat{V}_{k+1} - V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}}\|^{2} \\
= \|\Pi_{B_{V}} \left(\hat{V}_{k} + \beta_{k} G_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}')\right) - V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}}\|^{2} \\
\leq \|\hat{V}_{k} + \beta_{k} G_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}}\|^{2} \\
= \|\hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} + \beta_{k} \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) + \beta_{k} \left(G_{f}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\right) \\
+ V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} - V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}}\|^{2} \\
\leq \|\hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} + \beta_{k} \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\|^{2} \\
+ 2\beta_{k}^{2} \|G_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\|^{2} \\
+ 2 \|V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} - V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}}\|^{2} \\
+ \beta_{k} \langle \hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} + \beta_{k} \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}), G_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) \rangle \\
+ \langle \hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} + \beta_{k} \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}), V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} - V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}} \rangle, \tag{84}$$

where the first inequality follows from the fact that the projection to a convex set is non-expansive. To bound the first term of (84),

$$\begin{split} & \left\| \hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} + \beta_{k} \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) \right\|^{2} \\ & = \left\| \hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} \right\|^{2} + \beta_{k}^{2} \left\| \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) \right\|^{2} + 2\beta_{k} \langle \hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}, \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) \rangle \\ & = \left\| \hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} \right\|^{2} + \beta_{k}^{2} \left\| \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) - \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}) \right\|^{2} \end{split}$$

$$+ 2\beta_{k}\langle\hat{V}_{k} - V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}}, \bar{G}_{\tau_{k}}(x_{k},\theta_{k},\hat{V}_{k}) - \bar{G}_{\tau_{k}}(x_{k},\theta_{k},V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}})\rangle$$

$$= \|\hat{V}_{k} - V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}}\|^{2} + \beta_{k}^{2}\|\bar{G}_{\tau_{k}}(x_{k},\theta_{k},\hat{V}_{k}) - \bar{G}_{\tau_{k}}(x_{k},\theta_{k},V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}})\|^{2}$$

$$+ 2\beta_{k}\left(\hat{V}_{k} - V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}}\right)^{\top} (\gamma P^{\pi_{\theta_{k}}} - I)\left(\hat{V}_{k} - V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}}\right)$$

$$\leq \|\hat{V}_{k} - V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}}\|^{2} + L_{G}^{2}\beta_{k}^{2}\|\hat{V}_{k} - V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}}\|^{2} + 2(\gamma - 1)\beta_{k}\|\hat{V}_{k} - V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}}\|^{2}$$

$$\leq (1 - (1 - \gamma)\beta_{k})\varepsilon_{k}^{V},$$

$$(85)$$

where the second equation uses the relationship $\bar{G}_{\tau}(x,\theta,V_{\tau}^{x,\pi_{\theta}})=0$ for any τ,x,θ , and the last inequality follows from the step size condition $\beta_k \leq \frac{1-\gamma}{L_G^2}$.

For the third term of (84),

$$\begin{split} &2\|V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}}-V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}}\|^{2}\\ &\leq 4\|V_{\tau_{k}}^{x_{k+1},\pi_{\theta_{k+1}}}-V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}}\|^{2}+4\|V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}}-V_{\tau_{k}}^{x_{k+1},\pi_{\theta_{k+1}}}\|^{2}\\ &\leq 4|\mathcal{S}|\|V_{\tau_{k}}^{x_{k+1},\pi_{\theta_{k+1}}}-V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}}\|_{\infty}^{2}+8L_{V}^{2}\Big(\|x_{k}-x_{k+1}\|^{2}+\|\pi_{\theta_{k+1}}-\pi_{\theta_{k}}\|^{2}\Big)\\ &\leq 4|\mathcal{S}|\cdot\Big(\frac{\tau_{k}-\tau_{k+1}}{1-\gamma}\log|\mathcal{A}|\Big)^{2}+8L_{V}^{2}\Big(\frac{B_{D}^{2}\zeta_{k}^{2}}{w_{k}^{2}}+B_{F}^{2}\alpha_{k}^{2}\Big), \end{split}$$

where the second inequality follows from the Lipschitz continuity of the value function established in Lemma 3, and the third inequality applies Zeng et al. [2022a][Lemma 3]. Plugging in the bound on $\tau_k - \tau_{k+1}$ from Lemma 2, we get

$$2\|V_{\tau_{k}}^{x_{k},\pi_{\theta_{k}}} - V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}}\|^{2} \leq 4|\mathcal{S}| \cdot \left(\frac{8\tau_{k}}{3(1-\gamma)(k+1)}\log|\mathcal{A}|\right)^{2} + 8L_{V}^{2}\left(\frac{B_{D}^{2}\zeta_{k}^{2}}{w_{k}^{2}} + B_{F}^{2}\alpha_{k}^{2}\right)$$

$$\leq \frac{256|\mathcal{S}|\log^{2}|\mathcal{A}|\tau_{k}^{2}}{3(1-\gamma)^{2}(k+1)^{2}} + 16B_{F}^{2}L_{V}^{2}\alpha_{k}^{2}, \tag{86}$$

where we use the step size condition $\zeta_k \leq \frac{B_F \alpha_k w_k}{B_D}$.

For the fourth term of (84), we have in expectation

$$\mathbb{E}[\langle \hat{V}_k - V_{\tau_k}^{x_k, \pi_{\theta_k}} + \beta_k \bar{G}_{\tau_k}(x_k, \theta_k, \hat{V}_k), G_{\tau_k}(x_k, \theta_k, \hat{V}_k, s_k, a_k, s'_k) - \bar{G}_{\tau_k}(x_k, \theta_k, \hat{V}_k) \rangle]$$

$$= \mathbb{E}[\langle \hat{V}_k - V_{\tau_k}^{x_k, \pi_{\theta_k}} + \beta_k \bar{G}_{\tau_k}(x_k, \theta_k, \hat{V}_k), \mathbb{E}[G_{\tau_k}(x_k, \theta_k, \hat{V}_k, s_k, a_k, s'_k) - \bar{G}_{\tau_k}(x_k, \theta_k, \hat{V}_k) \mid \mathcal{F}_{k-1}] \rangle]$$

$$= 0.$$
(87)

Collecting the bounds from (85)-(87) and Lemma 15,

$$\begin{split} &\mathbb{E}[\varepsilon_{k+1}^{V}] \\ &= (1 - (1 - \gamma)\beta_{k})\mathbb{E}[\varepsilon_{k}^{V}] + 2\beta_{k}^{2}\mathbb{E}[\|G_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - \bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\|^{2}] \\ &+ \frac{256|\mathcal{S}|\log^{2}|\mathcal{A}|\tau_{k}^{2}}{3(1 - \gamma)^{2}(k + 1)^{2}} + 16B_{F}^{2}L_{V}^{2}\alpha_{k}^{2} \\ &+ \frac{(1 - \gamma)\beta_{k}}{2}\mathbb{E}[\|\hat{V}_{k} - V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}} + \beta_{k}\bar{G}_{\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\|^{2}] + \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}}\mathbb{E}[\varepsilon_{k}^{x}] \\ &+ \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{6L_{V}^{2}\alpha_{k}^{2}}{(1 - \gamma)\beta_{k}}\mathbb{E}[\|\nabla_{\theta}J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}})\|^{2}] + \frac{6L_{V}^{2}L_{F}^{2}\alpha_{k}^{2}}{(1 - \gamma)\beta_{k}}\mathbb{E}[\varepsilon_{k}^{V}] + \frac{6B_{F}^{2}L_{V}\tau_{0}\alpha_{k}^{2}}{\alpha_{0}} + \frac{32L_{V}^{2}\tau_{k}^{2}}{3(1 - \gamma)\beta_{k}(k + 1)^{2}} \\ &\leq \left(1 - (1 - \gamma)\beta_{k}\right)\mathbb{E}[\varepsilon_{k}^{V}] + \left(1 - (1 - \gamma)\beta_{k}\right)\frac{(1 - \gamma)\beta_{k}}{2}\mathbb{E}[\varepsilon_{k}^{V}] + \frac{6L_{V}^{2}L_{F}^{2}\alpha_{k}^{2}}{(1 - \gamma)\beta_{k}}\mathbb{E}[\varepsilon_{k}^{V}] + 8B_{G}\beta_{k}^{2} \\ &+ \frac{256|\mathcal{S}|\log^{2}|\mathcal{A}|\tau_{k}^{2}}{3(1 - \gamma)^{2}(k + 1)^{2}} + 16B_{F}^{2}L_{V}^{2}\alpha_{k}^{2} + \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1 - \gamma)\beta_{k}}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})\|^{2}] \end{aligned}$$

$$\begin{split} & + \frac{6L_{V}^{2}\alpha_{k}^{2}}{(1-\gamma)\beta_{k}}\mathbb{E}[\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}] + \frac{6B_{F}^{2}L_{V}\tau_{0}\alpha_{k}^{2}}{\alpha_{0}} + \frac{32L_{V}^{2}\tau_{k}^{2}}{3(1-\gamma)\beta_{k}(k+1)^{2}} \\ & \leq \left(1 - \frac{(1-\gamma)\beta_{k}}{4}\right)\mathbb{E}[\varepsilon_{k}^{V}] + \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] + \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ & + \frac{6L_{V}^{2}\alpha_{k}^{2}}{(1-\gamma)\beta_{k}}\mathbb{E}[\|\nabla_{\theta}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}})\|^{2}] + \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}}\mathbb{E}[\varepsilon_{k}^{x}] \\ & + \frac{22B_{F}^{2}L_{V}\tau_{0}\alpha_{k}^{2}}{\alpha_{0}} + \frac{64L_{V}^{2}\tau_{k}^{2}}{3(1-\gamma)\beta_{k}(k+1)^{2}} + 8B_{G}\beta_{k}^{2}, \end{split}$$

where the second inequality simplifies and combines terms based on the step size conditions $\frac{\alpha_k}{\beta_k} \leq \frac{1-\gamma}{2\sqrt{6}L_VL_F}$, $\beta_k \leq \frac{(1-\gamma)L_V^2}{8|\mathcal{S}|\log^2|\mathcal{A}|}$, and $\frac{\alpha_0}{\tau_0} \leq \frac{1}{L_V}$.

The bound on $\mathbb{E}[\varepsilon_{k+1}^{V,\mathcal{L}}]$ can be derived using an identical argument.

C.5 Proof of Proposition 2

The proof is almost identical to that of Proposition 1. We include the proof here for completeness.

From Lemma 10, we know that under a fixed regularization weight τ_0 , the objective Φ_{τ_0} has L_{Φ,τ_0} Lipschitz gradients, where we define

$$L_{\Phi,\tau_0} \triangleq \left(1 + \frac{2L_V}{C_L \tau_0}\right) \left(\frac{2L_f L_V}{C_L \tau_0} + \frac{2L_f L_V L_{V,2}}{\underline{\sigma} C_L \tau_0} + \frac{2L_f L_V^2 L_{V,2}}{\underline{\sigma}^2 C_L \tau_0} + \frac{2L_f L_V^2}{\underline{\sigma} C_L \tau}\right).$$

This implies

$$\Phi_{\tau_{0}}(x_{k+1}) - \Phi_{\tau_{0}}(x_{k})
\leq \langle \nabla_{x}\Phi_{\tau_{0}}(x_{k}), x_{k+1} - x_{k} \rangle + \frac{L_{\Phi,\tau_{0}}}{2} \|x_{k+1} - x_{k}\|^{2}
= -\zeta_{k} \langle \nabla_{x}\Phi_{\tau_{0}}(x_{k}), D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \rangle + \frac{L_{\Phi,\tau_{0}}\zeta_{k}^{2}}{2} \|D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \|^{2}
= -\zeta_{k} \|\nabla_{x}\Phi_{\tau_{0}}(x_{k})\|^{2} + \zeta_{k} \langle \nabla_{x}\Phi_{\tau_{0}}(x_{k}), \nabla_{x}\Phi_{\tau_{0}}(x_{k}) - D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \rangle
+ \frac{L_{\Phi,\tau_{0}}\zeta_{k}^{2}}{2} \|D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \|^{2}.$$
(88)

By the law of total expectation,

$$\begin{split} &\mathbb{E}[\Phi_{\tau_{0}}(x_{k+1}) - \Phi_{\tau_{0}}(x_{k})] \\ &\leq -\zeta_{k}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{0}}(x_{k})\|^{2}] + \frac{L_{\Phi,\tau_{0}}\zeta_{k}^{2}}{2}\mathbb{E}[\|D_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}},s_{k},a_{k},\bar{s}_{k},\bar{a}_{k},\xi_{k})\|^{2}] \\ &\quad + \zeta_{k}\mathbb{E}[\langle\nabla_{x}\Phi_{\tau_{0}}(x_{k}),\nabla_{x}\Phi_{\tau_{0}}(x_{k}) - \mathbb{E}[D_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}},s_{k},a_{k},\bar{s}_{k},\bar{a}_{k},\xi_{k}) \mid \mathcal{F}_{k-1}]\rangle] \\ &= -\zeta_{k}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{0}}(x_{k})\|^{2}] + \frac{L_{\Phi,\tau_{0}}\zeta_{k}^{2}}{2}\mathbb{E}[\|D_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}},s_{k},a_{k},\bar{s}_{k},\bar{a}_{k},\xi_{k})\|^{2}] \\ &\quad + \zeta_{k}\mathbb{E}[\langle\nabla_{x}\Phi_{\tau_{0}}(x_{k}),\nabla_{x}\Phi_{\tau_{0}}(x_{k}) - \bar{D}_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}})\rangle] \\ &\leq -\zeta_{k}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{0}}(x_{k})\|^{2}] + \frac{B_{D}^{2}L_{\Phi,\tau_{0}}\zeta_{k}^{2}}{2w_{k}^{2}} + \frac{\zeta_{k}}{2}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{0}}(x_{k})\|^{2}] + \frac{\zeta_{k}}{2}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{0}}(x_{k}) - \bar{D}_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}})\|^{2}] \\ &= -\frac{\zeta_{k}}{2}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{0}}(x_{k})\|^{2}] + \frac{B_{D}^{2}L_{\Phi,\tau_{0}}\zeta_{k}^{2}}{2w_{k}^{2}} \\ &\quad + \frac{\zeta_{k}}{2}\mathbb{E}\left[\|\left(\nabla_{x}\Phi_{\tau_{0}}(x_{k}) - \nabla_{x}\Phi_{w_{k},\tau_{0}}(x_{k})\right) + \left(\bar{D}_{w_{k}}(x_{k},\pi_{\tau_{0}}^{\star}(x_{k}),\pi_{w_{k},\tau_{0}}^{\star}(x_{k})) - \bar{D}_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}})\right)\|^{2}\right] \\ &\leq -\frac{\zeta_{k}}{2}\mathbb{E}[\|\nabla_{x}\Phi_{\tau_{0}}(x_{k})\|^{2}] + \zeta_{k}\mathbb{E}[\|\bar{D}_{w_{k}}(x_{k},\pi_{\tau_{0}}^{\star}(x_{k}),\pi_{w_{k},\tau_{0}}^{\star}(x_{k})) - \bar{D}_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}})\|^{2}] \end{split}$$

$$+ \zeta_k \mathbb{E}[\|\nabla_x \Phi_{\tau_0}(x_k) - \nabla_x \Phi_{w_k, \tau_0}(x_k)\|^2] + \frac{B_D^2 L_{\Phi, \tau_0} \zeta_k^2}{2w_b^2}, \tag{89}$$

where the second inequality applies Lemma 5 and the last equation follows from (32).

To bound the second term on the right hand side of (89), we apply Lemma 6

$$\|\bar{D}_{w_{k}}(x_{k}, \pi_{\tau_{0}}^{\star}(x_{k}), \pi_{w_{k}, \tau_{0}}^{\star}(x_{k})) - \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}})\|^{2}$$

$$\leq \frac{L_{D}^{2}}{w_{k}^{2}} \Big(\|\pi_{k} - \pi_{\tau_{0}}^{\star}(x_{k})\| + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{0}}^{\star}(x_{k})\| \Big)^{2}$$

$$\leq \frac{2L_{D}^{2}}{w_{k}^{2}} \|\pi_{k} - \pi_{\tau_{0}}^{\star}(x_{k})\|^{2} + \frac{2L_{D}^{2}}{w_{k}^{2}} \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{0}}^{\star}(x_{k})\|^{2}.$$

$$(90)$$

For the third term of (89), we have from Lemma 11

$$\|\nabla_x \Phi_{\tau_0}(x_k) - \nabla_x \Phi_{w_k, \tau_0}(x_k)\|^2 \le \left(\frac{4L_f L_V w_k}{C_L \underline{\sigma} \tau_0} (L_f + \frac{2L_f L_{V,2}}{C_L \tau_0})\right)^2 = C_{2, \tau_0} w_k^2, \tag{91}$$

where we define $C_{2,\tau_0}=\left(rac{4L_fL_V}{C_L\underline{\sigma} au_0}(L_f+rac{2L_fL_{V,2}}{C_L au_0})
ight)^2$.

Substituting (90) and (91) into (89),

$$\begin{split} & \mathbb{E}[\Phi_{\tau_0}(x_{k+1}) - \Phi_{\tau_0}(x_k)] \\ & \leq -\frac{\zeta_k}{2} \mathbb{E}[\|\nabla_x \Phi_{\tau_0}(x_k)\|^2] + \frac{2L_D^2 \zeta_k}{w_k^2} \mathbb{E}[\|\pi_k - \pi_{\tau_0}^{\star}(x_k)\|^2] + \frac{2L_D^2 \zeta_k}{w_k^2} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_0}^{\star}(x_k)\|^2] \\ & + C_{2,\tau_0} \zeta_k w_k^2 + \frac{B_D^2 L_{\Phi,\tau_0} \zeta_k^2}{2w_k^2}. \end{split}$$

D Proof of Supporting Results

D.1 Proof of Lemma 1

Uniqueness of $\pi^*(x)$. We first take the approach of proof by contradiction to show that $\pi^*(x)$ is unique for any x.

Given x, suppose that there exist two distinct optimal solutions of (5), which we denote by π_1, π_2 . From the definition of $\Pi^*(x)$, we know that π_1, π_2 satisfy

$$J(x,\pi_1) \ge J(x,\pi), \quad J(x,\pi_2) \ge J(x,\pi), \quad \forall \pi.$$
(92)

We construct another policy π' as follows, inspired by the proof of Theorem 1 in Zeng et al. [2023]. For all state s and action a, $\pi'(a \mid s)$ is expressed as

$$\pi'(a \mid s) = \frac{d_{\rho}^{\pi_1}(s)\pi_1(a \mid s) + d_{\rho}^{\pi_2}(s)\pi_2(a \mid s)}{d_{\rho}^{\pi_1}(s) + d_{\rho}^{\pi_2}(s)}.$$
(93)

Note that Assumption 1 guarantees $d_{\rho}^{\pi}(s) \geq (1 - \gamma)\rho_{\min} > 0$ for all s, ensuring π' is well-defined.

Our first step is to show that π' is also an optimal policy, i.e. $\pi' \in \Pi^*(x)$. To see this, we define a modified transition kernel \mathcal{P}_{γ} such that $\mathcal{P}_{\gamma}(s'\mid s,a) = (1-\gamma)\rho(s') + \gamma \mathcal{P}(s'\mid s,a)$ for any π , and P^{π}_{γ} such that $P^{\pi}_{\gamma}(s'\mid s) = (1-\gamma)\rho(s') + \gamma P^{\pi}(s'\mid s)$ for any π .

We also define a vector $B \in \mathbb{R}^{|\mathcal{S}|}$

$$B = \mathcal{P}_{\gamma}^{\pi'} \cdot \left(\frac{1}{2}d_{\rho}^{\pi_1} + \frac{1}{2}d_{\rho}^{\pi_2}\right).$$

We can express each entry of B in the following way

$$B(s') = \sum_{s,a} \mathcal{P}_{\gamma}(s' \mid s, a) \pi'(a \mid s) \left(\frac{1}{2} d_{\rho}^{\pi_1}(s) + \frac{1}{2} d_{\rho}^{\pi_2}(s) \right)$$

$$\begin{split} &= \sum_{s,a} \mathcal{P}_{\gamma}(s'\mid s,a) \frac{d_{\rho}^{\pi_{1}}(s)\pi_{1}(a\mid s) + d_{\rho}^{\pi_{2}}(s)\pi_{2}(a\mid s)}{d_{\rho}^{\pi_{1}}(s) + d_{\rho}^{\pi_{2}}(s)} \left(\frac{1}{2} d_{\rho}^{\pi_{1}}(s) + \frac{1}{2} d_{\rho}^{\pi_{2}}(s)\right) \\ &= \frac{1}{2} \sum_{s,a} \mathcal{P}_{\gamma}(s'\mid s,a) d_{\rho}^{\pi_{1}}(s)\pi_{1}(a\mid s) + \frac{1}{2} \sum_{s,a} \mathcal{P}(s'\mid s,a) d_{\rho}^{\pi_{2}}(s)\pi_{2}(a\mid s) \\ &= \frac{1}{2} \sum_{s,a} P_{\gamma}^{\pi_{1}}(s'\mid s) d_{\rho}^{\pi_{1}}(s) + \frac{1}{2} \sum_{s,a} P_{\gamma}^{\pi_{2}}(s'\mid s) d_{\rho}^{\pi_{2}}(s) \\ &= \frac{1}{2} d_{\rho}^{\pi_{1}}(s') + \frac{1}{2} d_{\rho}^{\pi_{2}}(s'), \end{split}$$

which leads to

$$P_{\gamma}^{\pi'} \cdot \left(\frac{1}{2}d_{\rho}^{\pi_1} + \frac{1}{2}d_{\rho}^{\pi_2}\right) = \frac{1}{2}d_{\rho}^{\pi_1} + \frac{1}{2}d_{\rho}^{\pi_2}.$$
 (94)

The Markov chain induced by P^π_γ is always ergodic, assumed in Assumption 1. Under ergodicity, it is known that d^π_ρ (properly normalized) is the unique eigenvector of P^π_γ associated with eigenvalue 1. Therefore, (94) implies that $\frac{1}{2}d^{\pi_1}_\rho+\frac{1}{2}d^{\pi_2}_\rho$ is the discounted visitation distribution induced by policy π' , i.e.

$$d_{\rho}^{\pi'} = \frac{1}{2} d_{\rho}^{\pi_1} + \frac{1}{2} d_{\rho}^{\pi_2}. \tag{95}$$

We use \hat{d}^{π}_{ρ} to denote the extend discounted visitation distribution over state and action such that

$$\hat{d}^{\pi}_{\rho}(s, a) = d^{\pi}_{\rho}(s)\pi(a \mid s).$$

We have from (95)

$$\begin{split} \hat{d}_{\rho}^{\pi'}(s,a) &= d_{\rho}^{\pi'}(s)\pi'(a\mid s) \\ &= \left(\frac{1}{2}d_{\rho}^{\pi_{1}}(s) + \frac{1}{2}d_{\rho}^{\pi_{2}}(s)\right)\frac{d_{\rho}^{\pi_{1}}(s)\pi_{1}(a\mid s) + d_{\rho}^{\pi_{2}}(s)\pi_{2}(a\mid s)}{d_{\rho}^{\pi_{1}}(s) + d_{\rho}^{\pi_{2}}\mu_{\pi_{2}}(s)} \\ &= \frac{1}{2}d_{\rho}^{\pi_{1}}(s)\pi_{1}(a\mid s) + \frac{1}{2}d_{\rho}^{\pi_{1}}(s)\pi_{2}(a\mid s) \\ &= \frac{1}{2}\hat{d}_{\rho}^{\pi_{1}}(s,a) + \frac{1}{2}\hat{d}_{\rho}^{\pi_{2}}(s,a). \end{split}$$

Note that the cumulative return J is linear in the space of extended discounted visitation distribution. We have

$$J(x,\pi') = \langle r_x, \hat{d}_{\rho}^{\pi'} \rangle = \frac{1}{2} \langle r_x, \frac{1}{2} \hat{d}_{\rho}^{\pi_1}(s,a) + \frac{1}{2} \hat{d}_{\rho}^{\pi_2}(s,a) \rangle = \frac{1}{2} J(x,\pi_1) + \frac{1}{2} J(x,\pi_2).$$

In view of (92), this implies that π' is an optimal policy, i.e. $\pi' \in \Pi^*(x)$.

Since π' is in the constraint set for the optimization problem in (5), we can create a contradiction that π_1, π_2 are the two distinct maximizers of (5) if π' has a larger weighted entropy. The entropy function E(x) is strictly concave for all state x for policies in the interior of the simplex (note that π_1, π_2 must in the interior of the simplex, as they cannot be the optimal solution of (5) otherwise). Recall from (93) that $\pi'(x)$ is a convex combination of $\pi_1(x)$ is $\pi_2(x)$, and by the property of strictly concave functions,

$$E(\pi',s) > \frac{d_{\rho}^{\pi_1}(s)E(\pi_1,s) + d_{\rho}^{\pi_2}(s)E(\pi_2,s)}{d_{\rho}^{\pi_1}(s) + d_{\rho}^{\pi_2}(s)}.$$

We denote $\mathcal{E}_{\pi} = \mathbb{E}_{s \sim d_{\alpha}^{\pi}}[E(\pi, s)].$

$$\begin{split} \mathcal{E}_{\pi'} &= \langle d_{\rho}^{\pi'}, E(\pi', \cdot) \rangle \\ &= \sum_{s} \left(\frac{1}{2} d_{\rho}^{\pi_1}(s) + \frac{1}{2} d_{\rho}^{\pi_2}(s) \right) E(\pi', s) \end{split}$$

$$\begin{split} &> \sum_{s} \left(\frac{1}{2} d_{\rho}^{\pi_{1}}(s) + \frac{1}{2} d_{\rho}^{\pi_{2}}(s)\right) \frac{d_{\rho}^{\pi_{1}}(s) E(\pi_{1}, s) + d_{\rho}^{\pi_{2}}(s) E(\pi_{2}, s)}{d_{\rho}^{\pi_{1}}(s) + d_{\rho}^{\pi_{2}}(s)} \\ &= \frac{1}{2} \langle d_{\rho}^{\pi_{1}}, E(\pi_{1}, s) \rangle + \frac{1}{2} \langle d_{\rho}^{\pi_{2}}, E(\pi_{2}, s) \rangle \\ &= \frac{1}{2} \mathcal{E}_{\pi_{1}} + \frac{1}{2} \mathcal{E}_{\pi_{2}}. \end{split}$$

This contradicts the condition that π_1, π_2 maximize the weighted entropy within the constraint set $\Pi^*(x)$ and concludes our proof on the uniqueness of $\pi^*(x)$.

Limit Point of $\{\pi_{\tau}^{\star}(x)\}_{\tau}$. We then show the limit point of $\{\pi_{\tau}^{\star}(x)\}_{\tau}$ is $\pi^{\star}(x)$ as $\tau \to 0$. As x is fixed here, we simply denote $\pi_{\tau} = \pi_{\tau}^{\star}(x)$. We define for simplicity $\bar{E}(\pi) \triangleq \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi}}[E(\pi,s)]$. By the Bolzano–Weierstrass theorem, as the sequence $\{\pi_{\tau}\}$ is bounded, it has a convergent subsequence. Let $\tau_n \to 0$ and $\pi_{\tau_n} \to \bar{\pi}$. We first need to show $\bar{\pi} \in \Pi^{\star}(x)$. By the definition of $\pi_{\tau}^{\star}(x)$,

$$J(x, \pi_{\tau}) + \tau \bar{E}(\pi_{\tau}) \ge J(x, \pi^{\star}(x)) + \tau \bar{E}(\pi^{\star}(x)), \tag{96}$$

leading to

$$J(x, \pi_{\tau}) \ge J(x, \pi^{\star}(x)) + \tau \Big(\bar{E}(\pi^{\star}(x)) - \bar{E}(\pi_{\tau})\Big).$$

As J is continuous, we take $n \to \infty$

$$\limsup J(x, \pi_{\tau}) \ge J(x, \pi^{\star}(x)),$$

implying $J(x, \bar{\pi})$. This means $\bar{\pi} \in \Pi^*$.

Then, to show $\bar{\pi}$ is the maximum entropy solution, we re-arrange the terms in (96)

$$\bar{E}(\pi_{\tau}) - \bar{E}(\pi^{\star}(x)) \ge \frac{J(x, \pi^{\star}(x)) - J(x, \pi_{\tau})}{\tau} \ge 0,$$

where the second inequality is due to the definition of $\pi^{\star}(x)$.

Taking the limit,

$$\lim \sup \bar{E}(\pi_{\tau}) \geq \bar{E}(\pi^{\star}(x)).$$

As the limit point $\bar{\pi}$ is in Π^* and we have $\pi^*(x) = \operatorname{argmax}_{\pi \in \Pi^*} \bar{E}(\pi)$, then it holds that $\limsup \bar{E}(\pi_\tau) = \bar{E}(\pi^*(x)) = \bar{E}(\bar{\pi})$. This allows us to conclude that $\pi_\tau^*(x) \to \pi^*(x)$ as $\tau \to 0$.

D.2 Proof of Lemma 2

We apply Zeng et al. [2022a][Lemma 7]

$$\tau_k - \tau_{k+1} = \frac{\tau_0}{(k+1)^{c_\tau}} - \frac{\tau_0}{(k+2)^{c_\tau}} \le \frac{8\tau_0}{3(k+1)^{c_\tau+1}} = \frac{8\tau_k}{3(k+1)}.$$

D.3 Proof of Lemma 3

We derive the inequalities on the value function $V^{x,\pi_{\theta}}_{\tau}$ and note that the ones on the cumulative return $J_{\tau}(x,\pi_{\theta})$ immediately follows as the cumulative return is simply an weighted average of the value function

$$J_{\tau}(x, \pi_{\theta}) = \mathbb{E}_{s \sim \rho}[V_{\tau}^{x, \pi_{\theta}}(s)].$$

Fixing x, we know that $V_{\tau}^{x,\pi}$ is the standard policy optimization objective (in a standard MDP) as a function of policy π . It is well-known that the following inequalities hold (for example, see Lemma B.5 of Zeng et al. [2021] and Lemma 5 of Zeng et al. [2022a])

$$||V^{x,\pi_{\theta}} - V^{x,\pi_{\theta'}}|| \le \frac{2}{(1-\gamma)^2} ||\theta - \theta'||,$$
 (97)

$$\|\nabla_{\theta} V^{x,\pi_{\theta}} - \nabla_{\theta} V^{x,\pi_{\theta'}}\| \le \frac{8}{(1-\gamma)^3} \|\theta - \theta'\|. \tag{98}$$

We define the shorthand notation

$$\mathcal{H}(\theta, s) \triangleq \mathbb{E}_{a_k \sim \pi(\cdot \mid s_k), s_{k+1} \sim \mathcal{P}(\cdot \mid s_k, a_k)} \left[\sum_{k=0}^{\infty} -\gamma^k \log \pi_{\theta}(a_k \mid s_k) \mid s_0 = s \right].$$

This implies $V_{\tau}^{x,\pi_{\theta}}(s) = V^{x,\pi_{\theta}}(s) + \tau \mathcal{H}(\theta,s)$. We also define the aggregate notation

$$\mathcal{H}(\theta) = [\mathcal{H}(\theta, s_1), \mathcal{H}(\theta, s_2), \cdots] \in \mathbb{R}^{|\mathcal{S}|}$$

Adapting Lemma 6 of Zeng et al. [2022a], we have for all $s \in \mathcal{S}$

$$\|\mathcal{H}(\theta, s) - \mathcal{H}(\theta', s)\| \le \frac{4 + 8\log|\mathcal{A}|}{(1 - \gamma)^3} \|\theta - \theta'\|,\tag{99}$$

$$\|\nabla_{\theta} \mathcal{H}(\theta, s) - \nabla_{\theta} \mathcal{H}(\theta', s)\| \le \frac{4 + 8\log|\mathcal{A}|}{(1 - \gamma)^3} \|\theta - \theta'\|. \tag{100}$$

We obviously have the following inequalities from (99) and (100)

$$\|\mathcal{H}(\theta) - \mathcal{H}(\theta')\| \le \frac{(4 + 8\log|\mathcal{A}|)\sqrt{|\mathcal{S}|}}{(1 - \gamma)^3} \|\theta - \theta'\|,\tag{101}$$

$$\|\nabla_{\theta} \mathcal{H}(\theta) - \nabla_{\theta} \mathcal{H}(\theta')\| \le \frac{(4 + 8\log|\mathcal{A}|)\sqrt{|\mathcal{S}|}}{(1 - \gamma)^3} \|\theta - \theta'\|. \tag{102}$$

Note that (100) also implies

$$\|\nabla_{\theta,\theta} \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta}}} [E(\pi_{\theta}, s)]\| \le \frac{4 + 8\log|\mathcal{A}|}{(1 - \gamma)^3}$$

hence leading to (42).

In addition, we have from (2)

$$|V_{\tau}^{x,\pi_{\theta}} - V_{\tau}^{x',\pi_{\theta}}| = |\mathbb{E}_{s' \sim d_{s}^{\pi}, a' \sim \pi(\cdot|s')}[r_{x}(s, a) - r_{x'}(s, a)]|$$

$$\leq \mathbb{E}_{s' \sim d_{s}^{\pi}, a' \sim \pi(\cdot|s')}[|r_{x}(s, a) - r_{x'}(s, a)|]$$

$$\leq \mathbb{E}_{s' \sim d_{s}^{\pi}, a' \sim \pi(\cdot|s')}[L_{r}||x - x'||]$$

$$\leq L_{r}||x - x'||, \qquad (103)$$

where the first inequality is a result of Jensen's inequality and the second inequality is due to Assumption 4.

We can express $\nabla_{\theta} V^{x,\pi_{\theta}}$ as follows [Agarwal et al., 2021]

$$\nabla_{\theta_{s',a'}} V^{x,\pi_{\theta}}(s) = \frac{1}{1-\gamma} d_s^{\pi_{\theta}}(s') \pi_{\theta}(a' \mid s') A^{x,\pi_{\theta}}(s',a').$$

This implies

$$\begin{split} |\nabla_{\theta_{s',a'}} V^{x,\pi_{\theta}}(s) - \nabla_{\theta_{s',a'}} V^{x',\pi_{\theta}}(s)| &= \frac{1}{1-\gamma} d_s^{\pi_{\theta}}(s') \pi_{\theta}(a' \mid s') |A^{x,\pi_{\theta}}(s',a') - A^{x',\pi_{\theta}}(s',a')| \\ &\leq \frac{1}{1-\gamma} |d_s^{\pi_{\theta}}(s')| |\pi_{\theta}(a' \mid s')| \cdot 2L_r ||x - x'|| \\ &\leq \frac{2}{1-\gamma} ||x - x'||, \end{split}$$

where the first inequality is due to the fact that the advantage function is $2L_r$ -Lipschitz with respect to x, since it is the difference of the value function and Q function, which are themselves L_r -Lipschitz with respect to x as can be seen from (103). Then, as the entropy regularizer is not a function of x, we have

$$\|\nabla_{\theta}V_{\tau}^{x,\pi_{\theta}} - \nabla_{\theta}V_{\tau}^{x,\pi_{\theta'}}\| = \|\nabla_{\theta}V^{x,\pi_{\theta}} - \nabla_{\theta}V^{x,\pi_{\theta'}}\|$$

$$\leq \frac{2L_r|\mathcal{S}||\mathcal{A}|}{1-\gamma} ||x-x'||. \tag{104}$$

Combining (97), (101), and (103), we have for all $\tau \leq 1$

$$||V_{\tau}^{x,\pi_{\theta}} - V_{\tau}^{x',\pi_{\theta'}}|| \leq ||V_{\tau}^{x,\pi_{\theta}} - V_{\tau}^{x',\pi_{\theta}}|| + ||V_{\tau}^{x',\pi_{\theta}} - V_{\tau}^{x',\pi_{\theta'}}||$$

$$\leq ||V_{\tau}^{x,\pi_{\theta}} - V_{\tau}^{x',\pi_{\theta}}|| + \tau ||\mathcal{H}(\theta) - \mathcal{H}(\theta')|| + ||V^{x',\pi_{\theta}} - V^{x',\pi_{\theta'}}||$$

$$\leq L_{r}||x - x'|| + \frac{(4 + 8\log|\mathcal{A}|)\sqrt{|\mathcal{S}|}}{(1 - \gamma)^{3}}||\theta - \theta'|| + \frac{2}{(1 - \gamma)^{2}}||\theta - \theta'||$$

$$\leq L_{r}||x - x'|| + \frac{(6 + 8\log|\mathcal{A}|)\sqrt{|\mathcal{S}|}}{(1 - \gamma)^{3}}||\theta - \theta'||.$$

This shows the Lipschitz continuity of the value function.

To show smoothness with respect to θ , we combine (98), (102), and (104)

$$\begin{split} \|\nabla_{\theta}V_{\tau}^{x,\pi_{\theta}} - \nabla_{\theta}V_{\tau}^{x',\pi_{\theta'}}\| &\leq \|\nabla_{\theta}V_{\tau}^{x,\pi_{\theta}} - \nabla_{\theta}V_{\tau}^{x',\pi_{\theta}}\| + \|\nabla_{\theta}V_{\tau}^{x',\pi_{\theta}} - \nabla_{\theta}V_{\tau}^{x',\pi_{\theta'}}\| \\ &\leq \|\nabla_{\theta}V_{\tau}^{x,\pi_{\theta}} - \nabla_{\theta}V_{\tau}^{x',\pi_{\theta}}\| + \tau \|\nabla_{\theta}\mathcal{H}(\theta) - \nabla_{\theta}\mathcal{H}(\theta')\| + \|\nabla_{\theta}V^{x',\pi_{\theta}} - \nabla_{\theta}V^{x',\pi_{\theta'}}\| \\ &\leq \frac{2L_{r}|\mathcal{S}||\mathcal{A}|}{1 - \gamma} \|x - x'\| + \frac{(4 + 8\log|\mathcal{A}|)\sqrt{|\mathcal{S}|}}{(1 - \gamma)^{3}} \|\theta - \theta'\| + \frac{8}{(1 - \gamma)^{3}} \|\theta - \theta'\| \\ &\leq \frac{2L_{r}|\mathcal{S}||\mathcal{A}|}{1 - \gamma} \|x - x'\| + \frac{(12 + 8\log|\mathcal{A}|)\sqrt{|\mathcal{S}|}}{(1 - \gamma)^{3}} \|\theta - \theta'\|. \end{split}$$

To show (39),

$$\begin{split} &\|\nabla_{x}J_{\tau}(x,\pi) - \nabla_{x}J_{\tau}(x',\pi')\| \\ &\leq \|\nabla_{x}J_{\tau}(x,\pi) - \nabla_{x}J_{\tau}(x,\pi')\| + \|\nabla_{x}J_{\tau}(x,\pi') - \nabla_{x}J_{\tau}(x',\pi')\| \\ &\leq \|\sum_{s,a}(d_{\rho}^{\pi}(s,a) - d_{\rho}^{\pi'}(s,a))\nabla_{x}r_{x}(s,a)\| + \mathbb{E}_{s\sim d_{\rho}^{\pi'},a\sim\pi'(\cdot|s)}[\|\nabla_{x}r_{x}(s,a) - \nabla_{x}r_{x'}(s,a)\|] \\ &\leq L_{r}\|d_{\rho}^{\pi} - d_{\rho}^{\pi'}\| + \mathbb{E}_{s\sim d_{\rho}^{\pi'},a\sim\pi'(\cdot|s)}[L_{r}\|x - x'\|] \\ &= \left\|(1-\gamma)\left((I-\gamma P^{\pi})^{-1} - (I-\gamma P^{\pi'})^{-1}\right)\rho\right\| + L_{r}\|x - x'\| \\ &\leq (1-\gamma) \cdot \frac{\gamma}{(1-\gamma)^{2}}\|\pi - \pi'\|\|\rho\| + L_{r}\|x - x'\| \\ &\leq L_{V}(\|\pi - \pi'\| + \|x - x'\|), \end{split}$$

where the fourth inequality is due to the fact that $(I - \gamma P^{\pi})$ is $\gamma/(1 - \gamma)^2$ Lipschitz in π .

Finally, we show $\nabla^2_{x,\theta}J_{\tau}(x,\pi_{\theta})$ and $\nabla^2_{\theta,\theta}J_{\tau}(x,\pi_{\theta})$ are Lipschitz. Since $\nabla_{\theta,\theta}\mathbb{E}_{s\sim d^{\pi_{\theta}}_{\rho}}[E(\pi_{\theta},s)]$ does not depend on x and can be shown to have Lipschitz Hessians by extending the argument in Mei et al. [2020][Lemma 14] (we skip showing the exact constant here), the problem reduces to showing that $\nabla^2_{x,\theta}J(x,\pi_{\theta})$ and $\nabla^2_{\theta,\theta}J(x,\pi_{\theta})$ are Lipschitz.

From (15), we have

$$\nabla_{x,\theta}^{2} J(x, \pi_{\theta}) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta}}, a \sim \pi_{\theta}(\cdot \mid s), s' \sim \mathcal{P}(\cdot \mid s, a)} \left[\left(\nabla_{x} r_{x}(s, a) + \gamma \nabla_{x} V_{\tau}^{x, \pi_{\theta}}(s') \right) \nabla_{\theta} \log \pi_{\theta}(a \mid s) \right]. \tag{105}$$

We define $\text{traj} = \{s_0, a_0, s_1, a_1, \dots\}$ and use $p(\pi, \text{traj})$ to denote the probability that the trajectory traj is generated under the policy π , i.e.

$$p(\pi, \operatorname{traj}) = \rho(s_0) \prod_{k=0}^{\infty} \pi(a_k \mid s_k) \mathcal{P}(s_{k+1} \mid s_k, a_k).$$

Adapting the result from Shen et al. [2019], we have

$$\nabla_{\theta,\theta}^2 J(x,\pi_{\theta}) = \mathbb{E}_{\text{traj} \sim p(\pi_{\theta},\cdot)} \Big[\nabla_{\theta} \phi(x,\pi,\text{traj}) \nabla_{\theta} p(\pi_{\theta},\text{traj})^{\top} + \nabla_{\theta,\theta}^2 \phi(x,\pi,\text{traj}) \Big], \quad (106)$$

where we define $\phi(x, \pi, \text{traj}) = \sum_{t=0}^{\infty} (\sum_{k=t}^{\infty} \gamma^k (r_x(s_k, a_k))^k) \log \pi(a_t \mid s_t)$.

As the right hand side expressions in (105) and (106) are the composition of Lipschitz functions, we know that $\nabla^2_{x,\theta}J(x,\pi_\theta)$ and $\nabla^2_{x,\theta}J(\theta,\pi_\theta)$ are Lipschitz.

D.4 Proof of Lemma 4

We can write the discounted visitation distribution as follows

$$d_{\rho}^{\pi} = (1 - \gamma)(I - \gamma P^{\pi})^{-1}\rho$$

which implies

$$\begin{split} \|d_{\rho}^{\pi} - d_{\rho}^{\pi'}\| &\leq (1 - \gamma) \|(I - \gamma P^{\pi})^{-1} \rho - (I - \gamma P^{\pi'})^{-1} \| \|\rho \| \\ &= (1 - \gamma) \|\rho\| \| (I - \gamma P^{\pi'})^{-1} \Big((I - \gamma P^{\pi'}) - (I - \gamma P^{\pi}) \Big) (I - \gamma P^{\pi})^{-1} \| \\ &\leq (1 - \gamma) \cdot 1 \cdot \| (I - \gamma P^{\pi})^{-1} \| \| (I - \gamma P^{\pi'})^{-1} \| \cdot \gamma \| P^{\pi} - P^{\pi'} \| \\ &\leq \frac{\gamma}{1 - \gamma} \|\pi - \pi' \|, \end{split}$$

where the third inequality follows from the standard result $||P^{\pi} - P^{\pi'}|| \leq \sqrt{|\mathcal{A}|} ||\pi - \pi'||$.

D.5 Proof of Lemma 5

By the definition of D_w in (26),

$$||D_{w}(x,\pi,\pi^{\mathcal{L}},s,a,\bar{s},\bar{a},\xi)|| = \left||\widetilde{\nabla}_{x}f(x,\pi^{\mathcal{L}},\xi) + \frac{1}{w}\left(\nabla_{x}r_{x}(s,a) + \nabla_{x}r_{x}(\bar{s},\bar{a})\right)\right||$$

$$\leq ||\widetilde{\nabla}_{x}f(x,\pi^{\mathcal{L}},\xi)|| + \frac{1}{w}||\nabla_{x}r_{x}(s,a) + \nabla_{x}r_{x}(\bar{s},\bar{a})||$$

$$\leq L_{f} + \frac{1}{w}(L_{r} + L_{r})$$

$$\leq \frac{3L_{r}}{w},$$

where the second inequality follows from Assumption 3 and the condition $\|\nabla_x r_x(s,a)\| \le L_r$, $\forall x, s, a$, which follows from Assumption 4, and the last inequality is due to $w \le \frac{L_r}{L_t}$.

To bound F, note that $\nabla_{\theta} \log \pi_{\theta}(a \mid s)$ has the following closed-form expression entry-wise

$$\frac{\partial \log \pi_{\theta}(a \mid s)}{\partial \theta_{s',a'}} = \mathbf{1} \left[s = s' \right] \left(\mathbf{1} \left[a = a' \right] - \pi_{\theta} \left(a' \mid s \right) \right). \tag{107}$$

This implies

$$\|\nabla_{\theta} \log \pi_{\theta}(a \mid s)\|_{2} < \|\nabla_{\theta} \log \pi_{\theta}(a \mid s)\|_{1} < 1 + 1 = 2. \tag{108}$$

By the definition of $F_{w,\tau}$ in (27),

$$||F_{w,\tau}(x,\theta,V,s,a,s',\xi)|| \le ||\nabla_{\theta}\log \pi_{\theta}(a \mid s)|| \Big| r_{x}(s,a) + \tau E(\pi_{\theta},s) + \gamma V(s') - V(s) \Big| + w ||\widetilde{\nabla}_{\theta} f(x,\pi_{\theta},\xi)||$$

$$\le 2(1 + \tau \log |\mathcal{A}| + \gamma B_{V} + B_{V}) + w L_{f}$$

$$\le 2(1 + \gamma)B_{V} + 2\log |\mathcal{A}| + 2 + L_{r},$$

where the second inequality applies (108), and the last inequality follows from $w \leq \frac{L_r}{L_f}$, $\tau \leq 1$.

Finally, by the definition of G in (28),

$$||G_{\tau}(x, \theta, V, s, a, s')|| \le ||e_{s}|| ||r_{x}(s, a) + \tau E(\pi_{\theta}, s) + \gamma V(s') - V(s)||$$

$$\le 1 \cdot (1 + \tau \log |\mathcal{A}| + \gamma B_{V} + B_{V})$$

$$\le (1 + \gamma)B_{V} + \log |\mathcal{A}| + 1.$$

42

D.6 Proof of Lemma 6

By the definition of \bar{D}_w in (29),

$$\begin{split} &\|\bar{D}_{w}(x_{1},\pi_{1},\pi_{1}^{\mathcal{L}}) - \bar{D}_{w}(x_{2},\pi_{2},\pi_{2}^{\mathcal{L}})\| \\ &= \left\| \mathbb{E}_{s \sim d_{\rho}^{\pi_{1}},a \sim \pi_{1}(\cdot|s),\bar{s} \sim d_{\rho}^{\pi_{1}^{\mathcal{L}}},\bar{a} \sim \pi_{1}^{\mathcal{L}}(\cdot|\bar{s})),\xi \sim \mu} [D_{w}(x_{1},\pi_{1},\pi_{1}^{\mathcal{L}},s,a,\bar{s},\bar{a},\xi)] \\ &- \mathbb{E}_{s \sim d_{\rho}^{\pi_{2}},a \sim \pi_{2}(\cdot|s),\bar{s} \sim d_{\rho}^{\pi_{2}^{\mathcal{L}}},\bar{a} \sim \pi_{2}^{\mathcal{L}}(\cdot|\bar{s})),\xi \sim \mu} [D_{w}(x_{2},\pi_{2},\pi_{2}^{\mathcal{L}},s,a,\bar{s},\bar{a},\xi)] \right\| \\ &= \left\| \sum_{s,a,\bar{s},\bar{a},\xi} \left(d_{\rho}^{\pi_{1}}(s)\pi_{1}(a \mid s) d_{\rho}^{\pi_{1}^{\mathcal{L}}}(\bar{s})\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) - d_{\rho}^{\pi_{2}}(s)\pi_{2}(a \mid s) d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s})\pi_{2}^{\mathcal{L}}(\bar{a} \mid \bar{s}) \right) \mu(\xi) D_{w}(x_{2},\pi_{2},\pi_{2}^{\mathcal{L}},s,a,\bar{s},\bar{a},\xi) \\ &+ \mathbb{E}_{s \sim d_{\rho}^{\pi_{1}},a \sim \pi_{1}(\cdot|s),\bar{s} \sim d_{\rho}^{\pi_{1}^{\mathcal{L}}},\bar{a} \sim \pi_{1}^{\mathcal{L}}(\cdot|\bar{s})),\xi \sim \mu} [D_{w}(x_{1},\pi_{1},\pi_{1}^{\mathcal{L}},s,a,\bar{s},\bar{a},\xi) - D_{w}(x_{2},\pi_{2},\pi_{2}^{\mathcal{L}},s,a,\bar{s},\bar{a},\xi)] \right\| \\ &\leq \left\| \mathbb{E}_{s \sim d_{\rho}^{\pi_{1}},a \sim \pi_{1}(\cdot|s),\bar{s} \sim d_{\rho}^{\pi_{1}^{\mathcal{L}}},\bar{a} \sim \pi_{1}^{\mathcal{L}}(\cdot|\bar{s})),\xi \sim \mu} [D_{w}(x_{1},\pi_{1},\pi_{1}^{\mathcal{L}},s,a,\bar{s},\bar{a},\xi) - D_{w}(x_{2},\pi_{2},\pi_{2}^{\mathcal{L}},s,a,\bar{s},\bar{a},\xi)] \right\| \\ &+ \frac{B_{D}}{w} \left\| \sum_{s,a,\bar{s},\bar{a}} \left(d_{\rho}^{\pi_{1}}(s)\pi_{1}(a \mid s) d_{\rho}^{\pi_{1}^{\mathcal{L}}}(\bar{s})\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) - d_{\rho}^{\pi_{2}}(s)\pi_{2}(a \mid s) d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s})\pi_{2}^{\mathcal{L}}(\bar{a} \mid \bar{s}) \right) \right\| \\ &+ \frac{B_{D}}{w} \left\| \sum_{s,a,\bar{s},\bar{a}} \left(d_{\rho}^{\pi_{1}}(s)\pi_{1}(a \mid s) - d_{\rho}^{\pi_{2}}(s)\pi_{2}(a \mid s) \right) d_{\rho}^{\pi_{1}^{\mathcal{L}}}(\bar{s})\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) \right\| \\ &+ \frac{B_{D}}{w} \left\| \sum_{s,a,\bar{s},\bar{a}} \left(d_{\rho}^{\pi_{1}^{\mathcal{L}}}(\bar{s})\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) - d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s})\pi_{2}^{\mathcal{L}}(\bar{a} \mid \bar{s}) \right) d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s})\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) \right\| \\ &+ \frac{B_{D}}{w} \left\| \sum_{s,a,\bar{s},\bar{a}} \left(d_{\rho}^{\pi_{1}^{\mathcal{L}}}(\bar{s})\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) - d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s})\pi_{2}^{\mathcal{L}}(\bar{a} \mid \bar{s}) \right) d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s})\pi_{2}^{\mathcal{L}}(\bar{a} \mid \bar{s}) \right\| \\ &+ \frac{B_{D}}{w} \left\| \sum_{s,a,\bar{s},\bar{a}} \left(d_{\rho}^{\pi_{1}^{\mathcal{L}}}(\bar{s})\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) - d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s})\pi_{2}^{\mathcal{L}}(\bar{a} \mid \bar{s}) \right) d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s} \mid \bar{s}) \right\| \\$$

To bound the first term of (109), note that by Jensen's inequality it suffices to bound the norm of the term within the expectation

$$\begin{split} & \|D_{w}(x_{1}, \pi_{1}, \pi_{1}^{\mathcal{L}}, s, a, \bar{s}, \bar{a}, \xi) - D_{w}(x_{2}, \pi_{2}, \pi_{2}^{\mathcal{L}}, s, a, \bar{s}, \bar{a}, \xi)\| \\ & = \left\| \widetilde{\nabla}_{x} f(x_{1}, \pi_{1}^{\mathcal{L}}, \xi) + \frac{1}{w} \left(\nabla_{x} r_{x_{1}}(s, a) - \nabla_{x} r_{x_{1}}(\bar{s}, \bar{a}) \right) - \widetilde{\nabla}_{x} f(x_{2}, \pi_{2}^{\mathcal{L}}, \xi) - \frac{1}{w} \left(\nabla_{x} r_{x_{2}}(s, a) - \nabla_{x} r_{x_{2}}(\bar{s}, \bar{a}) \right) \right\| \\ & \leq \|\widetilde{\nabla}_{x} f(x_{1}, \pi_{1}^{\mathcal{L}}, \xi) - \widetilde{\nabla}_{x} f(x_{2}, \pi_{2}^{\mathcal{L}}, \xi)\| + \frac{1}{w} \|\nabla_{x} r_{x_{1}}(s, a) - \nabla_{x} r_{x_{2}}(s, a)\| + \frac{1}{w} \|\nabla_{x} r_{x_{1}}(\bar{s}, \bar{a}) - \nabla_{x} r_{x_{2}}(\bar{s}, \bar{a})\| \\ & \leq L_{f} \left(\|x_{1} - x_{2}\| + \|\pi_{1}^{\mathcal{L}} - \pi_{2}^{\mathcal{L}}\| \right) + \frac{2L_{r}}{w} \|x_{1} - x_{2}\|. \end{split}$$

$$(110)$$

To bound the second term of (109),

$$\begin{split} \left| \sum_{s,a} (d_{\rho}^{\pi_{1}}(s)\pi_{1}(a \mid s) - d_{\rho}^{\pi_{2}}(s)\pi_{2}(a \mid s)) \right| + \left| \sum_{\bar{s},\bar{a}} (d_{\rho}^{\pi_{1}^{\mathcal{L}}}(\bar{s})\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) - d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s})\pi_{2}^{\mathcal{L}}(\bar{a} \mid \bar{s})) \right| \\ \leq \left| \sum_{s,a} (d_{\rho}^{\pi_{1}}(s) - d_{\rho}^{\pi_{2}}(s))\pi_{1}(a \mid s) \right| + \left| \sum_{s,a} (\pi_{1}(a \mid s) - \pi_{2}(a \mid s))d_{\rho}^{\pi_{2}}(s) \right| \\ + \left| \sum_{\bar{s},\bar{a}} (d_{\rho}^{\pi_{1}^{\mathcal{L}}}(\bar{s}) - d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s}))\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) \right| + \left| \sum_{\bar{s},\bar{a}} (\pi_{1}^{\mathcal{L}}(\bar{a} \mid \bar{s}) - \pi_{2}^{\mathcal{L}}(\bar{a} \mid \bar{s}))d_{\rho}^{\pi_{2}^{\mathcal{L}}}(\bar{s}) \right| \\ \leq \|d_{\rho}^{\pi_{1}} - d_{\rho}^{\pi_{2}}\| + \|\pi_{1} - \pi_{2}\| + \|d_{\rho}^{\pi_{1}^{\mathcal{L}}} - d_{\rho}^{\pi_{2}^{\mathcal{L}}}\| + \|\pi_{1}^{\mathcal{L}} - \pi_{2}^{\mathcal{L}}\| \\ \leq \frac{1}{1 - \gamma} \|\pi_{1} - \pi_{2}\| + \frac{1}{1 - \gamma} \|\pi_{1}^{\mathcal{L}} - \pi_{2}^{\mathcal{L}}\|, \end{split} \tag{111}$$

where the final inequality follows from Lemma 4.

Substituting (110) and (111) into (109), we have

$$\begin{split} &\|\bar{D}_{w}(x_{1}, \pi_{1}, \pi_{1}^{\mathcal{L}}) - \bar{D}_{w}(x_{2}, \pi_{2}, \pi_{2}^{\mathcal{L}})\| \\ &\leq L_{f} \Big(\|x_{1} - x_{2}\| + \|\pi_{1}^{\mathcal{L}} - \pi_{2}^{\mathcal{L}}\| \Big) + \frac{2L_{r}}{w} \|x_{1} - x_{2}\| + \frac{B_{D}}{(1 - \gamma)w} \|\pi_{1} - \pi_{2}\| + \frac{B_{D}}{(1 - \gamma)w} \|\pi_{1}^{\mathcal{L}} - \pi_{2}^{\mathcal{L}}\| \\ &= (L_{f} + \frac{2L_{r}}{w}) \|x_{1} - x_{2}\| + (L_{f} + \frac{B_{D}}{(1 - \gamma)w}) \|\pi_{1}^{\mathcal{L}} - \pi_{2}^{\mathcal{L}}\| + \frac{B_{D}}{(1 - \gamma)w} \|\pi_{1} - \pi_{2}\| \\ &\leq \frac{3L_{r}}{w} \|x_{1} - x_{2}\| + \frac{2B_{D}}{(1 - \gamma)w} \|\pi_{1}^{\mathcal{L}} - \pi_{2}^{\mathcal{L}}\| + \frac{B_{D}}{(1 - \gamma)w} \|\pi_{1} - \pi_{2}\|, \end{split}$$

$$(112)$$

where the second inequality simplifies terms under the condition $w \leq \min\{\frac{L_r}{L_f}, \frac{B_D}{(1-\gamma)L_f}\}$.

Similarly, by the definition of $\bar{F}_{w,\tau}$ in (30),

$$\begin{split} &\|\bar{F}_{w,\tau}(x_{1},\theta_{1},V_{1}) - \bar{F}_{w,\tau}(x_{2},\theta_{2},V_{2})\| \\ &= \left\| \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta_{1}}},a \sim \pi_{\theta_{1}}(\cdot|s),s' \sim \mathcal{P}(\cdot|s,a),\xi \sim \mu} [F_{w,\tau}(x_{1},\theta_{1},V_{1},s,a,s',\xi)] \right\| \\ &- \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta_{2}}},a \sim \pi_{\theta_{2}}(\cdot|s),s' \sim \mathcal{P}(\cdot|s,a),\xi \sim \mu} [F_{w,\tau}(x_{2},\theta_{2},V_{2},s,a,s',\xi)] \right\| \\ &= \left\| \sum_{s,a,s',\xi} \left(d_{\rho}^{\pi_{\theta_{1}}}(s) \pi_{\theta_{1}}(a \mid s) \mathcal{P}(s' \mid s,a) - d_{\rho}^{\pi_{\theta_{2}}}(s) \pi_{\theta_{2}}(a \mid s) \mathcal{P}(s' \mid s,a) \right) \mu(\xi) F_{w,\tau}(x_{2},\theta_{2},V_{2},s,a,s',\xi) \right\| \\ &+ \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta_{1}}},a \sim \pi_{\theta_{1}}(\cdot|s),s' \sim \mathcal{P}(\cdot|s,a),\xi \sim \mu} [F_{w,\tau}(x_{1},\theta_{1},V_{1},s,a,s',\xi) - F_{w,\tau}(x_{2},\theta_{2},V_{2},s,a,s',\xi)] \right\| \\ &\leq \left\| \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta_{1}}},a \sim \pi_{\theta_{1}}(\cdot|s),s' \sim \mathcal{P}(\cdot|s,a),\xi \sim \mu} [F_{w,\tau}(x_{1},\theta_{1},V_{1},s,a,s',\xi) - F_{w,\tau}(x_{2},\theta_{2},V_{2},s,a,s',\xi)] \right\| \\ &+ B_{F} \left| \sum_{s,a,s'} \left(d_{\rho}^{\pi_{\theta_{1}}}(s) \pi_{\theta_{1}}(a \mid s) \mathcal{P}(s' \mid s,a) - d_{\rho}^{\pi_{\theta_{2}}}(s) \pi_{\theta_{2}}(a \mid s) \mathcal{P}(s' \mid s,a) \right) \right| \\ &= \left\| \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta_{1}}},a \sim \pi_{\theta_{1}}(\cdot|s),s' \sim \mathcal{P}(\cdot|s,a),\xi \sim \mu} [F_{w,\tau}(x_{1},\theta_{1},V_{1},s,a,s',\xi) - F_{w,\tau}(x_{2},\theta_{2},V_{2},s,a,s',\xi)] \right\| \\ &+ B_{F} \left| \sum_{s,a} \left(d_{\rho}^{\pi_{\theta_{1}}}(s) \pi_{\theta_{1}}(a \mid s) - d_{\rho}^{\pi_{\theta_{2}}}(s) \pi_{\theta_{2}}(a \mid s) \right) \right|. \end{aligned}$$

$$(113)$$

To bound the first term of (109), note that by Jensen's inequality it suffices to bound the norm of the term within the expectation

$$\begin{aligned} &\|F_{w,\tau}(x_{1},\theta_{1},V_{1},s,a,s',\xi) - F_{w,\tau}(x_{2},\theta_{2},V_{2},s,a,s',\xi)\| \\ &= \left\| \left(r_{x_{1}}(s,a) + \tau E(\pi_{\theta_{1}},s) + \gamma V_{1}(s') - V_{1}(s) \right) \nabla_{\theta} \log \pi_{\theta_{1}}(a \mid s) - w \widetilde{\nabla}_{\theta} f(x_{1},\pi_{\theta_{1}},\xi) \right. \\ &- \left(r_{x_{2}}(s,a) + \tau E(\pi_{\theta_{2}},s) + \gamma V_{2}(s') - V_{2}(s) \right) \nabla_{\theta} \log \pi_{\theta_{2}}(a \mid s) + w \widetilde{\nabla}_{\theta} f(x_{2},\pi_{\theta_{2}},\xi) \right\| \\ &\leq w \|\widetilde{\nabla}_{\theta} f(x_{1},\pi_{\theta_{1}},\xi) - \widetilde{\nabla}_{\theta} f(x_{2},\pi_{\theta_{2}},\xi)\| \\ &+ \|\nabla_{\theta} \log \pi_{\theta_{1}}(a \mid s) - \nabla_{\theta} \log \pi_{\theta_{2}}(a \mid s) \| \left| r_{x_{1}}(s,a) + \tau E(\pi_{\theta_{1}},s) + \gamma V_{1}(s') - V_{1}(s) \right| \\ &+ \|\nabla_{\theta} \log \pi_{\theta_{2}}(a \mid s) \| \left| r_{x_{1}}(s,a) - r_{x_{2}}(s,a) + \tau E(\pi_{\theta_{1}},s) - \tau E(\pi_{\theta_{2}},s) + \gamma V_{1}(s') - \gamma V_{2}(s') - V_{1}(s) + V_{2}(s) \right| \\ &\leq L_{f} w(\|x_{1} - x_{2}\| + \|\pi_{\theta_{1}} - \pi_{\theta_{2}}\|) \\ &+ \|r_{x_{1}}(s,a) + \tau E(\pi_{\theta_{1}},s) + \gamma V_{1}(s') - V_{1}(s) \|\theta_{1} - \theta_{2}\| \\ &+ \|\nabla_{\theta} \log \pi_{\theta_{2}}(a \mid s) \| \left(L_{r} \|x_{1} - x_{2}\| + \log |\mathcal{A}| \tau \|\theta_{1} - \theta_{2}\| + (1 + \gamma) \|V_{1} - V_{2}\| \right) \\ &\leq L_{r} (\|x_{1} - x_{2}\| + \|\theta_{1} - \theta_{2}\|) + (1 + \log |\mathcal{A}| + \frac{\gamma}{1 - \gamma} + \frac{1}{1 - \gamma}) \|\theta_{1} - \theta_{2}\| \\ &+ 2 \left(L_{r} \|x_{1} - x_{2}\| + \log |\mathcal{A}| \|\theta_{1} - \theta_{2}\| + (1 + \gamma) \|V_{1} - V_{2}\| \right) \\ &\leq 3L_{r} \|x_{1} - x_{2}\| + (L_{r} + 2\log |\mathcal{A}| + \frac{2}{1 - \gamma} + 1) \|\theta_{1} - \theta_{2}\| + 4 \|V_{1} - V_{2}\|, \end{aligned}$$

$$(114)$$

where the second inequality is due to the 1-Lipschitz continuity of $\nabla_{\theta} \log \pi_{\theta}$ and the $\log |\mathcal{A}|$ -Lipschitz continuity of the entropy function with respect to the softmax parameter, and the third inequality is due to $w \leq \frac{L_r}{L_f}$, $\tau \leq 1$ and the fact that $\|\nabla_{\theta} \log \pi_{\theta}(a \mid s)\|_2 \leq 2$ for all θ .

To bound the second term of (113),

$$\left| \sum_{s,a} (d_{\rho}^{\pi_{\theta_{1}}}(s) \pi_{\theta_{1}}(a \mid s) - d_{\rho}^{\pi_{\theta_{2}}}(s) \pi_{\theta_{2}}(a \mid s)) \right| \leq \left| \sum_{s,a} (d_{\rho}^{\pi_{\theta_{1}}}(s) - d_{\rho}^{\pi_{\theta_{2}}}(s)) \pi_{\theta_{1}}(a \mid s) \right| + \left| \sum_{s,a} (\pi_{\theta_{1}}(a \mid s) - \pi_{\theta_{2}}(a \mid s)) d_{\rho}^{\pi_{\theta_{2}}}(s) \right| \leq \|d_{\rho}^{\pi_{\theta_{1}}} - d_{\rho}^{\pi_{\theta_{2}}}\| + \|\pi_{\theta_{1}} - \pi_{\theta_{2}}\| \leq \frac{1}{1 - \gamma} \|\pi_{\theta_{1}} - \pi_{\theta_{2}}\| \leq \frac{1}{1 - \gamma} \|\theta_{1} - \theta_{2}\|,$$
(115)

where the third inequality follows from Lemma 4.

Substituting (114) and (115) into (111), we have

$$\begin{split} &\|\bar{F}_{w,\tau}(x_1,\theta_1,V_1) - \bar{F}_{w,\tau}(x_2,\theta_2,V_2)\| \\ &\leq 3L_r\|x_1 - x_2\| + (L_r + 2\log|\mathcal{A}| + \frac{2}{1-\gamma} + 1)\|\theta_1 - \theta_2\| + 4\|V_1 - V_2\| + \frac{B_F}{1-\gamma}\|\theta_1 - \theta_2\| \\ &= 3L_r\|x_1 - x_2\| + (L_r + 2\log|\mathcal{A}| + \frac{2+B_F}{1-\gamma} + 1)\|\theta_1 - \theta_2\| + 4\|V_1 - V_2\|. \end{split}$$

By the definition of \bar{G}_{τ} in (31),

$$\begin{split} & \| \bar{G}_{\tau}(x_{1}, \theta_{1}, V_{1}) - \bar{G}_{\tau}(x_{2}, \theta_{2}, V_{2}) \| \\ & = \left\| \mathbb{E}_{s \sim d_{\rho}^{\pi\theta_{1}}, a \sim \pi_{\theta_{1}}(\cdot|s), s' \sim \mathcal{P}(\cdot|s, a)} [G_{\tau}(x_{1}, \theta_{1}, V_{1}, s, a, s')] \right\| \\ & - \mathbb{E}_{s \sim d_{\rho}^{\pi\theta_{2}}, a \sim \pi_{\theta_{2}}(\cdot|s), s' \sim \mathcal{P}(\cdot|s, a)} [G_{\tau}(x_{2}, \theta_{2}, V_{2}, s, a, s')] \right\| \\ & = \left\| \sum_{s, a, s'} \left(d_{\rho}^{\pi\theta_{1}}(s) \pi_{\theta_{1}}(a \mid s) \mathcal{P}(s' \mid s, a) - d_{\rho}^{\pi\theta_{2}}(s) \pi_{\theta_{2}}(a \mid s) \mathcal{P}(s' \mid s, a) \right) G_{\tau}(x_{2}, \theta_{2}, V_{2}, s, a, s', \xi) \\ & + \mathbb{E}_{s \sim d_{\rho}^{\pi\theta_{1}}, a \sim \pi_{\theta_{1}}(\cdot|s), s' \sim \mathcal{P}(\cdot|s, a)} [G_{\tau}(x_{1}, \theta_{1}, V_{1}, s, a, s') - G_{\tau}(x_{2}, \theta_{2}, V_{2}, s, a, s')] \right\| \\ & \leq \left\| \mathbb{E}_{s \sim d_{\rho}^{\pi\theta_{1}}, a \sim \pi_{\theta_{1}}(\cdot|s), s' \sim \mathcal{P}(\cdot|s, a)} [G_{\tau}(x_{1}, \theta_{1}, V_{1}, s, a, s') - G_{\tau}(x_{2}, \theta_{2}, V_{2}, s, a, s')] \right\| \\ & + B_{G} \left| \sum_{s, a, s'} \left(d_{\rho}^{\pi\theta_{1}}(s) \pi_{\theta_{1}}(a \mid s) \mathcal{P}(s' \mid s, a) - d_{\rho}^{\pi\theta_{2}}(s) \pi_{\theta_{2}}(a \mid s) \mathcal{P}(s' \mid s, a) \right) \right| \\ & = \left\| \mathbb{E}_{s \sim d_{\rho}^{\pi\theta_{1}}, a \sim \pi_{\theta_{1}}(\cdot|s), s' \sim \mathcal{P}(\cdot|s, a), \xi \sim \mu} [F_{w, \tau}(x_{1}, \theta_{1}, V_{1}, s, a, s', \xi) - F_{w, \tau}(x_{2}, \theta_{2}, V_{2}, s, a, s', \xi) \right] \right\| \\ & + B_{G} \left| \sum_{s, a} \left(d_{\rho}^{\pi\theta_{1}}(s) \pi_{\theta_{1}}(a \mid s) - d_{\rho}^{\pi\theta_{2}}(s) \pi_{\theta_{2}}(a \mid s) \right) \right| \\ & \leq \left\| \mathbb{E}_{s \sim d_{\rho}^{\pi\theta_{1}}, a \sim \pi_{\theta_{1}}(\cdot|s), s' \sim \mathcal{P}(\cdot|s, a), \xi \sim \mu} [F_{w, \tau}(x_{1}, \theta_{1}, V_{1}, s, a, s', \xi) - F_{w, \tau}(x_{2}, \theta_{2}, V_{2}, s, a, s', \xi) \right] \right\| \\ & + \frac{B_{G}}{1 - \rho} \|\theta_{1} - \theta_{2}\|, \end{split}$$
(116)

where the last inequality follows from (115).

The first of (116) can be bounded as follows. Again, by Jensen's inequality it suffices to bound the norm of the term within the expectation

$$||G_{\tau}(x_1, \theta_1, V_1, s, a, s') - G_{\tau}(x_2, \theta_2, V_2, s, a, s')||$$

$$= \left\| \left(r_{x_1}(s, a) + \tau E(\pi_{\theta_1}, s) + \gamma V_1(s') - V_1(s) \right) e_s - \left(r_{x_2}(s, a) + \tau E(\pi_{\theta_2}, s) + \gamma V_2(s') - V_2(s) \right) e_s \right\|$$

$$\leq L_r \|x_1 - x_2\| + \log |\mathcal{A}|\tau \|\theta_1 - \theta_2\| + (1 + \gamma) \|V_1 - V_2\|,$$
(117)

where in the last inequality we again use the $\log |\mathcal{A}|$ -Lipschitz continuity of the entropy function with respect to the softmax parameter.

Combining (116) and (117), we have under τ

$$\begin{split} \|\bar{G}_{\tau}(x_{1},\theta_{1},V_{1}) - \bar{G}_{\tau}(x_{2},\theta_{2},V_{2})\| &\leq L_{r}\|x_{1} - x_{2}\| + \log|\mathcal{A}|\tau\|\theta_{1} - \theta_{2}\| + (1+\gamma)\|V_{1} - V_{2}\| \\ &+ \frac{B_{G}}{1-\gamma}\|\theta_{1} - \theta_{2}\| \\ &\leq L_{r}\|x_{1} - x_{2}\| + (\frac{B_{G}}{1-\gamma} + \log|\mathcal{A}|)\|\theta_{1} - \theta_{2}\| + 2\|V_{1} - V_{2}\|. \end{split}$$

D.7 Proof of Lemma 7

The proof proceeds in a manner similar to Kwon et al. [2023][Lemma 3.2], with strong convexity replaced by the PL condition.

First, we consider a fixed τ . Recall the definition of $\pi_{w,\tau}^{\star}$ in Section 3.1. Let $\theta_{w_1,\tau}^{\star}(x_1)$ denote a softmax parameter that encodes $\pi_{w_1,\tau}^{\star}(x_1)$. The optimality condition of $\theta_{w_1,\tau}^{\star}(x_1)$ indicates

$$\nabla_{\theta} \mathcal{L}_{w_1,\tau}(x_1, \pi_{\theta_{w_1,\tau}^{\star}(x_1)}) = \nabla_{\theta} f(x_1, \pi_{\theta_{w_1,\tau}^{\star}(x_1)}) - \frac{1}{w_1} \nabla_{\pi} J_{\tau}(x_1, \pi_{\theta_{w_1,\tau}^{\star}(x_1)}) = 0, \quad (118)$$

which obviously implies

$$\|\nabla_{\theta} J_{\tau}(x_1, \pi_{\theta_{w_1, \tau}^{\star}(x_1)})\| = w_1 \|\nabla_{\theta} f(x_1, \pi_{\theta_{w_1, \tau}^{\star}(x_1)})\| \le L_f w_1.$$
(119)

Applying the relationship in (118), we have

$$\begin{split} &\nabla_{\theta} \mathcal{L}_{w_{2},\tau}(x_{2}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}) \\ &= \nabla_{\theta} f(x_{2}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}) - \frac{1}{w_{2}} \nabla_{\pi} J_{\tau}(x_{2}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}) \\ &= \left(\nabla_{\theta} f(x_{2}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}) - \nabla_{\theta} f(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})})\right) - \frac{1}{w_{2}} \left(\nabla_{\theta} J_{\tau}(x_{2}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}) - \nabla_{\theta} J_{\tau}(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})})\right) \\ &+ \nabla_{\theta} f(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}) - \frac{1}{w_{2}} \nabla_{\pi} J_{\tau}(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}) \\ &= \left(\nabla_{\theta} f(x_{2}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}) - \nabla_{\theta} f(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})})\right) - \frac{1}{w_{2}} \left(\nabla_{\theta} J_{\tau}(x_{2}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}) - \nabla_{\theta} J_{\tau}(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})})\right) \\ &- \left(\frac{1}{w_{2}} - \frac{1}{w_{1}}\right) \nabla_{\theta} J_{\tau}(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})}). \end{split}$$

Taking the norm,

$$\|\nabla_{\theta} \mathcal{L}_{w_{2},\tau}(x_{2}, \pi_{\theta_{w_{1},\tau}(x_{1})})\|$$

$$\leq \|\nabla_{\theta} f(x_{2}, \pi_{w_{1},\tau}^{\star}(x_{1})) - \nabla_{\pi} f(x_{1}, \pi_{w_{1},\tau}^{\star}(x_{1}))\| + \frac{1}{w_{2}} \|\nabla_{\pi} J_{\tau}(x_{2}, \pi_{w_{1},\tau}^{\star}(x_{1})) - \nabla_{\theta} J_{\tau}(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})})\|$$

$$+ \left| \frac{1}{w_{2}} - \frac{1}{w_{1}} \right| \|\nabla_{\theta} J_{\tau}(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})})\|$$

$$\leq L_{f} \|x_{1} - x_{2}\| + \frac{L_{V}}{w_{2}} \|x_{1} - x_{2}\| + \left| \frac{1}{w_{2}} - \frac{1}{w_{1}} \right| \|\nabla_{\theta} J_{\tau}(x_{1}, \pi_{\theta_{w_{1},\tau}^{\star}(x_{1})})\|$$

$$\leq (L_{f} + \frac{L_{V}}{w_{2}}) \|x_{1} - x_{2}\| + \frac{L_{f} |w_{1} - w_{2}|}{w_{2}}, \tag{120}$$

where the second inequality follows from the Lipschitz continuous gradients of f and J_{τ} , and the third inequality plugs in (119).

Due to (25), we have

$$\|\nabla_{\theta} \mathcal{L}_{w_2,\tau}(x, \pi_{\theta_{w_1,\tau}^{\star}(x_1)})\| \ge \frac{C_L \tau}{2w_2} \|\pi_{w_1,\tau}^{\star}(x_1) - \pi_{w_2,\tau}^{\star}(x_2)\|.$$
 (121)

Combining (120) and (121),

$$\frac{C_L \tau}{2w_2} \|\pi_{w_1,\tau}^{\star}(x_1) - \pi_{w_2,\tau}^{\star}(x_2)\| \le (L_f + \frac{L_V}{w_2}) \|x_1 - x_2\| + \frac{L_f |w_1 - w_2|}{w_2}.$$

For any $w_2 > 0$, this simplifies to

$$\|\pi_{w_1,\tau}^{\star}(x_1) - \pi_{w_2,\tau}^{\star}(x_2)\| \le \left(\frac{2L_f w_2}{C_L \tau} + \frac{2L_V}{C_L \tau}\right) \|x_1 - x_2\| + \frac{2L_f |w_1 - w_2|}{C_L \tau}.$$
 (122)

Recognizing $\pi_{\tau}^{\star}(x) = \lim_{w \to 0^{+}} \pi_{w,\tau}^{\star}(x)$, we have from (122)

$$\|\pi_{\tau}^{\star}(x) - \pi_{w,\tau}^{\star}(x)\| \le \frac{2L_f w}{C_L \tau}.$$

Now, we fix w, x and show the bound on $\|\pi_{w,\tau_1}^{\star}(x) - \pi_{w,\tau_2}^{\star}(x)\|$. We use $\theta_{w,\tau_1}(x)^{\star}$ to denote a softmax parameter for $\pi_{w,\tau_1}(x)^{\star}$. The optimality condition of $\theta_{w,\tau_1}(x)^{\star}$ indicates

$$\nabla_{\theta} \mathcal{L}_{w,\tau_1}(x, \pi_{\theta_{w,\tau_1}^{\star}(x)}) = \nabla_{\theta} f(x, \pi_{\theta_{w,\tau_1}^{\star}(x)}) - \frac{1}{w} \nabla_{\theta} J_{\tau_1}(x, \pi_{\theta_{w,\tau_1}^{\star}(x)}) = 0.$$

Applying the equation, we get

$$\nabla_{\theta} \mathcal{L}_{w,\tau_{2}}(x, \pi_{\theta_{w,\tau_{1}}^{\star}(x)}) = \nabla_{\theta} f(x, \pi_{\theta_{w,\tau_{1}}^{\star}(x)}) - \frac{1}{w} \nabla_{\theta} J_{\tau_{2}}(x, \pi_{\theta_{w,\tau_{1}}^{\star}(x)})$$

$$= \frac{1}{w} \Big(\nabla_{\theta} J_{\tau_{1}}(x, \pi_{\theta_{w,\tau_{1}}^{\star}(x)}) - \nabla_{\theta} J_{\tau_{2}}(x, \pi_{\theta_{w,\tau_{1}}^{\star}(x)}) \Big). \tag{123}$$

The regularized RL objective has a closed-form expression (see Mei et al. [2020][Lemma 10])

$$\frac{\partial J_{\tau}(x, \pi_{\theta})}{\partial \theta(s, a)} = \frac{d_{\rho}^{\pi_{\theta}}(s)}{1 - \gamma} \cdot \pi_{\theta}(a \mid s) \cdot A_{\tau}^{x, \pi_{\theta}}(s, a),$$

which in combination with (123) implies

$$\|\nabla_{\theta} \mathcal{L}_{w,\tau_2}(x, \pi_{\theta_{w,\tau_1}^{\star}(x)})\| \leq \|\nabla_{\theta} \mathcal{L}_{w,\tau_2}(x, \pi_{\theta_{w,\tau_1}^{\star}(x)})\|_1$$

$$\leq \frac{1}{(1-\gamma)w} \sum_{s,a} \pi_{\theta_{w,\tau_1}^{\star}(x)}(a \mid s) \Big(A_{\tau_1}^{x,\pi_{\theta_{w,\tau_1}^{\star}(x)}}(s,a) - A_{\tau_2}^{x,\pi_{\theta_{w,\tau_1}^{\star}(x)}}(s,a) \Big). \tag{124}$$

Due to Zeng et al. [2022a] [Lemma 3], we have for any s

$$|V_{\tau_1}^{x,\pi}(s) - V_{\tau_2}^{x,\pi}(s)| \le |\tau_1 - \tau_2| \log |\mathcal{A}|.$$

The definitions of the Q function and advantage function in (23) imply

$$|Q_{\tau_1}^{x,\pi}(s,a) - Q_{\tau_2}^{x,\pi}(s,a)| \le \gamma |V_{\tau_1}^{x,\pi}(s) - V_{\tau_2}^{x,\pi}(s)| \le \gamma |\tau_1 - \tau_2| \log |\mathcal{A}|,$$

and

$$\left| \sum_{a} \pi(a \mid s) A_{\tau_{1}}^{x,\pi}(s,a) - A_{\tau_{2}}^{x,\pi}(s,a) \right|$$

$$\leq \sum_{a} \pi(a \mid s) |Q_{\tau_{1}}^{x,\pi}(s,a) - Q_{\tau_{2}}^{x,\pi}(s,a)| + |V_{\tau_{1}}^{x,\pi}(s) - V_{\tau_{2}}^{x,\pi}(s)| + |\tau_{1} - \tau_{2}| E(\pi,s)$$

$$\leq 3|\tau_{1} - \tau_{2}| \log |\mathcal{A}|. \tag{125}$$

Plugging (125) into (124),

$$\|\nabla_{\theta} \mathcal{L}_{w,\tau_2}(x, \pi_{\theta_{w,\tau_1}^{\star}(x)})\| \leq \frac{3|\tau_1 - \tau_2||\mathcal{S}|\log|\mathcal{A}|}{(1-\gamma)w}.$$

Again, due to (25),

$$\frac{C_L \tau_2}{2w} \|\pi_{w,\tau_1}^{\star}(x) - \pi_{w,\tau_2}^{\star}(x)\| \le \|\nabla_{\pi} \mathcal{L}_{w,\tau_2}(x, \pi_{w,\tau_1}^{\star}(x))\| \le \frac{3|\tau_1 - \tau_2||\mathcal{S}|\log|\mathcal{A}|}{(1 - \gamma)w}.$$

This leads to

$$\|\pi_{w,\tau_1}^{\star}(x) - \pi_{w,\tau_2}^{\star}(x)\| \le \frac{6|\tau_1 - \tau_2||\mathcal{S}|\log|\mathcal{A}|}{(1 - \gamma)C_L\tau_2}.$$
 (126)

Putting (122) and (126) together,

$$\begin{split} \|\pi_{w_1,\tau_1}^{\star}(x_1) - \pi_{w_2,\tau_2}^{\star}(x_2)\| &\leq \|\pi_{w_1,\tau_1}^{\star}(x_1) - \pi_{w_2,\tau_1}^{\star}(x_2)\| + \|\pi_{w_2,\tau_2}^{\star}(x_2) - \pi_{w_2,\tau_1}^{\star}(x_2)\| \\ &\leq (\frac{2L_f w_2}{C_L \tau_1} + \frac{2L_V}{C_L \tau_1}) \|x_1 - x_2\| + \frac{2L_f |w_1 - w_2|}{C_L \tau_1} + \frac{6|\tau_1 - \tau_2||\mathcal{S}|\log|\mathcal{A}|}{(1 - \gamma)C_L \tau_1}. \end{split}$$

D.8 Proof of Lemma 8

We know from Lemma 1 that $\pi^*(x)$, defined in (5), is the limit point of $\pi^*_{\tau}(x)$ as $\tau \to 0$ Let $\theta^*_{\tau}(x)$ denote a softmax parameter for $\pi^*_{\tau}(x)$. By the first-order optimality condition, we have

$$\nabla_{\theta} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) = 0.$$

We further differentiate with respect to τ . Due to the differentiation chain rule,

$$\frac{d}{d\tau} \nabla_{\theta} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) = \nabla_{\tau, \theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) + \nabla_{\theta, \theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) \cdot \frac{d\theta_{\tau}^{\star}(x)}{d\tau} = 0.$$

As $\|\nabla^2_{\theta,\theta}J_{\tau}(x,\pi_{\theta_{\tau}^*(x)})\|$ is lower bounded by $\underline{\sigma}$ due to Assumption 2, we have for any $\tau\geq 0$

$$\left\| \frac{d\theta_{\tau}^{\star}(x)}{d\tau} \right\| = \left\| -\left(\|\nabla_{\theta,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) \right)^{-1} \nabla_{\tau,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) \right\|$$

$$\leq \left\| \left(\|\nabla_{\theta,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) \right)^{-1} \right\| \left\| \nabla_{\tau,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) \right\|$$

$$\leq \frac{1}{\underline{\sigma}} \|\nabla_{\tau,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) \right\|. \tag{127}$$

It is clear from (3)

$$\nabla_{\tau} J_{\tau}(x, \pi_{\theta}) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta}}, a \sim \pi_{\theta}(\cdot | s)} [E(\pi_{\theta}, s)].$$

Therefore,

$$\nabla_{\tau,\theta}^2 J_{\tau}(x,\pi_{\theta}) = \frac{1}{1-\gamma} \nabla_{\theta} \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta}}, a \sim \pi_{\theta}(\cdot|s)} [E(\pi_{\theta}, s)]. \tag{128}$$

Zeng et al. [2022a][Lemma 6] shows that $\mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta}}, a \sim \pi_{\theta}(\cdot|s)}[E(\pi_{\theta}, s)]$ is Lipschitz with constant $\frac{4+8\log|\mathcal{A}|}{(1-\gamma)^3}$, which is equivalent to

$$\|\nabla_{\theta} \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta}}, a \sim \pi_{\theta}(\cdot|s)} [E(\pi_{\theta}, s)]\| \le \frac{4 + 8\log|\mathcal{A}|}{(1 - \gamma)^3}, \quad \forall \theta.$$
 (129)

Combining (127)-(129), we have

$$\left\| \frac{d\theta_{\tau}^{\star}(x)}{d\tau} \right\| = \frac{1}{\underline{\sigma}} \cdot \frac{1}{1-\gamma} \cdot \frac{4+8\log|\mathcal{A}|}{(1-\gamma)^3} = \frac{4+8\log|\mathcal{A}|}{\underline{\sigma}(1-\gamma)^4} = L_{\star}.$$

This implies $\theta_{\tau}^{\star}(x)$ is L_{\star} -Lipschitz with respect to τ for $\tau \geq 0$.

48

D.9 Proof of Lemma 9

It can be seen from (12) that the gradients of the Lagrangian function have the following closed-form expressions

$$\nabla_{x}\mathcal{L}_{w,\tau}(x,\pi) = \nabla_{x}f(x,\pi) + \frac{1}{w} \Big(\nabla_{x}J_{\tau}(x,\pi_{\tau}^{\star}(x)) + \nabla_{x}\pi_{\tau}^{\star}(x)\nabla_{\pi}J_{\tau}(x,\pi_{\tau}^{\star}(x)) - \nabla_{x}J_{\tau}(x,\pi) \Big)$$

$$= \nabla_{x}f(x,\pi) + \frac{1}{w} \Big(\nabla_{x}J_{\tau}(x,\pi_{\tau}^{\star}(x)) - \nabla_{x}J_{\tau}(x,\pi) \Big), \tag{130}$$

$$\nabla_{\pi}\mathcal{L}_{w,\tau}(x,\pi) = \nabla_{\pi}f(x,\pi) - \frac{1}{w}\nabla_{\pi}J_{\tau}(x,\pi), \tag{131}$$

where the equation (130) is due to the optimality condition of J_{τ} at $\pi_{\tau}^{\star}(x)$.

According to (130), we have

$$\begin{split} &\|\nabla_{x}\mathcal{L}_{w,\tau}(x,\pi) - \nabla_{x}\mathcal{L}_{w,\tau}(x',\pi')\| \\ &= \|\nabla_{x}f(x,\pi) + \frac{1}{w}\Big(\nabla_{x}J_{\tau}(x,\pi_{\tau}^{\star}(x)) - \nabla_{x}J_{\tau}(x,\pi)\Big) - \nabla_{x}f(x',\pi') - \frac{1}{w}\Big(\nabla_{x}J_{\tau}(x',\pi_{\tau}^{\star}(x')) - \nabla_{x}J_{\tau}(x',\pi')\Big)\| \\ &\leq \|\nabla_{x}f(x,\pi) - \nabla_{x}f(x',\pi')\| + \frac{1}{w}\|\nabla_{x}J_{\tau}(x,\pi_{\tau}^{\star}(x)) - \nabla_{x}J_{\tau}(x',\pi_{\tau}^{\star}(x'))\| \\ &+ \frac{1}{w}\|\nabla_{x}J_{\tau}(x,\pi) - \nabla_{x}J_{\tau}(x',\pi')\| \\ &\leq L_{f}(\|x-x'\| + \|\pi-\pi'\|) + \frac{L_{V}}{w}(\|x-x'\| + \|\pi_{\tau}^{\star}(x) - \pi_{\tau}^{\star}(x')\|) + \frac{L_{V}}{w}(\|x-x'\| + \|\pi-\pi'\|). \end{split}$$

Recognizing $\pi_{ au}^{\star}(x)=\lim_{w\to 0^{+}}\pi_{w, au}^{\star}(x),$ we have from Lemma 7

$$\|\pi_{\tau}^{\star}(x) - \pi_{\tau}^{\star}(x')\| \le \frac{2L_V}{C_L \tau} \|x - x'\|.$$

Combining the two inequalities above and imposing the condition $w, \tau \leq 1$, we get

$$\|\nabla_{x}\mathcal{L}_{w,\tau}(x,\pi) - \nabla_{x}\mathcal{L}_{w,\tau}(x',\pi')\| \leq (L_{V} + L_{f} + \frac{L_{V}(C_{L} + 2L_{V})}{C_{L}})\frac{1}{w\tau}\|x - x'\| + \frac{L_{f} + L_{V}}{w}\|\pi - \pi'\|$$

$$\leq \frac{L_{L}}{w\tau}\|x - x'\| + \frac{L_{L}}{w}\|\theta - \theta'\|.$$

According to (131), we have

$$\begin{split} &\|\nabla_{\theta} \mathcal{L}_{w,\tau}(x,\pi_{\theta}) - \nabla_{\theta} \mathcal{L}_{w,\tau}(x',\pi_{\theta'})\| \\ &= \|\nabla_{\theta} f(x,\pi_{\theta}) - \frac{1}{w} \nabla_{\theta} J_{\tau}(x,\pi_{\theta}) - \nabla_{\theta} f(x',\pi_{\theta'}) - \frac{1}{w} \nabla_{\theta} J_{\tau}(x',\pi_{\theta'})\| \\ &\leq \|\nabla_{\theta} f(x,\pi_{\theta}) - \nabla_{\theta} f(x',\pi_{\theta'})\| + \frac{1}{w} \|\nabla_{\theta} J_{\tau}(x,\pi_{\theta}) - \nabla_{\theta} J_{\tau}(x',\pi_{\theta'})\| \\ &\leq L_{f}(\|x - x'\| + \|\pi_{\theta} - \pi_{\theta'}\|) + \frac{L_{V}}{w}(\|x - x'\| + \|\pi_{\theta} - \pi_{\theta'}\|) \\ &\leq \frac{L_{L}}{w} \|x - x'\| + \frac{L_{L}}{w} \|\theta - \theta'\|. \end{split}$$

D.10 Proof of Lemma 10

Let $\theta_{\tau}^{\star}(x)$ denote a softmax parameter of $\pi_{\tau}^{\star}(x)$. By the definition of ℓ_{τ} , we have

$$\nabla \ell_{\tau}(x) = \nabla_{x} J_{\tau}(x, \pi_{\theta^{\star}(x)}),$$

due to $\nabla_{\theta} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) = 0.$

As J_{τ} is L_V -smooth (from Lemma 3) this implies

$$\begin{split} \|\nabla \ell_{\tau}(x_{1}) - \nabla \ell_{\tau}(x_{2})\| &= \|\nabla_{x} J_{\tau}(x_{1}, \pi_{\theta_{\tau}^{\star}(x_{1})}) - \nabla_{x} J_{\tau}(x_{2}, \pi_{\theta_{\tau}^{\star}(x_{2})})\| \\ &\leq L_{V} \|x_{1} - x_{2}\| + L_{V} \|\pi_{\tau}^{\star}(x_{1}) - \pi_{\tau}^{\star}(x_{2})\| \\ &\leq L_{V} \|x_{1} - x_{2}\| + L_{V} \cdot \frac{2L_{V}}{C_{L}\tau} \|x_{1} - x_{2}\| \\ &\leq \left(L_{V} + \frac{2L_{V}^{2}}{C_{L}\tau}\right) \|x_{1} - x_{2}\|, \end{split}$$

where the second inequality follows from Lemma 7 by recognizing that $\pi_{\tau}^{\star}(x) = \lim_{w \to 0} \pi_{w,\tau}^{\star}(x)$.

We next show the smoothness of $\Phi_{w,\tau}$. From (13) it can be seen

$$\begin{split} &\|\nabla_{x}\Phi_{w,\tau}(x_{1}) - \nabla_{x}\Phi_{w,\tau}(x_{2})\| \\ &\leq \left\|\nabla_{x}f(x_{1},\pi_{\tau}^{\star}(x_{1})) - \nabla_{x}f(x_{2},\pi_{\tau}^{\star}(x_{2}))\right\| + \frac{1}{w}\left\|\nabla_{x}J_{\tau}(x_{1},\pi_{\tau}^{\star}(x_{1})) - \nabla_{x}J_{\tau}(x_{2},\pi_{\tau}^{\star}(x_{2}))\right\| \\ &\quad + \frac{1}{w}\left\|\nabla_{x}J_{\tau}(x_{1},\pi_{w,\tau}^{\star}(x_{1})) - \nabla_{x}J_{\tau}(x_{2},\pi_{w,\tau}^{\star}(x_{2}))\right\| \\ &\leq L_{f}\|x_{1} - x_{2}\| + \left(L_{f} + \frac{L_{V}}{w}\right)\|\pi_{\tau}^{\star}(x_{1}) - \pi_{\tau}^{\star}(x_{2})\| + \frac{L_{V}}{w}\|\pi_{w,\tau}^{\star}(x_{1}) - \pi_{w,\tau}^{\star}(x_{2})\| \\ &\leq L_{f}\|x_{1} - x_{2}\| + \left(L_{f} + \frac{L_{V}}{w}\right) \cdot \frac{2L_{V}}{C_{L}\tau}\|x_{1} - x_{2}\| + \frac{L_{V}}{w} \cdot \left(\frac{2L_{f}w_{2}}{C_{L}\tau} + \frac{2L_{V}}{C_{L}\tau}\right)\|x_{1} - x_{2}\| \\ &\leq \left(L_{f} + \frac{4L_{f}L_{V}}{C_{L}\tau} + \frac{4L_{V}^{2}}{C_{L}w\tau}\right)\|x_{1} - x_{2}\|. \end{split}$$

Finally, we show the smoothness of Φ_{τ} . Let $\theta_{\tau}^{\star}(x)$ denote one of the softmax parameters that encodes $\pi_{\tau}^{\star}(x)$. We can express the hyper-gradient of Φ_{τ} as follows

$$\nabla_x \Phi_{\tau}(x) = \nabla_x f(x, \pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x,\theta}^2 J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) \nabla_{\theta,\theta}^2 J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)})^{-1} \nabla_{\theta} f(x, \pi_{\theta_{\tau}^{\star}(x)}).$$

This implies

$$\|\nabla_{x}\Phi_{\tau}(x_{1}) - \nabla_{x}\Phi_{\tau}(x_{2})\|$$

$$\leq \underbrace{\|\nabla_{x}f(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})}) - \nabla_{x}f(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})\|}_{T_{1}}$$

$$+ \underbrace{\|\nabla_{x,\theta}^{2}J_{\tau}(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})})\nabla_{\theta,\theta}^{2}J_{\tau}(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})})^{-1}\nabla_{\theta}f(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})}) - \nabla_{x,\theta}^{2}J_{\tau}(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})\nabla_{\theta,\theta}^{2}J_{\tau}(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})})^{-1}\nabla_{\theta}f(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})})\|}_{T_{2}}$$

$$+ \underbrace{\|\nabla_{x,\theta}^{2}J_{\tau}(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})\nabla_{\theta,\theta}^{2}J_{\tau}(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})})^{-1}\nabla_{\theta}f(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})}) - \nabla_{x,\theta}^{2}J_{\tau}(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})\nabla_{\theta,\theta}^{2}J_{\tau}(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})^{-1}\nabla_{\theta}f(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})})\|}_{T_{3}}$$

$$+ \underbrace{\|\nabla_{x,\theta}^{2}J_{\tau}(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})\nabla_{\theta,\theta}^{2}J_{\tau}(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})^{-1}\nabla_{\theta}f(x_{1},\pi_{\theta_{\tau}^{\star}(x_{1})}) - \nabla_{x,\theta}^{2}J_{\tau}(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})\nabla_{\theta,\theta}^{2}J_{\tau}(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})^{-1}\nabla_{\theta}f(x_{2},\pi_{\theta_{\tau}^{\star}(x_{2})})\|}_{T_{4}}$$

We treat each term of (132) individually. First, we bound T_1 using the smoothness property of f

$$T_{1} \leq L_{f} \left(\|x_{1} - x_{2}\| + \|\pi_{\tau}^{\star}(x_{1}) - \pi_{\tau}^{\star}(x_{2})\| \right)$$

$$\leq \left(L_{f} + \frac{2L_{f}L_{V}}{C_{L}\tau} \right) \|x_{1} - x_{2}\|, \tag{133}$$

where to derive the second inequality, we plug in the result from Lemma 7 to get $\|\pi_{\tau}^{\star}(x_1) - \pi_{\tau}^{\star}(x_2)\| \le \frac{2L_V}{C_L \tau} \|x_1 - x_2\|$ (note that $\pi_{\tau}^{\star}(x) = \lim_{w \to 0} \pi_{w,\tau}^{\star}(x)$).

As we have $\nabla_{\theta} f(x, \pi_{\theta}) \leq L_f$ from Assumption 3 and $\|\nabla^2_{\theta, \theta} J_{\tau}(x, \pi_{\theta_{\tau}^*(x)})^{-1}\| \leq \frac{1}{\underline{\sigma}}$ due to Assumption 2,

$$T_2 \leq \|\nabla^2_{x,\theta} J_{\tau}(x_1, \pi_{\theta_{\tau}^{\star}(x_1)}) - \nabla^2_{x,\theta} J_{\tau}(x_2, \pi_{\theta_{\tau}^{\star}(x_2)})\| \|\nabla^2_{\theta,\theta} J_{\tau}(x_1, \pi_{\theta_{\tau}^{\star}(x_1)})^{-1} \nabla_{\theta} f(x_1, \pi_{\theta_{\tau}^{\star}(x_1)})\|$$

$$\leq \frac{L_f}{\underline{\sigma}} \cdot L_{V,2} \Big(\|x_1 - x_2\| + \|\pi_{\tau}^{\star}(x_1) - \pi_{\tau}^{\star}(x_2)\| \Big)
\leq \Big(\frac{L_f L_{V,2}}{\underline{\sigma}} + \frac{2L_f L_V L_{V,2}}{\underline{\sigma} C_L \tau} \Big) \|x_1 - x_2\|, \tag{134}$$

where second inequality is due to Lemma 3, and the last inequality follows from an argument similar to the one in (133).

Similarly, for T_3

$$T_{3} \leq \|\nabla_{x,\theta}^{2} J_{\tau}(x_{2}, \pi_{\theta_{\tau}^{\star}(x_{2})})\| \|\nabla_{\theta,\theta}^{2} J_{\tau}(x_{1}, \pi_{\theta_{\tau}^{\star}(x_{1})})^{-1} - \nabla_{\theta,\theta}^{2} J_{\tau}(x_{2}, \pi_{\theta_{\tau}^{\star}(x_{2})})^{-1} \| \|\nabla_{\theta} f(x_{1}, \pi_{\theta_{\tau}^{\star}(x_{1})})\|$$

$$\leq L_{f} L_{V} \|\nabla_{\theta,\theta}^{2} J_{\tau}(x_{1}, \pi_{\theta_{\tau}^{\star}(x_{1})})^{-1} \| \|\nabla_{\theta,\theta}^{2} J_{\tau}(x_{2}, \pi_{\theta_{\tau}^{\star}(x_{2})}) - \nabla_{\theta,\theta}^{2} J_{\tau}(x_{1}, \pi_{\theta_{\tau}^{\star}(x_{1})})^{-1} \| \|\nabla_{\theta,\theta}^{2} J_{\tau}(x_{2}, \pi_{\theta_{\tau}^{\star}(x_{2})})^{-1} \|$$

$$\leq \frac{L_{f} L_{V}}{\underline{\sigma}^{2}} \cdot L_{V,2} \Big(\|x_{1} - x_{2}\| + \|\pi_{\tau}^{\star}(x_{1}) - \pi_{\tau}^{\star}(x_{2})\| \Big)$$

$$\leq (\frac{L_{f} L_{V} L_{V,2}}{\underline{\sigma}^{2}} + \frac{2L_{f} L_{V}^{2} L_{V,2}}{\underline{\sigma}^{2} C_{L} \tau}) \|x_{1} - x_{2}\|.$$

$$(135)$$

For the final term, we have

$$T_{4} \leq \|\nabla_{x,\theta}^{2} J_{\tau}(x_{2}, \pi_{\theta_{\tau}^{\star}(x_{2})}) \nabla_{\theta,\theta}^{2} J_{\tau}(x_{2}, \pi_{\theta_{\tau}^{\star}(x_{2})})^{-1} \| \|\nabla_{\theta} f(x_{1}, \pi_{\theta_{\tau}^{\star}(x_{1})}) - \nabla_{\theta} f(x_{2}, \pi_{\theta_{\tau}^{\star}(x_{2})}) \|$$

$$\leq \frac{L_{V}}{\underline{\sigma}} \|\nabla_{\theta} f(x_{1}, \pi_{\theta_{\tau}^{\star}(x_{1})}) - \nabla_{\theta} f(x_{2}, \pi_{\theta_{\tau}^{\star}(x_{2})}) \|$$

$$\leq \frac{L_{V}}{\underline{\sigma}} \cdot L_{f} \Big(\|x_{1} - x_{2}\| + \|\pi_{\tau}^{\star}(x_{1}) - \pi_{\tau}^{\star}(x_{2})\| \Big)$$

$$\leq \Big(\frac{L_{f} L_{V}}{\underline{\sigma}} + \frac{2L_{f} L_{V}^{2}}{\underline{\sigma} C_{L} \tau} \Big) \|x_{1} - x_{2}\|.$$

$$(136)$$

We combine (133)-(136)

$$\begin{split} \|\nabla_x \Phi_\tau(x_1) - \nabla_x \Phi_\tau(x_2)\| &\leq (L_f + \frac{2L_f L_V}{C_L \tau}) \|x_1 - x_2\| + (\frac{L_f L_{V,2}}{\underline{\sigma}} + \frac{2L_f L_V L_{V,2}}{\underline{\sigma} C_L \tau}) \|x_1 - x_2\| \\ &\quad + (\frac{L_f L_V L_{V,2}}{\underline{\sigma}^2} + \frac{2L_f L_V^2 L_{V,2}}{\underline{\sigma}^2 C_L \tau}) \|x_1 - x_2\| + (\frac{L_f L_V}{\underline{\sigma}} + \frac{2L_f L_V^2}{\underline{\sigma} C_L \tau}) \|x_1 - x_2\| \\ &\leq (1 + \frac{2L_V}{C_L \tau}) \Big(\frac{2L_f L_V}{C_L \tau} + \frac{2L_f L_V L_{V,2}}{\underline{\sigma} C_L \tau} + \frac{2L_f L_V^2 L_{V,2}}{\underline{\sigma}^2 C_L \tau} + \frac{2L_f L_V^2}{\underline{\sigma} C_L \tau} \Big) \|x_1 - x_2\|. \end{split}$$

Imposing the step size condition $\tau \leq \frac{2L_V}{C_L}$, we get

$$\|\nabla_x \Phi_{\tau}(x_1) - \nabla_x \Phi_{\tau}(x_2)\| \le \left(\frac{4L_f L_V}{C_L \tau} + \frac{4L_f L_V L_{V,2}}{\underline{\sigma} C_L \tau} + \frac{4L_f L_V^2 L_{V,2}}{\underline{\sigma}^2 C_L \tau} + \frac{4L_f L_V^2}{\underline{\sigma} C_L \tau}\right) \|x_1 - x_2\|$$

$$\le \frac{L_{\Phi}}{\tau} \|x_1 - x_2\|.$$

D.11 Proof of Lemma 11

Let $\theta_{\tau}^{\star}(x)$ denote a parameter representing $\pi_{\tau}^{\star}(x)$ through the softmax function. Define for $\tau > 0$

$$\nabla \Phi_{\tau}(x) = \nabla_{x} f(x, \pi_{\tau}^{\star}(x)) - \nabla_{x,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) \nabla_{\theta,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)})^{-1} \nabla_{\theta} f(x, \pi_{\theta_{\tau}^{\star}(x)}).$$

We consider the following decomposition

$$\|\nabla_x \Phi(x) - \nabla_x \Phi_{w,\tau}(x)\| \le \|\nabla_x \Phi_{\tau}(x) - \nabla_x \Phi_{w,\tau}(x)\| + \|\nabla_x \Phi(x) - \nabla_x \Phi_{\tau}(x)\|. \tag{137}$$

We first bound the first term of (137).

To derive the bound on $\|\nabla_x \Phi_{\tau}(x) - \nabla_x \Phi_{w,\tau}(x)\|$, we take an argument similar to Kwon et al. [2023][Lemma A.2], which we adapt to the case of a non-convex lower level objective. Note that λ

in Kwon et al. [2023] plays the same role as our 1/w. Kwon et al. [2023][Lemma A.2] is still valid without lower level convexity, with the lower bound on $\|\nabla^2_{\theta,\theta}J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)})\|$ changed from the strong convexity coefficient to σ . This allows us to write

$$\begin{split} & \left\| \nabla_{x} \Phi_{\tau}(x) - \nabla_{x} \mathcal{L}_{w,\tau}(x,\pi_{\theta}) + \nabla_{x,\theta}^{2} J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)}) \nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)})^{-1} \nabla_{\theta} \mathcal{L}_{w,\tau}(x,\pi_{\theta}) \right\| \\ & \leq \frac{2L_{V}}{\underline{\sigma}} \|\pi_{\theta} - \pi_{\tau}^{\star}(x)\| \left(L_{f} + \frac{L_{V,2}}{w} \|\pi_{\theta} - \pi_{\tau}^{\star}(x)\| \right). \end{split}$$

Recognizing $\nabla_{\theta} \mathcal{L}_{w,\tau}(x, \pi_{\theta_{w,\tau}^{\star}(x)}) = 0$, we have

$$\|\nabla_{x}\Phi_{\tau}(x) - \nabla_{x}\Phi_{w,\tau}(x)\| = \|\nabla_{x}\Phi_{\tau}(x) - \nabla_{x}\mathcal{L}_{w,\tau}(x, \pi_{w,\tau}^{\star}(x))\|$$

$$\leq \frac{2L_{V}}{\underline{\sigma}} \|\pi_{w,\tau}^{\star}(x) - \pi_{\tau}^{\star}(x)\| \left(L_{f} + \frac{L_{V,2}}{w} \|\pi_{w,\tau}^{\star}(x) - \pi_{\tau}^{\star}(x)\|\right)$$

$$\leq \frac{2L_{V}}{\underline{\sigma}} \cdot \frac{2L_{f}w}{C_{L}\tau} \left(L_{f} + \frac{L_{V,2}}{w} \cdot \frac{2L_{f}w}{C_{L}\tau}\right)$$

$$= \frac{4L_{f}L_{V}w}{C_{L}\underline{\sigma}\tau} (L_{f} + \frac{2L_{f}L_{V,2}}{C_{L}\tau}), \tag{138}$$

where the second inequality follows from Lemma 7.

Next, we bound $\|\nabla_x \Phi_{\tau}(x) - \nabla_x \Phi(x)\|$.

$$\nabla_{x}\Phi_{\tau}(x) - \nabla_{x}\Phi(x) = \underbrace{\nabla_{x}f(x,\pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x}f(x,\pi_{\theta^{\star}(x)})}_{T_{1}} + \underbrace{\nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta^{\star}(x)}) - \nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}_{\tau}(x)})\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}_{\tau}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta^{\star}_{\tau}(x)})}_{T_{2}} + \underbrace{\nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}_{\tau}(x)})\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}_{\tau}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta^{\star}_{\tau}(x)}) - \nabla_{x,\theta}^{2}J_{\tau}(x,\pi_{\theta^{\star}_{\tau}(x)})\nabla_{\theta,\theta}^{2}J_{\tau}(x,\pi_{\theta^{\star}_{\tau}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta^{\star}_{\tau}(x)})}_{T_{3}}.$$

$$(139)$$

To treat T_1 , we have from the Lipschitz continuity of f

$$||T_1|| \le L_f ||\pi_{\tau}^{\star}(x) - \pi^{\star}(x)||.$$

For T_2 ,

$$||T_{2}|| \leq ||\nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)}) - \nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})|||\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1}||||\nabla_{\theta}f(x,\pi_{\theta^{\star}(x)})|||$$

$$+ ||\nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})||||\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1} - \nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1}|||\nabla_{\theta}f(x,\pi_{\theta^{\star}(x)})|||$$

$$+ ||\nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})||||\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1}||||\nabla_{\theta}f(x,\pi_{\theta^{\star}(x)}) - \nabla_{\theta}f(x,\pi_{\theta^{\star}(x)})|||$$

$$\leq L_{V,2}||\pi^{\star}(x) - \pi^{\star}(x)|| \cdot \frac{1}{\underline{\sigma}} \cdot L_{f}$$

$$+ L_{V} \cdot ||\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1}||||\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)}) - \nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})|||||\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1}||$$

$$+ L_{V} \cdot \frac{1}{\underline{\sigma}} \cdot L_{f}||\pi^{\star}_{\tau}(x) - \pi^{\star}(x)||$$

$$\leq \frac{L_{f}(L_{V} + L_{V,2})}{\underline{\sigma}}||\pi^{\star}_{\tau}(x) - \pi^{\star}(x)|| + L_{V} \cdot \frac{1}{\underline{\sigma}} \cdot L_{V,2}||\pi^{\star}_{\tau}(x) - \pi^{\star}(x)|| \cdot \frac{1}{\underline{\sigma}}$$

$$\leq \frac{L_{f}L_{V} + L_{f}L_{V,2} + L_{V}L_{V,2}}{\underline{\sigma}^{2}}||\pi^{\star}_{\tau}(x) - \pi^{\star}(x)||.$$

We then bound T_3 . Note that

$$J_{\tau}(x,\pi) - J(x,\pi) = \frac{\tau}{1-\gamma} \mathbb{E}_{s \sim d_{\rho}^{\pi}} [E(\pi,s)],$$

which is independent of x. This implies $\nabla^2_{x,\theta}J(x,\pi_\theta)=\nabla^2_{x,\theta}J_\tau(x,\pi_\theta)$. Using this relationship, we have

$$||T_{3}|| = \left\| \nabla_{x,\theta}^{2} J(x,\pi_{\theta}) \left(\nabla_{\theta,\theta}^{2} J(x,\pi_{\theta_{\tau}^{\star}(x)})^{-1} \nabla_{\theta} f(x,\pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)})^{-1} \nabla_{\theta} f(x,\pi_{\theta_{\tau}^{\star}(x)}) \right) \right\|$$

$$\leq L_{V,2} L_{f} ||\nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)})^{-1} |||\nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{\theta,\theta}^{2} J(x,\pi_{\theta_{\tau}^{\star}(x)}) |||\nabla_{\theta,\theta}^{2} J(x,\pi_{\theta_{\tau}^{\star}(x)})^{-1} ||$$

$$\leq L_{f} L_{V,2} \cdot \frac{1}{\underline{\sigma}} \cdot \frac{\tau}{1-\gamma} ||\nabla_{\theta,\theta}^{2} \mathbb{E}_{s \sim d_{\rho}^{\pi_{\tau}^{\star}(x)}} [E(\pi_{\tau}^{\star}(x),s)] || \cdot \frac{1}{\underline{\sigma}}$$

$$\leq \frac{L_{f} L_{V,2} (4+8 \log |\mathcal{A}|)\tau}{(1-\gamma)^{4}\underline{\sigma}^{2}},$$

where the third inequality follows from the fact that $\mathbb{E}_{s \sim d_{\rho}^{\pi}}[E(\pi, s)]$ is $\frac{4+8 \log |\mathcal{A}|}{(1-\gamma)^3}$ -Lipschitz (see, for example, Lemma 6 of Zeng et al. [2022a]).

Collecting the bounds on T_1 - T_3 and substituting them into (139),

$$\|\nabla_{x}\Phi_{\tau}(x) - \nabla_{x}\Phi_{w,\tau}(x)\| \leq \|T_{1}\| + \|T_{2}\| + \|T_{3}\|$$

$$\leq \frac{L_{f}(L_{V}+1) + L_{f}L_{V,2} + L_{V}L_{V,2}}{\underline{\sigma}^{2}} \|\pi_{\tau}^{\star}(x) - \pi^{\star}(x)\| + \frac{L_{f}L_{V,2}(4 + 8\log|\mathcal{A}|)\tau}{(1 - \gamma)^{4}\underline{\sigma}^{2}}$$

$$\leq \frac{L_{\star}L_{f}(L_{V}+1) + L_{\star}L_{f}L_{V,2} + L_{\star}L_{V}L_{V,2}}{\underline{\sigma}^{2}} \tau + \frac{L_{f}L_{V,2}(4 + 8\log|\mathcal{A}|)\tau}{(1 - \gamma)^{4}\underline{\sigma}^{2}}$$

$$\leq \frac{L_{\star}L_{f}(L_{V}+1) + L_{\star}L_{f}L_{V,2} + L_{\star}L_{V}L_{V,2} + L_{f}L_{V,2}(4 + 8\log|\mathcal{A}|)}{(1 - \gamma)^{4}\underline{\sigma}^{2}} \tau,$$

$$(140)$$

where the third inequality follows from Lemma 8.

Substituting (138) and (140) into (137),

$$\begin{split} \|\nabla_{x}\Phi(x) - \nabla_{x}\Phi_{w,\tau}(x)\| &\leq \|\nabla_{x}\Phi_{\tau}(x) - \nabla_{x}\Phi_{w,\tau}(x)\| + \|\nabla_{x}\Phi(x) - \nabla_{x}\Phi_{\tau}(x)\| \\ &\leq \frac{4L_{f}L_{V}w}{C_{L}\underline{\sigma}\tau} (L_{f} + \frac{2L_{f}L_{V,2}}{C_{L}\tau}) \\ &+ \frac{L_{\star}L_{f}(L_{V}+1) + L_{\star}L_{f}L_{V,2} + L_{\star}L_{V}L_{V,2} + L_{f}L_{V,2}(4+8\log|\mathcal{A}|)}{(1-\gamma)^{4}\sigma^{2}}\tau. \end{split}$$

D.12 Proof of Lemma 12

Let $\theta^{\star}(x)$ and $\theta^{\star}_{\tau}(x)$ denote one of the softmax parameters that encodes $\pi^{\star}(x)$ and $\pi^{\star}_{\tau}(x)$. Recall the gradient expression $\nabla_x \Phi_{\tau}(x)$ from (9). We can similarly write

$$\nabla_x \Phi(x) = \nabla_x f(x, \pi^*(x)) - \nabla_{x,\pi}^2 J(x, \pi^*(x)) \nabla_{\pi,\pi}^2 J(x, \pi^*(x))^{-1} \nabla_{\pi} f(x, \pi^*(x)). \tag{141}$$

Combining (9) and (141),

$$\|\nabla_{x}\Phi_{\tau}(x) - \nabla_{x}\Phi(x)\|$$

$$\leq \underbrace{\|\nabla_{x}f(x,\pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x}f(x,\pi_{\theta_{\tau}^{\star}(x)})\|}_{T_{1}}$$

$$+ \underbrace{\|\nabla_{x,\theta}^{2}J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)})\nabla_{\theta,\theta}^{2}J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})\nabla_{\theta,\theta}^{2}J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta_{\tau}^{\star}(x)})\|}_{T_{2}}$$

$$+ \underbrace{\|\nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})\nabla_{\theta,\theta}^{2}J_{\tau}(x,\pi_{\theta_{\tau}^{\star}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta_{\tau}^{\star}(x)})\|}_{T_{3}}$$

$$+ \underbrace{\|\nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})\nabla_{\theta,\theta}^{2}J(x,\pi_{\theta^{\star}(x)})^{-1}\nabla_{\theta}f(x,\pi_{\theta^{\star}(x)})\|}_{T_{4}}.$$

$$(142)$$

We treat each term of (142) individually. First, we bound T_1 using the smoothness property of f

$$T_1 \le L_f \|\pi_{\tau}^{\star}(x) - \pi^{\star}(x)\| \le L_{\star} L_f \tau,$$
 (143)

where the second inequality follows from Lemma 8.

We have $\|\nabla_{\theta} f(x, \pi_{\theta})\| \le L_f$ from Assumption 3 and $\|\nabla_{\theta, \theta}^2 J_{\tau}(x, \pi_{\theta_{\tau}^*(x)})^{-1}\| \le \frac{1}{\underline{\sigma}}$ due to Assumption 2. This allows us to bound T_2 as follows

$$T_{2} \leq \|\nabla_{x,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x,\theta}^{2} J(x, \pi_{\theta_{\tau}^{\star}(x)}) \| \|\nabla_{\theta,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)})^{-1} \nabla_{\theta} f(x, \pi_{\theta_{\tau}^{\star}(x)}) \|$$

$$\leq \frac{L_{f}}{\underline{\sigma}} \Big(\|\nabla_{x,\theta}^{2} J_{\tau}(x, \pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x,\theta}^{2} J(x, \pi_{\theta_{\tau}^{\star}(x)}) \| + \|\nabla_{x,\theta}^{2} J(x, \pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x,\theta}^{2} J(x, \pi_{\theta^{\star}(x)}) \| \Big)$$

$$= \frac{L_{f}}{\underline{\sigma}} \|\nabla_{x,\theta}^{2} J(x, \pi_{\theta_{\tau}^{\star}(x)}) - \nabla_{x,\theta}^{2} J(x, \pi_{\theta^{\star}(x)}) \|$$

$$\leq \frac{L_{f}}{\underline{\sigma}} \cdot L_{V,2} \|\pi_{\theta_{\tau}^{\star}(x)} - \pi_{\theta^{\star}(x)} \|$$

$$= \frac{L_{f}}{\underline{\sigma}} \cdot L_{V,2} \|\pi_{\tau}^{\star}(x) - \pi^{\star}(x) \|$$

$$\leq \frac{L_{\star} L_{f} L_{V,2\tau}}{\underline{\sigma}}, \tag{144}$$

where the first equation is due to the fact that $J_{\tau}(x,\pi) - J(x,\pi)$ is independent of x, so the derivative with respect to x is zero. The last inequality plugs in (143).

Similarly, for T_3

$$T_{3} \leq \|\nabla_{x,\theta}^{2} J(x,\pi_{\theta^{*}(x)})\| \|\nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta^{*}_{\tau}(x)})^{-1} - \nabla_{\theta,\theta}^{2} J(x,\pi_{\theta^{*}(x)})^{-1}\| \|\nabla_{\theta} f(x,\pi_{\theta^{*}_{\tau}(x)})\|$$

$$\leq L_{f} L_{V} \|\nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta^{*}_{\tau}(x)})^{-1}\| \|\nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta^{*}_{\tau}(x)}) - \nabla_{\theta,\theta}^{2} J(x,\pi_{\theta^{*}(x)})\| \|\nabla_{\theta,\theta}^{2} J(x,\pi_{\theta^{*}_{\tau}(x)})^{-1}\|$$

$$\leq \frac{L_{f} L_{V}}{\underline{\sigma}^{2}} \Big(\|\nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta^{*}_{\tau}(x)}) - \nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta^{*}(x)})\| + \|\nabla_{\theta,\theta}^{2} J_{\tau}(x,\pi_{\theta^{*}(x)}) - \nabla_{\theta,\theta}^{2} J(x,\pi_{\theta^{*}(x)})\| \Big)$$

$$\leq \frac{L_{f} L_{V}}{\underline{\sigma}^{2}} \Big(L_{V,2} \|\pi_{\tau}^{*}(x) - \pi^{*}(x)\| + L_{V} \tau \Big)$$

$$\leq \frac{L_{x} L_{f} L_{V} L_{V,2} \tau}{\underline{\sigma}^{2}} + \frac{L_{f} L_{V}^{2} \tau}{\underline{\sigma}^{2}},$$

$$(145)$$

where the third inequality again follows from Assumption 2, and the fourth inequality is due to (42) of Lemma 3 (note that $J_{\tau}(x,\pi) - J(x,\pi) = \tau \mathbb{E}_{s \sim d_{\rho}^{\pi}}[E(\pi,s)]$). The last inequality again plugs in (143).

For the final term, we have

$$T_{4} \leq \|\nabla_{x,\theta}^{2} J(x, \pi_{\theta^{\star}(x)}) \nabla_{\theta,\theta}^{2} J(x, \pi_{\theta^{\star}(x)})^{-1} \| \|\nabla_{\theta} f(x, \pi_{\theta^{\star}(x)}) - \nabla_{\theta} f(x, \pi_{\theta^{\star}(x)}) \|$$

$$\leq \frac{L_{V}}{\underline{\sigma}} \|\nabla_{\theta} f(x, \pi_{\theta^{\star}(x)}) - \nabla_{\theta} f(x, \pi_{\theta^{\star}(x)}) \|$$

$$\leq \frac{L_{V}}{\underline{\sigma}} \|\nabla_{\pi} f(x, \pi^{\star}(x)) - \nabla_{\pi} f(x, \pi^{\star}(x)) \|$$

$$\leq \frac{L_{V}}{\underline{\sigma}} \|\pi^{\star}_{\tau}(x) - \pi^{\star}(x) \|$$

$$\leq \frac{L_{\star} L_{V} \tau}{\underline{\sigma}}, \tag{146}$$

where the third inequality is due to the 1-Lipschitz continuity of softmax function, and the fourth inequality is due to Assumption 3.

Combining (143)-(146) leads to

$$\|\nabla_x \Phi_{\tau}(x) - \nabla_x \Phi(x)\| \le L_{\star} L_f \tau + \frac{L_{\star} L_f L_{V,2} \tau}{\sigma} + \frac{L_{\star} L_f L_V L_{V,2} \tau}{\sigma^2} + \frac{L_{\star} L_V^2 \tau}{\sigma} + \frac{L_{\star} L_V \tau}{\sigma}$$

$$\leq \frac{L_{\star}L_{f}L_{V}^{2}L_{V,2}\tau}{\underline{\sigma}^{2}}$$

D.13 Proof of Lemma 13

As $\nabla_x J_{\tau_k}(x_k, \pi_{\theta_k}) - \nabla_x J_{\tau_k}(x_k, \pi_{\theta_{\tau_k}^*}(x_k))$ does not depend on the randomness at iteration k, we have

$$\begin{split} &\mathbb{E}[-\langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) - \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}),x_{k+1} - x_{k}\rangle\rangle] \\ &= \zeta_{k}\mathbb{E}[\langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) - \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}),\mathbb{E}[D_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}},s_{k},a_{k},\bar{s}_{k},\bar{a}_{k},\xi_{k}) \mid \mathcal{F}_{k-1}]\rangle] \\ &= \zeta_{k}\mathbb{E}[\langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) - \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}),\bar{D}_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}})\rangle] \\ &= \zeta_{k}\mathbb{E}[\langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) - \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}),\nabla_{x}\Phi_{w_{k},\tau_{k}}(x_{k})\rangle] \\ &- \zeta_{k}\mathbb{E}[\langle \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{k}}) - \nabla_{x}J_{\tau_{k}}(x_{k},\pi_{\theta_{\tau_{k}}^{\star}(x_{k})}),\bar{D}_{w_{k}}(x_{k},\pi_{\tau_{k}^{\star}}(x_{k}),\pi_{w_{k},\tau_{k}}^{\star}(x_{k})) - \bar{D}_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}})\rangle], \end{split}$$

where the third equation is from (32).

By Young's inequality,

$$\begin{split} &\mathbb{E} \big[- \langle \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}) - \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{\tau_{k}}^{\star}(x_{k})}), x_{k+1} - x_{k} \rangle \rangle \big] \\ &\leq \frac{C_{L}^{2} \alpha_{k} \tau_{k}^{2}}{128 L_{V}^{2}} \mathbb{E} \big[\| \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}) - \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{\tau_{k}}^{\star}(x_{k})}) \|^{2} \big] + \frac{32 L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} \tau_{k}^{2}} \mathbb{E} \big[\varepsilon_{k}^{x} \big] \\ &\quad + \frac{C_{L}^{2} \alpha_{k} \tau_{k}^{2}}{128 L_{V}^{2}} \mathbb{E} \big[\| \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}) - \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{\tau_{k}}^{\star}(x_{k})}) \|^{2} \big] \\ &\quad + \frac{32 L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} \tau_{k}^{2}} \mathbb{E} \big[\| \bar{D}_{w_{k}}(x_{k}, \pi_{\tau_{k}^{\star}}(x_{k}), \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})) - \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}) \|^{2} \big] \\ &\quad \leq \frac{C_{L}^{2} \alpha_{k} \tau_{k}^{2}}{64 L_{V}^{2}} \mathbb{E} \big[\| \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}) - \nabla_{x} J_{\tau_{k}}(x_{k}, \pi_{\theta_{\tau_{k}}^{\star}(x_{k})}) \|^{2} \big] + \frac{32 L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E} \big[\| \pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k}) \|^{2} + \| \pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k}) \|^{2} \big] \\ &\quad \leq \frac{C_{L}^{2} \alpha_{k} w_{k}^{2}}{64} \mathbb{E} \big[\| \pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k}) \|^{2} \big] \\ &\quad + \frac{64 L_{D}^{2} L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E} \big[\| \pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k}) \|^{2} \big] \\ &\quad + \frac{64 L_{D}^{2} L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E} \big[\| \pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k}) \|^{2} \big] \\ &\quad + \frac{64 L_{D}^{2} L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E} \big[\| \pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k}) \|^{2} \big] \\ &\quad + \frac{64 L_{D}^{2} L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E} \big[\| \pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k}) \|^{2} \big] \\ &\quad + \frac{64 L_{D}^{2} L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E} \big[\| \pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k}) \|^{2} \big] \\ &\quad + \frac{64 L_{D}^{2} L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E} \big[\| \pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k}) \|^{2} \big] \\ &\quad + \frac{64 L_{D}^{2} L_{V}^{2} \zeta_{k}^{2}}{C_{L}^{2} \alpha_{k} w_{k}^{2} \tau_{k}^{2}} \mathbb{E} \big[\| \pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k}) \|^{2} \big] \\ &\quad + \frac{64 L_{D}^{2} L_{V}^{2} \zeta_{k}^{$$

where the second inequality employs the Lipschitz continuity of \bar{D}_{w_k} established in Lemma 6.

D.14 Proof of Lemma 14

As $\nabla_x \mathcal{L}^{\text{reweight}}_{w_k, \tau_k}(x_k, \pi_{\theta_k^{\mathcal{L}}}) - \nabla_x \mathcal{L}^{\text{reweight}}_{w_k, \tau_k}(x_k, \pi_{\theta_{w_k, \tau_k}^{\star}}(x_k))$ does not depend on the randomness at iteration k, we have

$$\begin{split} &\mathbb{E}[\langle \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{w_{k},\tau_{k}}}(x_{k})), x_{k+1} - x_{k} \rangle] \\ &= -\zeta_{k} \mathbb{E}[\langle \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{w_{k},\tau_{k}}}(x_{k})), \mathbb{E}[D_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}_{k}, \xi_{k}) \mid \mathcal{F}_{k-1}] \rangle] \\ &= -\zeta_{k} \mathbb{E}[\langle \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{w_{k},\tau_{k}}}(x_{k})), \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}) \rangle] \\ &= -\zeta_{k} \mathbb{E}[\langle \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{w_{k},\tau_{k}}}(x_{k})), \nabla_{x} \Phi_{w_{k},\tau_{k}}(x_{k}) \rangle] \\ &+ \zeta_{k} \mathbb{E}[\langle \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{x} \mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k}, \pi_{\theta_{w_{k},\tau_{k}}}(x_{k})), \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}) \rangle], \\ &\bar{D}_{w_{k}}(x_{k}, \pi_{\tau_{k}}^{\mathcal{L}}(x_{k}), \pi_{w_{k},\tau_{k}}^{\mathcal{L}}(x_{k})) - \bar{D}_{w_{k}}(x_{k}, \pi_{k}, \pi_{k}^{\mathcal{L}}) \rangle], \end{split}$$

where the third equation is from (32).

By Young's inequality,

$$\begin{split} &\mathbb{E}[\langle \nabla_{x}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\mathcal{L}}}) - \nabla_{x}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{w_{k},\tau_{k}}^{\star}}(x_{k})),x_{k+1} - x_{k}\rangle] \\ &\leq \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{128L_{L}^{2}}\mathbb{E}[\|\nabla_{x}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\star}}) - \nabla_{x}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{w_{k},\tau_{k}}^{\star}}(x_{k}))\|^{2}] + \frac{32L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\|\nabla_{x}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\star}}) - \nabla_{x}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{w_{k},\tau_{k}}^{\star}}(x_{k}))\|^{2}] \\ &+ \frac{32L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\|\bar{D}_{w_{k}}(x_{k},\pi_{\tau_{k}}^{\star}(x_{k}),\pi_{w_{k},\tau_{k}}^{\star}(x_{k})) - \bar{D}_{w_{k}}(x_{k},\pi_{k},\pi_{k}^{\mathcal{L}})\|^{2}] \\ &\leq \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64L_{L}^{2}}\mathbb{E}[\|\nabla_{x}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{k}^{\star}}) - \nabla_{x}\mathcal{L}_{w_{k},\tau_{k}}^{\text{reweight}}(x_{k},\pi_{\theta_{w_{k},\tau_{k}}^{\star}}(x_{k}))\|^{2}] + \frac{32L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}] \\ &+ \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &\leq \frac{C_{L}^{2}\alpha_{k}\tau_{k}^{2}}{64}\mathbb{E}[\|\pi_{k}^{\mathcal{L}} - \pi_{w_{k},\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2}] \\ &+ \frac{64L_{D}^{2}L_{L}^{2}\zeta_{k}^{2}}{C_{L}^{2}\alpha_{k}w_{k}^{2}\tau_{k}^{2}}\mathbb{E}[\varepsilon_{k}^{x}], \end{split}$$

where the second inequality employs the Lipschitz continuity of \bar{D}_{w_k} established in Lemma 6.

D.15 Proof of Lemma 15

Within the proof of this lemma, we employ the shorthand notation $z_k = [x_k, \theta_k, \tau_k], \ell(z_k) = V_{\tau_k}^{x_k, \pi_{\theta_k}}$, and

$$y_k = \hat{V}_k - V_{\tau_k}^{x_k, \pi_{\theta_k}} + \beta_k \bar{G}_{\tau_k}(x_k, \theta_k, \hat{V}_k).$$

As $V_{\tau}^{x,\pi_{\theta}}$ is smooth in x,θ,τ , we have from the mean-value theorem that there exists $z_{k+1}^m=mz_k+(1-m)z_{k+1}$ for some scalar $m\in[0,1]$ such that

$$\ell(z_{k}) - \ell(z_{k+1}) \\
= \nabla_{z}\ell(z_{k+1}^{m})^{\top} \left(z_{k} - z_{k+1}\right) \\
= \left(\nabla_{x}V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}^{m}}\right)^{\top} \left(x_{k} - x_{k+1}\right) + \left(\nabla_{\theta}V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}^{m}}\right)^{\top} \left(\theta_{k} - \theta_{k+1}\right) \\
+ \left(\nabla_{\tau}V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}^{m}}\right)^{\top} \left(\tau_{k} - \tau_{k+1}\right) \\
= \zeta_{k} \left(\nabla_{x}V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}^{m}}\right)^{\top} \bar{D}_{w_{k}}(x_{k},\pi_{\theta_{k}},\pi_{\theta_{k}^{\mathcal{L}}}) \\
+ \zeta_{k} \left(\nabla_{x}V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}^{m}}\right)^{\top} \left(D_{w_{k}}(x_{k},\pi_{\theta_{k}},\pi_{\theta_{k}^{\mathcal{L}}},s_{k},a_{k},\bar{s}_{k},\bar{a}_{k},\xi_{k}) - \bar{D}_{w_{k}}(x_{k},\pi_{\theta_{k}},\pi_{\theta_{k}^{\mathcal{L}}}) \\
+ \alpha_{k} \left(\nabla_{\theta}V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}^{m}}\right)^{\top} \bar{F}_{0,\tau_{k}}(x_{k},\theta_{k},\hat{V}_{k}) \\
+ \alpha_{k} \left(\nabla_{\theta}V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}^{m}}\right)^{\top} \left(F_{0,\tau_{k}}(x_{k},\theta_{k},\hat{V}_{k},s_{k},a_{k},s_{k}') - \bar{F}_{0,\tau_{k}}(x_{k},\theta_{k},\hat{V}_{k})\right) \\
+ \left(\nabla_{\tau}V_{\tau_{k+1}}^{x_{k+1},\pi_{\theta_{k+1}}^{m}}\right)^{\top} \left(\tau_{k} - \tau_{k+1}\right), \tag{147}$$

where we denote $x_{k+1}^m = mx_k + (1-m)x_{k+1}, \theta_{k+1}^m = m\theta_k + (1-m)\theta_{k+1}, \tau_{k+1}^m = m\tau_k + (1-m)\tau_{k+1}$.

Plugging (147) into the cross term of interest, we have

$$\langle \hat{V}_k - V_{\tau_k}^{x_k, \pi_{\theta_k}} + \beta_k \bar{G}_{\tau_k}(x_k, \theta_k, \hat{V}_k), V_{\tau_k}^{x_k, \pi_{\theta_k}} - V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}} \rangle$$

$$= \langle y_{k}, \ell(z_{k}) - \ell(x_{k+1}) \rangle$$

$$= \alpha_{k} \langle y_{k}, \left(\nabla_{\theta} V_{\tau_{k+1}^{m}}^{x_{k+1}^{m}, \pi_{\theta_{k+1}^{m}}} \right)^{\top} \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) \rangle$$

$$+ \alpha_{k} \langle y_{k}, \left(\nabla_{\theta} V_{\tau_{k+1}^{m}}^{x_{k+1}^{m}, \pi_{\theta_{k+1}^{m}}} \right)^{\top} \left(F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) \right) \rangle$$

$$+ \zeta_{k} \langle y_{k}, \left(\nabla_{x} V_{\tau_{k+1}^{m}}^{x_{k+1}^{m}, \pi_{\theta_{k+1}^{m}}} \right)^{\top} \bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}}) \rangle$$

$$+ \zeta_{k} \langle y_{k}, \left(\nabla_{x} V_{\tau_{k+1}^{m}, \pi_{\theta_{k+1}^{m}}}^{x_{\theta_{k+1}^{m}}, \pi_{\theta_{k+1}^{m}}} \right)^{\top} \left(D_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}, \xi_{k}) - \bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}}) \right) \rangle$$

$$+ \langle y_{k}, \left(\nabla_{\tau} V_{\tau_{k+1}^{m}}^{x_{k+1}, \pi_{\theta_{k+1}^{m}}} \right)^{\top} \left(\tau_{k} - \tau_{k+1} \right) \rangle. \tag{148}$$

We bound each term of (148) individually. First, by Young's inequality

$$\alpha_{k} \langle y_{k}, \left(\nabla_{\theta} V_{\tau_{k+1}^{m}}^{x_{k+1}, \pi_{\theta_{k+1}^{m}}}\right)^{\top} \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) \rangle
\leq L_{V} \alpha_{k} \|y_{k}\| \|\bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) \|
\leq \frac{(1-\gamma)\beta_{k}}{12} \|y_{k}\|^{2} + \frac{3L_{V}^{2}\alpha_{k}^{2}}{(1-\gamma)\beta_{k}} \|\bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\|^{2}
\leq \frac{(1-\gamma)\beta_{k}}{12} \|y_{k}\|^{2} + \frac{6L_{V}^{2}\alpha_{k}^{2}}{(1-\gamma)\beta_{k}} \|\bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}) \|^{2}
+ \frac{6L_{V}^{2}\alpha_{k}^{2}}{(1-\gamma)\beta_{k}} \|\bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}) \|^{2}
= \frac{(1-\gamma)\beta_{k}}{12} \|y_{k}\|^{2} + \frac{6L_{V}^{2}\alpha_{k}^{2}}{(1-\gamma)\beta_{k}} \|\nabla_{\theta} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}, V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}) \|^{2}
+ \frac{6L_{V}^{2}\alpha_{k}^{2}}{(1-\gamma)\beta_{k}} \|\bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}) - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}) \|^{2}
\leq \frac{(1-\gamma)\beta_{k}}{12} \|y_{k}\|^{2} + \frac{6L_{V}^{2}\alpha_{k}^{2}}{(1-\gamma)\beta_{k}} \|\nabla_{\theta} J_{\tau_{k}}(x_{k}, \pi_{\theta_{k}}) \|^{2} + \frac{6L_{V}^{2}L_{F}^{2}\alpha_{k}^{2}}{(1-\gamma)\beta_{k}} \varepsilon_{k}^{V}, \tag{149}$$

where the equation follows from the fact that $\bar{F}_{0,\tau}(x,\theta,V_{\tau}^{x,\pi_{\theta}}) = \nabla_{\theta}J_{\tau}(x,\pi_{\theta})$ for any x,θ,τ (see (33)), and the final inequality is due to the Lipschitz continuity of $\bar{F}_{0,\tau_{k}}$.

For the second term of (148), we take the expectation

$$\alpha_{k}\mathbb{E}\left[\langle y_{k}, \left(\nabla_{\theta}V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}^{m}}\right)^{\top} \left(F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\right)\rangle\right]$$

$$= \alpha_{k}\mathbb{E}[\langle y_{k}, \left(\nabla_{\theta}V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}^{m}} - \nabla_{\theta}V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}\right)^{\top} \left(F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\right)]$$

$$+ \alpha_{k}\mathbb{E}[\langle y_{k}, \left(\nabla_{\theta}V_{\tau_{k}}^{x_{k+1}, \pi_{\theta_{k+1}}^{m}} - \nabla_{\theta}V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}\right)^{\top} \left(F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\right)]$$

$$= \alpha_{k}\mathbb{E}[\langle y_{k}, \left(\nabla_{\theta}V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}^{m}} - \nabla_{\theta}V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}\right)^{\top} \left(F_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k}, s_{k}, a_{k}, s_{k}') - \bar{F}_{0,\tau_{k}}(x_{k}, \theta_{k}, \hat{V}_{k})\right)]$$

$$\leq 2B_{F}\alpha_{k}\mathbb{E}[\|y_{k}\|\|\nabla_{\theta}V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}^{m}} - \nabla_{\theta}V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}\|]]$$

$$\leq 2B_{F}L_{V}\alpha_{k}\mathbb{E}[\|y_{k}\|\left(\|\pi_{\theta_{k+1}}^{m} - \pi_{\theta_{k}}\| + \|x_{k+1}^{m} - x_{k}\| + |\tau_{k+1}^{m} - \tau_{k}|\right)]$$

$$\leq 2B_{F}L_{V}\alpha_{k}\mathbb{E}[\|y_{k}\|]\left(\alpha_{k}B_{F} + \frac{\zeta_{k}B_{D}}{w_{k}} + \frac{8\tau_{k}}{3(k+1)}\right)$$

$$\leq 2B_{F}L_{V}\alpha_{k}\mathbb{E}[\|y_{k}\|] \cdot \frac{3\tau_{0}B_{F}\alpha_{k}}{\alpha_{0}}$$

$$\leq \frac{3B_{F}^{2}L_{V}\tau_{0}\alpha_{k}^{2}}{\alpha_{0}}\mathbb{E}[\|y_{k}\|^{2}] + \frac{3B_{F}^{2}L_{V}\tau_{0}\alpha_{k}^{2}}{\alpha_{0}},$$
(150)

where the fifth inequality follows from the step size conditions $\zeta_k \leq \frac{B_F \alpha_k w_k}{B_D}$ and $\alpha_0 \leq \min\{\tau_0, \frac{3B_F}{8}\}$, and the second equation follows from

$$\mathbb{E}[\langle y_k, \left(\nabla_{\theta} V_{\tau_k}^{x_k, \pi_{\theta_k}}\right)^{\top} \left(F_{0, \tau_k}(x_k, \theta_k, \hat{V}_k, s_k, a_k, s'_k) - \bar{F}_{0, \tau_k}(x_k, \theta_k, \hat{V}_k)\right)]$$

$$= \mathbb{E}[\langle y_k, \left(V_{\tau_k}^{x_k, \pi_{\theta_k}}\right)^{\top} \mathbb{E}\left[F(\theta_k, \omega_k, \hat{\mu}_k, \hat{V}_{f,k}, s_k, a_k, b_k, s'_k) - \bar{F}(\theta_k, \omega_k, \hat{\mu}_k, \hat{V}_{f,k}) \mid \mathcal{F}_{k-1}\right]]$$

$$= 0.$$

The third term of (148) can be bounded similar to the first term,

$$\zeta_{k}\langle y_{k}, \left(\nabla_{x}V_{\tau_{k+1}}^{x_{m+1}^{*}, \pi_{\theta_{k+1}^{*}}}\right)^{\top} \bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}})\rangle
\leq L_{V}\zeta_{k} \|y_{k}\| \|\bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}})\|
\leq \frac{(1-\gamma)\beta_{k}}{12} \|y_{k}\|^{2} + \frac{3L_{V}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}} \|\bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}})\|^{2}
\leq \frac{(1-\gamma)\beta_{k}}{12} \|y_{k}\|^{2} + \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}} \|\bar{D}_{w_{k}}(x_{k}, \pi_{\tau_{k}}^{\star}(x_{k}), \pi_{w_{k}, \tau_{k}}^{\star}(x_{k}))\|^{2}
+ \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}} \|\bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}}) - \bar{D}_{w_{k}}(x_{k}, \pi_{\tau_{k}}^{\star}(x_{k}), \pi_{w_{k}, \tau_{k}}^{\star}(x_{k}))\|^{2}
= \frac{(1-\gamma)\beta_{k}}{12} \|y_{k}\|^{2} + \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}} \varepsilon_{k}^{x} + \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}} \|\bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}}) - \bar{D}_{w_{k}}(x_{k}, \pi_{\tau_{k}}^{\star}(x_{k}), \pi_{w_{k}, \tau_{k}}^{\star}(x_{k}))\|^{2}
\leq \frac{(1-\gamma)\beta_{k}}{12} \|y_{k}\|^{2} + \frac{6L_{V}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}} \varepsilon_{k}^{x}
+ \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}} \|\pi_{k} - \pi_{\tau_{k}}^{\star}(x_{k})\|^{2} + \frac{12L_{V}^{2}L_{D}^{2}\zeta_{k}^{2}}{(1-\gamma)\beta_{k}} \|\pi_{k}^{\mathcal{L}} - \pi_{w_{k}, \tau_{k}}^{\star}(x_{k})\|^{2}, \tag{151}$$

where the equation is due to the condition in (32).

For the fourth term of (148), we again take the expectation and use the technique in (150)

$$\zeta_{k}\mathbb{E}[\langle y_{k}, \left(\nabla_{x}V_{\tau_{k+1}}^{x_{k+1}}\pi_{\theta_{k+1}}^{m_{k}}\right)^{\top} \left(D_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}, \xi_{k}) - \bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}})\right)\rangle] \\
= \zeta_{k}\mathbb{E}[\langle y_{k}, \left(\nabla_{x}V_{\tau_{k+1}}^{x_{m+1}}\pi_{\theta_{k+1}}^{m_{k}} - \nabla_{x}V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}\right)^{\top} \left(D_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}, \xi_{k}) - \bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}})\right)] \\
+ \zeta_{k}\mathbb{E}[\langle y_{k}, \left(\nabla_{x}V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}\right)^{\top} \left(D_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}, \xi_{k}) - \bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}})\right)] \\
= \zeta_{k}\mathbb{E}[\langle y_{k}, \left(\nabla_{x}V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}} - \nabla_{x}V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}\right)^{\top} \left(D_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}}, s_{k}, a_{k}, \bar{s}_{k}, \bar{a}, \xi_{k}) - \bar{D}_{w_{k}}(x_{k}, \pi_{\theta_{k}}, \pi_{\theta_{k}^{\mathcal{L}}})\right)] \\
\leq 2B_{D}\zeta_{k}\mathbb{E}[\|y_{k}\|\|\nabla_{x}V_{\tau_{k+1}}^{x_{m+1}, \pi_{\theta_{k+1}}} - \nabla_{x}V_{\tau_{k}}^{x_{k}, \pi_{\theta_{k}}}\|]] \\
\leq 2B_{D}L_{V}\zeta_{k}\mathbb{E}[\|y_{k}\|\left(\|\pi_{\theta_{k+1}} - \pi_{\theta_{k}}\| + \|x_{k+1} - x_{k}\| + |\tau_{k+1} - \tau_{k}|\right)] \\
\leq 2B_{D}L_{V}\zeta_{k}\mathbb{E}[\|y_{k}\|\left(\|\pi_{\theta_{k+1}} - \pi_{\theta_{k}}\| + \|x_{k+1} - x_{k}\| + |\tau_{k+1} - \tau_{k}|\right)] \\
\leq 2B_{D}L_{V}\zeta_{k}\mathbb{E}[\|y_{k}\|\right] \cdot \frac{3\tau_{0}B_{F}\alpha_{k}}{\alpha_{0}} \\
\leq \frac{3B_{F}^{2}L_{V}\tau_{0}\alpha_{k}^{2}}{\alpha_{0}}\mathbb{E}[\|y_{k}\|^{2}] + \frac{3B_{F}^{2}L_{V}\tau_{0}\alpha_{k}^{2}}{\alpha_{0}}, \qquad (152)$$

where the last inequality follows from the step size condition $\frac{\zeta_k}{\alpha_k} \leq \frac{B_F}{B_D}$.

For the fifth term of (148), we apply Lemma 2

$$\langle y_{k}, \left(\nabla_{\tau} V_{\tau_{k+1}^{m}}^{x_{k+1}^{m}, \pi_{\theta_{k}^{m}}}\right)^{\top} \left(\tau_{k} - \tau_{k+1}\right) \rangle \leq L_{V} \|y_{k}\| |\tau_{k} - \tau_{k+1}|$$

$$\leq L_{V} \|y_{k}\| \cdot \frac{8\tau_{k}}{3(k+1)}$$

$$\leq \frac{(1-\gamma)\beta_{k}}{6} \|y_{k}\|^{2} + \frac{32L_{V}^{2}\tau_{k}^{2}}{3(1-\gamma)\beta_{k}(k+1)^{2}}.$$
(153)

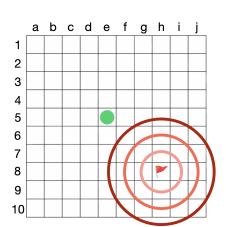
Collecting the results in (149)-(153) and substituting into (148),

$$\begin{split} &\mathbb{E}[\langle \hat{V}_k - V_{\tau_k}^{x_k, \pi_{\theta_k}} + \beta_k \bar{G}_{\tau_k}(x_k, \theta_k, \hat{V}_k), V_{\tau_k}^{x_k, \pi_{\theta_k}} - V_{\tau_{k+1}}^{x_{k+1}, \pi_{\theta_{k+1}}} \rangle] \\ &\leq \frac{(1 - \gamma)\beta_k}{12} \mathbb{E}[\|y_k\|^2] + \frac{6L_V^2 \alpha_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\nabla_{\theta} J_{\tau_k}(x_k, \pi_{\theta_k})\|^2] + \frac{6L_V^2 L_F^2 \alpha_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\varepsilon_k^V] \\ &\quad + \frac{3B_F^2 L_V \tau_0 \alpha_k^2}{\alpha_0} \mathbb{E}[\|y_k\|^2] + \frac{3B_F^2 L_V \tau_0 \alpha_k^2}{\alpha_0} \\ &\quad + \frac{(1 - \gamma)\beta_k}{12} \mathbb{E}[\|y_k\|^2] + \frac{6L_V^2 \zeta_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\varepsilon_k^x] \\ &\quad + \frac{12L_V^2 L_D^2 \zeta_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^*(x_k)\|^2] + \frac{12L_V^2 L_D^2 \zeta_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^*(x_k)\|^2] \\ &\quad + \frac{3B_F^2 L_V \tau_0 \alpha_k^2}{\alpha_0} \mathbb{E}[\|y_k\|^2] + \frac{3B_F^2 L_V \tau_0 \alpha_k^2}{\alpha_0} \\ &\quad + \frac{(1 - \gamma)\beta_k}{6} \mathbb{E}[\|y_k\|^2] + \frac{32L_V^2 \tau_k^2}{3(1 - \gamma)\beta_k (k + 1)^2} \\ &\leq \frac{(1 - \gamma)\beta_k}{3} \mathbb{E}[\|y_k\|^2] + \frac{6B_F^2 L_V \tau_0 \alpha_k^2}{\alpha_0} \mathbb{E}[\|y_k\|^2] + \frac{6L_V^2 \zeta_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\varepsilon_k^X] \\ &\quad + \frac{6L_V^2 \alpha_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\nabla_{\theta} J_{\tau_k}(x_k, \pi_{\theta_k})\|^2] + \frac{6L_V^2 L_F^2 \alpha_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^*(x_k)\|^2] \\ &\quad + \frac{12L_V^2 L_D^2 \zeta_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^*(x_k)\|^2] + \frac{12L_V^2 L_D^2 \zeta_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^*(x_k)\|^2] \\ &\quad + \frac{6B_F^2 L_V \tau_0 \alpha_k^2}{\alpha_0} + \frac{32L_V^2 \tau_k^2}{3(1 - \gamma)\beta_k (k + 1)^2} \\ &\leq \frac{(1 - \gamma)\beta_k}{(1 - \gamma)\beta_k} \mathbb{E}[\|y_k\|^2] + \frac{6L_V^2 \zeta_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\nabla_{\theta} J_{\tau_k}(x_k, \pi_{\theta_k})\|^2] + \frac{6L_V^2 L_F^2 \alpha_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\nabla_{\theta} J_{\tau_k}(x_k, \pi_{\theta_k})\|^2] + \frac{6L_V^2 L_F^2 \alpha_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^*(x_k)\|^2] \\ &\quad + \frac{12L_V^2 L_D^2 \zeta_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k - \pi_{\tau_k}^*(x_k)\|^2] + \frac{12L_V^2 L_D^2 \zeta_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^*(x_k)\|^2] \\ &\quad + \frac{6B_F^2 L_V \tau_0 \alpha_k^2}{\alpha_0} + \frac{32L_V^2 \tau_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^*(x_k)\|^2] \\ &\quad + \frac{6B_F^2 L_V \tau_0 \alpha_k^2}{\alpha_0} + \frac{32L_V^2 \tau_k^2}{(1 - \gamma)\beta_k} \mathbb{E}[\|\pi_k^{\mathcal{L}} - \pi_{w_k, \tau_k}^*(x_k)\|^2] \\ &\quad + \frac{6B_F^2 L_V \tau_0 \alpha_k^2}{\alpha_0} + \frac{32L_V^2 \tau_k^2}{(1 - \gamma)\beta_k (k + 1)^2}, \end{aligned}$$

where the terms are combined in the last inequality under the step size condition $\alpha_k \leq \beta_k$ and $\frac{\alpha_k}{\beta_k} \leq \frac{1-\gamma}{36B_F^2 L_V \tau_0}$.

E Simulation Details

The lower-level MDP is defined on a 10×10 grid, where each state corresponds to a position on the grid. At every state, the agent can choose from four possible actions: $\mathcal{A} = \{\text{UP}, \text{DOWN}, \text{LEFT}, \text{RIGHT}\}$. Each action moves the agent to the adjacent cell in the corresponding direction. If the current position lies on the boundary and the action



would move the agent outside the grid, the state remains unchanged. The upper-level decision variable x sets a goal state for the lower-level problem, shown by the flag in Figure 2. Let $x = (\text{coor}_1(\text{goal}), \text{coor}_2(\text{goal}))$. The reward of state $s = (\text{coor}_1(s), \text{coor}_2(s))$ is

$$r(s) = -\left(\operatorname{coor}_{1}(s) - \operatorname{coor}_{1}(\operatorname{goal})\right)^{2}$$
$$-\left(\operatorname{coor}_{2}(s) - \operatorname{coor}_{2}(\operatorname{goal})\right)^{2}.$$

We choose the upper-level objective f such that $f(x,\pi)$ when x is close to the center of the grid, indicated by green circle in Figure 2 and that π has bias towards DOWN and RIGHT actions. Specifically, with the coordinate of the center cell denoted by $(\text{coor}_1(\text{center}), \text{coor}_2(\text{center}))$, we consider

$$\begin{split} f(x,\pi) &= \Big(\operatorname{coor}_1(\operatorname{goal}) - \operatorname{coor}_1(\operatorname{center}) \Big)^2 \\ &+ \Big(\operatorname{coor}_2(\operatorname{goal}) - \operatorname{coor}_2(\operatorname{center}) \Big)^2 \\ &- \lambda \sum_s \Big(\pi(\operatorname{DOWN} \mid s) + \pi(\operatorname{RIGHT} \mid s) \Big), \end{split}$$

where λ is a weight parameter.

By setting λ sufficiently large, the optimal solution to the bi-level problem is to set the goal on the bottom right corner. This is indeed the learned solution from Algorithm 1.