

# WASSERSTEIN PROXIMAL POLICY GRADIENT: IMPLICIT POLICIES, ENTROPY REGULARIZATION AND LINEAR CONVERGENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We revisit Wasserstein Proximal Policy Gradient (WPPG) for continuous control in infinite-horizon discounted reinforcement learning. By projecting the iterate of Wasserstein proximal gradient onto a parametric policy family with respect to the Wasserstein distance, we derive a new WPPG update that eliminates the need for policy densities or score functions. This makes our method directly applicable to implicit stochastic policies. We prove a linear convergence rate for the WPPG iterate under entropy regularization and [Talagrand’s transport-entropy inequality](#) on the policy class, for both exact and approximate value function estimates. Empirically, our algorithm is simple to implement and achieves competitive performance on standard benchmarks.

## 1 INTRODUCTION

Reinforcement learning (RL) has become a powerful paradigm for solving complex sequential decision-making problems, powering landmark achievements from superhuman performance in strategic games (Silver et al., 2016; 2017) and advanced robotic control (Levine et al., 2016) to the training of Large Language Models (Guo et al., 2025). At the heart of many of these successes are policy gradient (PG) methods (Williams, 1992; Sutton et al., 1999), which iteratively update a parameterized policy to maximize expected rewards.

The geometry underlying policy updates plays an important role in the learning process. Standard policy gradient methods use the Euclidean geometry of parameter space, while the natural policy gradient (Kakade, 2001) and trust-region methods such as TRPO (Schulman et al., 2017a) and PPO (Schulman et al., 2017b) instead exploit the information geometry of policies via Kullback–Leibler (KL) divergence. These methods are supported by a growing body of analysis, with recent results establishing fast global convergence rates in finite action spaces (Agarwal et al., 2020; Lan, 2021; Xiao, 2022; Cen et al., 2022; Bhandari & Russo, 2024).

Recent work explores an alternative paradigm that formulates policy optimization in distribution space under the Wasserstein metric. This perspective builds on the theory of gradient flows in probability spaces (Zhang et al., 2018; Moskovitz et al., 2021; Ziesche & Rozo, 2023; Pfau et al., 2025). In contrast to KL-based methods, which treat actions as independent categories, Wasserstein-based approaches inherently respect the geometry of the action space, capturing meaningful notions of proximity between actions (Pacchiano et al., 2020; Moskovitz et al., 2021; Song et al., 2023). The resulting stochastic policy updates rely on the gradient of the action-value function with respect to the action, drawing a close connection to deterministic policy gradients (Pfau et al., 2025).

While Wasserstein policy optimization offers a compelling alternative to KL-based methods, its theoretical foundations are far less developed. Using the JKO framework (Jordan et al., 1998) for Wasserstein gradient flows, Zhang et al. (2018) established asymptotic convergence of the entropy-regularized problem when the policy distribution is approximated by particles. For finite action spaces, Song et al. (2023) proved global convergence as the Wasserstein penalty coefficient vanishes. Beyond these results, however, convergence guarantees in the more general setting of continuous action spaces—particularly for parametric policies beyond particle approximations (e.g., mixtures of Gaussians)—remain, to the best of our knowledge, an open question.

In this paper, we introduce a version of the Wasserstein policy optimization, which we term Wasserstein Proximal Policy Gradient (WPPG). Our approach projects the theoretical Wasserstein gradient flow onto the parametric policy family, using the Wasserstein metric as the projection criterion. This

contrasts with Pfau et al. (2025), which effectively relies on the KL divergence to project the flow onto the parametric policy manifold. Our main findings are as follows.

- Our WPPG update introduces a new scheme for optimizing stochastic policies. Unlike most existing approaches, it relies solely on the gradient of the action–value function with respect to the action, without requiring access to the policy distribution’s (log-)density or score function. This enables a novel approach to policy optimization with *implicit policies* (Tang & Agrawal, 2018). Empirically, the resulting algorithm is simple to implement and demonstrates competitive performance on standard continuous-control benchmarks.
- We establish a linear convergence rate of WPPG for the entropy-regularized problem, under some regularity and approximate realizability assumptions that can be satisfied by certain implicit policy classes. Our analysis applies to both the exact and approximate value function estimation.

## 1.1 RELATED WORKS

**On Wasserstein policy update** Zhang et al. (2018) formulate continuous-time entropy-regularized policy optimization as a gradient flow of an energy functional, and derive its discrete-time counterpart via the JKO scheme (Jordan et al., 1998), parameterizing policies with particles or energy-based models. Related gradient-flow perspectives also appear in Richemond & Maginnis (2018), and in Ziesche & Rozo (2023) for mixtures of Gaussian policies. Moskovitz et al. (2021) introduce a Wasserstein trust-region policy update and develop efficient kernel-based estimators. Song et al. (2023) study Wasserstein and Sinkhorn trust-region updates in finite action spaces, with an extension to one-dimensional continuous control, and derive closed-form policy updates via duality. More recently, Pfau et al. (2025) propose a Wasserstein gradient flow–inspired update by projecting the Wasserstein gradient flow on parametric manifold using KL divergence.

Most of the above-mentioned Wasserstein policy updates rely on the (log-)density of the policy distribution (or probability mass function for finite action spaces) and/or its score functions, with the exception of Moskovitz et al. (2021), whose kernel-based method instead requires gradients of the kernel for implicit models. In particular, our update, like those of Zhang et al. (2018); Pfau et al. (2025), depends on the gradient of the action-value function with respect to the action. However, unlike both of these two approaches, our method does not rely on the (log-)density of the policy distribution. This is achieved by handling the entropy term through Gaussian noise injection, rather than working directly with the density as in Zhang et al. (2018), and by projecting under the Wasserstein metric instead of the KL divergence, as in Pfau et al. (2025).

We remark that other related approaches leverage Wasserstein geometry in different ways. For instance, Pacchiano et al. (2020) compare policies in latent behavior spaces but do not focus on explicit policy updates, while Abdullah et al. (2019) use Wasserstein distances for robust uncertainty modeling, treating Wasserstein as a constraint on the transition dynamics of the environment rather than as a tool for defining learning dynamics.

**On convergence analysis** Zhang et al. (2018) establish asymptotic convergence of Wasserstein policy optimization based on the JKO scheme when policies are approximated by particles. Song et al. (2023) proves linear convergence analysis for Wasserstein trust-region optimization on finite action spaces, but requires vanishing Wasserstein penalty coefficients. Our convergence proof strategy parallels KL-based analyses such as Lan (2021) for mirror policy gradient in finite action spaces. However, instead of relying on KL-specific tools (e.g., the three-point identity), we develop new analyses tailored to Wasserstein geometry and adopt different assumptions on the problem.

## 2 PRELIMINARIES

### 2.1 REINFORCEMENT LEARNING

**Markov decision processes.** We consider an infinite-horizon discounted MDP  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \rho, \gamma)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $\gamma \in (0, 1)$  the discount factor,  $\mathbb{P}(\cdot \mid s, a)$  the transition kernel,  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  the reward function, and  $\rho$  the initial-state distribution. A policy is a Markov kernel  $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$  so that  $\pi(\cdot \mid s)$  is a probability distribution on  $\mathcal{A}$  for each  $s \in \mathcal{S}$ . Our results apply to the case where  $\mathcal{A}$  is a general metric space.

For a policy  $\pi$ , the value and action-value (Q) functions are

$$V^\pi(s) := \mathbb{E}_{\substack{a_t \sim \pi(\cdot|s_t) \\ s_{t+1} \sim P(\cdot|s_t, a_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right],$$

$$Q^\pi(s, a) := \mathbb{E}_{\substack{a_t \sim \pi(\cdot|s_t) \\ s_{t+1} \sim P(\cdot|s_t, a_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right].$$

We have the relationship that  $V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^\pi(a, s)]$ . The advantage function is  $A^\pi(s, a) := Q^\pi(s, a) - V^\pi(s)$ . Given an initial distribution  $\rho$ , the performance (expected return) is  $J_\rho(\pi) := \mathbb{E}_{s_0 \sim \rho} [V^\pi(s_0)]$ . The discounted state visitation probability is defined as  $d_\rho^\pi(s) := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}_\pi(s_t = s \mid s_0 \sim \rho)$ , with  $\sum_s d_\rho^\pi(s) = 1$ .

**Entropy regularization.** Entropy regularization is widely used in reinforcement learning to prevent premature policy collapse and encourage exploration. It smooths the optimization landscape, reduces gradient variance, and promotes robust, generalizable strategies. These benefits make it a standard component in modern algorithms like TRPO, PPO, and SAC. Define the negative entropy of a policy  $\pi$  at state  $s$  as  $H^\pi(s) := \int_{a \in \mathcal{A}} \pi(a|s) \log \pi(a|s) da$ . Define the modified reward  $r_\tau(s, a) := r(s, a) - \tau \log \pi(a|s)$ . Then  $V_\tau^\pi$  is the value of  $\pi$  under  $r_\tau$ . The corresponding *soft-V* and *soft-Q* functions satisfy the soft Bellman recursion

$$V_\tau^\pi(s) := \mathbb{E}^\pi \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) - \tau H^\pi(s_t)) \mid s_0 = s \right], \quad (1)$$

$Q_\tau^\pi(s, a) := r(s, a) - \tau H^\pi(s) + \gamma \mathbb{E} [V_\tau^\pi(s') \mid s, a]$ ,  $V_\tau^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [Q_\tau^\pi(s, a)]$ , where the temperature  $\tau \geq 0$  trades off exploitation and exploration.

**Policy Gradient** Policy gradient methods form a core class of reinforcement learning algorithms, particularly effective in high-dimensional or continuous settings. They optimize parameterized policies  $\pi = \pi_\theta$  directly through gradient ascent on the parameter  $\theta$ . The geometry underlying these updates significantly influences their behavior. Vanilla policy gradient uses Euclidean steps in parameter space, updating parameters via simple gradient ascent; natural policy gradient incorporates information geometry by equipping the parameter space with the Fisher information metric, yielding updates that are invariant to smooth reparameterizations and closely related to trust-region methods; Mirror or proximal policy gradient methods operate directly in the space of probability distributions, updating policies through mirror ascent with respect to a Bregman divergence such as KL.

## 2.2 WASSERSTEIN PROXIMAL GRADIENT

We next introduce a proximal gradient scheme formulated in the 2-Wasserstein geometry. To stay consistent with our reinforcement learning framework, we present it as a maximization problem and emphasize the underlying intuition; for a rigorous mathematical treatment, we refer the reader to Ambrosio et al. (2008); Santambrogio (2015).

Consider the problem of maximizing an entropy-regularized functional

$$\max_q F_0(q) - \tau H(q), \quad (2)$$

where  $H(q) = \mathbb{E}_q[\log q]$ . At iteration  $k$ , consider the Wasserstein proximal gradient update defined as

$$q_{k+1} = \operatorname{argmax}_q \langle \frac{\delta F_0}{\delta q}[q_k], q \rangle - \tau H(q) - \frac{1}{2\eta} W_2^2(q, q_k), \quad (3)$$

where  $\frac{\delta F}{\delta q}[q_k]$  denotes the first variation of  $F$  at  $q_k$ , and  $W_2$  denotes the 2-Wasserstein metric, and we write  $\langle f, \mu \rangle := \int f d\mu$  for the function-measure pairing. For simplicity, we take the transport cost on the action space  $\mathcal{A}$  to be the norm  $\|\cdot - \cdot\|$ , although our results extend to more general cost functions. Analogous to the Euclidean setting, the Wasserstein proximal gradient update can be interpreted as selecting, among all candidates close to the current iterate, the one that most improves the local linearization of the objective.

The continuous-time limit of the Wasserstein proximal gradient scheme above is given by the Fokker-Planck equation

$$\partial_t q_t(a) = -\operatorname{div} \cdot \left( q_t(a) \nabla \frac{\delta F_0}{\delta q}[q_t](a) \right) + \tau \Delta q_t(a),$$

where the diffusion (Laplacian) term arises from the entropy regularization and the fact that

$$-\operatorname{div} \cdot \left( q(a) \nabla \frac{\delta H}{\delta q}(a) \right) = \operatorname{div} \cdot \left( q(a) \nabla \log q(a) \right) = \Delta q(a).$$

Using Euler-Maruyama discretization, the entropy-induced diffusion can be implemented by injecting Gaussian noise at each iteration. Specifically, a sample  $a \sim q_k$  is updated as

$$a \mapsto a + \eta \nabla \frac{\delta F_0}{\delta q}[q_k](a) + \sqrt{2\tau\eta}\xi, \quad \xi \sim \mathcal{N}(0, I). \quad (4)$$

This update corresponds to one step of noisy Wasserstein gradient ascent, where the drift term pushes samples toward regions of higher value of  $\frac{\delta F_0}{\delta q}[q_k](a)$  and the diffusion term encourages exploration with strength controlled by  $\tau$ .

### 3 WASSERSTEIN PROXIMAL POLICY GRADIENT

We now apply the Wasserstein proximal gradient scheme described in Section 2.2 to the reinforcement learning setting. Define  $\tilde{V}_\tau^\pi(s) = V_\tau^\pi(s) + \tau H^\pi(s)$ . We rewrite the soft value function (1) as

$$\tilde{V}_\tau^\pi(s) - \tau H^\pi(s), \quad (5)$$

which mirrors the structure of (2). For each state  $s$ , the first variation of  $\tilde{V}_\tau^\pi(s)$  with respect to the policy  $\pi$  is given by (see Lemma 3 in Appendix):

$$\frac{\delta \tilde{V}_\tau^\pi}{\delta \pi}(a|s) = \frac{1}{1-\gamma} d_{s_0}^\pi(s) Q_\tau^\pi(s, a).$$

Apply the Wasserstein proximal gradient update (3) to problem (5), we obtain the following update rule for the policy:

$$\pi_{k+1}(\cdot|s) = \operatorname{argmax}_{q(\cdot|s) \in \Pi(s)} \langle Q_\tau^{\pi_k}(s, \cdot), q(\cdot|s) \rangle - \tau H^q(s) - \frac{1}{2\eta} W_2^2(q(\cdot|s), \pi_k(\cdot|s)), \quad (\text{WPPG})$$

where  $\Pi(s)$  denotes the policy class over which we optimize. Equivalently, this can be viewed as the solution to the following Wasserstein policy proximal gradient problem with Wasserstein metric weighted by the state visitation distribution:

$$\max_{\pi \in \Pi} \mathbb{E}_{s \sim d_{s_0}^{\pi_k}, a \sim \pi(\cdot|s)} [A_\tau^{\pi_k}(s, a)] - \frac{1}{2\eta} \mathbb{E}_{s \sim d_{s_0}^{\pi_k}} [W_2^2(\pi(\cdot|s), \pi_k(\cdot|s))],$$

where the advantage function  $A_\tau^{\pi_k}(s, a) = Q_\tau^{\pi_k}(s, a) - V_\tau^{\pi_k}(s)$  replaces the action-value function  $Q_\tau^{\pi_k}(s, a)$  since  $V_\tau^{\pi_k}(s)$  does not depend on the optimization variable  $\pi(\cdot|s)$ . When  $\tau = 0$ , this reduces to the formulation considered in Song et al. (2023); Pfau et al. (2025). When the Wasserstein distance is replaced with Bregman’s distance, (WPPG) becomes the mirror descent policy optimization in Lan (2021).

#### 3.1 PRACTICAL ALGORITHM WITH IMPLICIT POLICIES

Suppose now that the (stochastic) policy class is parameterized as an implicit generative model

$$a = g_\theta(s, Z), \quad \theta \in \Theta,$$

where the latent variable  $Z \sim \nu$  is drawn from an easy-to-sample distribution. In this parameterization, using the drift-diffusion update (4), WPPG at state  $s$  corresponds to transporting each action sample  $a = g_{\theta_k}(s, Z)$  to

$$\tilde{a}_s(Z, \xi) := g_{\theta_k}(s, Z) + \eta \nabla_a Q_\tau^{\pi_k}(s, g_{\theta_k}(s, Z)) + \sqrt{2\tau\eta}\xi, \quad \xi \sim \mathcal{N}(0, I).$$

Let  $\tilde{\mu}_s$  be the distribution of  $\tilde{a}_s$ . Since the exact transported distribution generally does not lie in the parametric family  $\{g_\theta(s, \cdot)_{\#}\nu : \theta \in \Theta\}$ , we update the policy parameters by projection under the Wasserstein metric—in contrast to Pfau et al. (2025) in which information geometry is involved in the projection. Specifically, we first project the drift component onto the parametric family

$$\min_{\theta} W_2^2(g_\theta(s, \cdot)_{\#}\nu, \tilde{\mu}_s). \quad (6)$$

When the step size  $\eta$  is small, the shared-latent coupling

$$(g_\theta(s, Z), \tilde{a}_s(Z, \xi)), \quad Z \sim \nu$$

induced by sharing the same latent variable  $Z \sim \nu$ , is first-order optimal for computing the Wasserstein distance (6). This follows from the geometry of optimal transport: the 2-Wasserstein distance between  $q$  and its displacement  $(\operatorname{id} + v)_{\#}q$  along a vector field  $\eta v$  is given by  $W_2^2(q, (\operatorname{id} + \eta v)_{\#}q) = \eta^2 \mathbb{E}_{a \sim q} [\|v(a)\|^2] + o(\eta^2)$  (Ambrosio et al., 2008, Theorem 8.4.6); and the coupling  $(a, a + \eta v(a))$  with  $a \sim q$  achieves this value up to  $o(\eta^2)$ . Therefore, to update  $\theta$ , we can minimize the expected squared distance under the shared-latent coupling

$$\theta_{k+1} = \operatorname{argmin}_{\theta} \mathbb{E}_{s \sim d^{\pi_k}, Z \sim \nu} \left[ \|g_\theta(s, Z) - g_{\theta_k}(s, Z) - \eta \nabla_a Q_\tau^{\pi_k}(s, g_{\theta_k}(s, Z)) - \sqrt{2\tau\eta}\xi\|^2 \right]. \quad (7)$$

This yields a principled and tractable implementation of Wasserstein proximal policy gradient within the implicit policy class, where actions are nudged by the critic’s gradient and simultaneously diffused by Gaussian noise to encourage exploration. In practical implementation, we ignore the entropy term in  $\nabla_a Q_\tau$ , in a way similar to the approximation in SAC Haarnoja et al. (2018). For a rigorous treatment of Wasserstein information geometry, we refer readers to Chen & Li (2020); here, we offer a more intuitive argument. We also note that while Moskovitz et al. (2021) considers Wasserstein geometry, they solve the optimal coupling problem using kernel methods.

#### 4 CONVERGENCE ANALYSIS

This section derives convergence guarantees for (WPPG), considering both exact (Section 4.1) and inexact (Section 4.2) Q-functions. Our analysis follows the roadmap in Lan (2021), but we provide special treatment for the Wasserstein metric in place of the Bregman divergence used therein.

##### 4.1 EXACT Q-FUNCTION

We begin by stating the main assumptions. Let  $\mathcal{P}_2(\mathbb{R}^d)$  denote the space of probability measures on  $\mathbb{R}^d$  with finite second moments, equipped with the 2-Wasserstein metric. A probability distribution  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$  is said to satisfy  $T_2$  transportation-information inequality, if for every  $\nu \in \mathcal{P}_2(\mathbb{R}^d)$  it holds that

$$W_2^2(\nu, \mu) \leq \frac{2}{\lambda} \text{KL}(\nu \parallel \mu). \quad (T_2(\lambda))$$

**Assumption 1.** *For every state  $s$ , the followings hold:*

- (i) (Uniform  $T_2(\lambda)$ ) *There exists a constant  $\lambda > 0$  such that  $\pi(\cdot|s)$  satisfies  $T_2(\lambda)$  for every  $\pi(\cdot|s) \in \Pi(s)$ .*
- (ii) (Boundedness) *The reward is uniformly bounded and the action space  $\mathcal{A}$  is bounded by  $R$ .*
- (iii) (Approximate realizability) *There exists  $\delta > 0$  such that for any iterate  $\pi_k(\cdot|s)$  of (WPPG), there exists a  $\pi_+(\cdot|s) \in \Pi(s)$  such that  $\text{KL}(\pi(\cdot|s) \parallel \pi_+(\cdot|s)) \leq \delta^2$ , where*

$$\pi(\cdot|s) \in \arg\max_{q(\cdot|s) \in \mathcal{P}(\mathcal{A})} \left\{ \langle Q^\pi(s, \cdot), q(\cdot|s) \rangle - \tau H^q(s) - \frac{1}{2\eta_k} W_2^2(q(\cdot|s), \pi_k(\cdot|s)) \right\}. \quad (8)$$

This assumption requires the policy class to be sufficiently expressive and regular. Consider, for example, a class of implicit policies defined in Section 3.1 with  $g_\theta$  being a neural network and  $\nu$  being a standard Gaussian distribution. Then, (i) can be ensured by designing the neural network to be Lipschitz continuous with respect to the latent variable inputs—this is because the standard Gaussian distribution satisfies  $T_2(1)$  and Lipschitz transformations preserve  $T_2$  inequalities (Ledoux, 2001). (iii) requires that the policy class to be sufficiently expressive to approximate the unconstrained Wasserstein proximal policy update (8) up to a small KL error.

**Remark 1.** *The bounded action space can be replaced by the condition that the policy  $\pi(\cdot|s)$  has uniformly bounded second moment; see Lemma 4 and the comment afterwards in Appendix A.*

For a policy  $\pi(\cdot|s)$  with initial state distribution  $\rho$ , we define

$$J_\rho(\pi) := \mathbb{E}_{s_0 \sim \rho} [V_\tau^\pi(s_0)].$$

Recall the definition of the soft value function (1), and by the soft Bellman optimality conditions, there exists an optimal policy  $\pi^*$  such that

$$V_\tau^{\pi^*}(s) \geq V_\tau^\pi(s) \quad \text{for all } s \in \mathcal{S}, \quad \text{and any } \pi.$$

Hence, optimizing  $V^\pi(\cdot)$  state-wise is equivalent to optimizing any strictly positive weighted average of it. In particular, for any probability weights  $\rho \in \mathcal{P}(\mathcal{S})$  with full support,

$$\pi^* \in \arg\max_{\pi} \mathbb{E}_{s \sim \rho} [V^\pi(s)] \quad \text{s.t.} \quad \pi(\cdot|s) \in \mathcal{P}(\mathcal{A}), \quad \forall s \in \mathcal{S}.$$

While the initial distribution  $\rho$  can be chosen arbitrarily, we follow Lan (2021) and set  $\rho = \nu^*$ —the stationary distribution induced by the optimal policy  $\pi^*$ —to ease the proof. Notably, our (WPPG) algorithm designed to optimize the objective (9) does not require access to  $\nu^*$ . We thereby define the objective function as

$$J(\pi) := J_{\nu^*}(\pi) = \mathbb{E}_{s \sim \nu^*} [V^\pi(s)], \quad (9)$$

and our goal is to maximize  $J_{\nu^*}(\pi)$  over all admissible policies:

$$\max_{\pi} J(\pi) \quad \text{s.t.} \quad \pi(\cdot|s) \in \mathcal{P}_{\mathcal{A}}, \quad \forall s \in \mathcal{S}. \quad (10)$$

Our main result in this subsection is as follows.

**Theorem 1** (Linear Convergence). *Suppose Assumption 1 holds and set the step size  $\eta_k = \eta \geq \frac{1}{\gamma\lambda\tau}$ . Then for any  $k \geq 0$ , the iterates of (WPPG) satisfy*

$$J(\pi_*) - J(\pi_k) + \lambda\tau\mathcal{D}(\pi_k, \pi^*) \leq \gamma^k [J(\pi^*) - J(\pi_0) + \lambda\tau\mathcal{D}(\pi_0, \pi^*)] + \mathcal{O}(\delta + \tau)$$

where  $J$  is defined in (9), and  $\mathcal{D}(\pi_k, \pi^*) := \mathbb{E}_{s \sim \nu^*} [\frac{1}{2} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s))]$ .

Consequently, in order to achieve an error of  $\mathcal{O}(\varepsilon + \delta + \tau)$ , the required iteration complexity is

$$\mathcal{O}\left(\frac{1}{1-\gamma} \log \frac{J(\pi^*) - J(\pi_0) + \lambda\tau\mathcal{D}(\pi_0, \pi^*)}{\varepsilon}\right).$$

This result shares some similarities with the linear convergence of mirror descent policy optimization presented in Lan (2021), but there are notable differences. First, our analysis is based on the Wasserstein distance, whereas theirs relies on the KL divergence, leading to distinct constants in the results. This discrepancy arises because the three-point lemma for KL/Bregman divergence does not apply in our setting. Instead, we leverage the geometry of the Wasserstein distance and the transportation-information inequality to establish new inequalities with different constants. From a methodological perspective, our work addresses the continuous action and state setting, explicitly accounting for the approximation error introduced by the implicit policy class. In contrast, their analysis is confined to finite states and finite actions, where explicit probability mass representations of policies enable exact solutions without the need to consider approximation error.

Another related result is Song et al. (2023, Theorem 5). In terms of the result, we focus on the entropy-regularized problem on continuous action spaces, whereas they consider the unregularized case on finite action spaces (with an extension to one-dimensional continuous spaces). The proof techniques are quite different. Their analysis builds on a bounded difference between Wasserstein policy update and policy-iteration using a uniform bound on the Wasserstein distance, but their analysis does not exploit any specific geometric property of the Wasserstein distance. As a result, their method requires an increasing step size schedule (corresponding to a decreasing Lagrangian multiplier  $\beta$  in their paper), which implies that an  $\mathcal{O}(1/\varepsilon)$  step size is required to obtain an  $\mathcal{O}(\varepsilon)$  solution. In contrast, in our method, if we set the policy realizability error  $\delta = 0$  (matching their setting), we can achieve the same accuracy with a constant step size independent of  $\varepsilon$ .

## 4.2 INEXACT Q-FUNCTION

In practice, the exact action-value function  $Q^{\pi_k}$  is rarely available, since computing it requires either full knowledge of the environment dynamics or an infinite number of Monte Carlo samples. Instead, one typically constructs a stochastic estimator  $Q^{\pi_k, \xi_k}$  from finite trajectories, temporal-difference updates, or function approximation.

In this case, the (WPPG) update is defined by substituting the exact value function  $Q^{\pi_k}$  in (29) with its stochastic estimator  $Q^{\pi_k, \xi_k}$ . Formally, the update rule is given by

$$\pi_{k+1}(\cdot|s) \in \operatorname{argmax}_{q(\cdot|s) \in \Pi(s)} \langle Q^{\pi_k, \xi_k}(s, \cdot), q \rangle - \tau H^\pi(s) - \frac{1}{2\eta_k} W_2^2(q, \pi_k(\cdot|s)). \quad (11)$$

Such an estimator inevitably introduces both variance and bias, which can accumulate across iterations and significantly affect policy updates. To ensure that our analysis remains tractable while still capturing realistic scenarios, we impose mild conditions on the stochastic approximation and estimation error.

**Assumption 2.** *For each iteration  $k \geq 0$ , the stochastic estimator  $Q^{\pi_k, \xi_k}$  satisfies*

$$\mathbb{E}_{\xi_k} [Q^{\pi_k, \xi_k}] = \bar{Q}^{\pi_k}, \quad (12)$$

$$\|\bar{Q}^{\pi_k} - Q^{\pi_k}\|_\infty \leq \epsilon_k, \quad (13)$$

$$\mathbb{E}_{\xi_k} \left[ \|\nabla_a Q^{\pi_k, \xi_k} - \nabla_a Q^{\pi_k}\|_{2, \infty}^2 \right] \leq \sigma_k^2. \quad (14)$$

Where  $\|\cdot\|_\infty$  is the uniform norm over  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , and for any function  $f$ ,  $\|\cdot\|_{2, \infty}$  is defined as

$$\|f\|_{2, \infty} := \sup_{(s, a) \in \mathcal{S} \times \mathcal{A}} \|f(s, a)\|_2.$$

This assumption is similar to Lan (2021), except that (14) concerns the action-gradient of the  $Q$ -function, which is needed in our Wasserstein policy update. We have the following convergence result.

**Theorem 2** (Linear Convergence). Suppose Assumptions 1 and 2 hold, and for all  $k \geq 0$ ,  $\epsilon_k \leq \epsilon$ ,  $\sigma_k \leq \sigma$  and  $\|Q^{\pi_k, \xi_k}\|_\infty \leq B$ . Then the iterates of (11) using step size  $\eta_k = \eta \geq \frac{1}{\gamma\lambda\tau}$  satisfies

$$\mathbb{E}_{\xi_{0:k-1}} [J(\pi^*) - J(\pi_k) + \lambda\tau D(\pi_k, \pi^*)] \leq \gamma^k [J(\pi^*) - J(\pi_0) + \lambda\tau D(\pi_0, \pi^*)] + \mathcal{O}(\delta + \tau + \epsilon + \sigma) \quad (15)$$

where  $J$  is defined in (9), and  $\mathcal{D}(\pi_k, \pi^*) := \mathbb{E}_{s \sim \nu^*} [\frac{1}{2} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s))]$ . Consequently, in order to achieve an error of  $\mathcal{O}(\epsilon + \tau + \delta + \epsilon + \sigma)$  in expectation, the required iteration complexity is

$$\mathcal{O}\left(\frac{1}{1-\gamma} \log \frac{J(\pi^*) - J(\pi_0) + \lambda\tau \mathcal{D}(\pi_0, \pi^*)}{\epsilon}\right).$$

Compared with Theorem 1, this result involves additional bias arising from the estimation of the  $Q$ -function and variance from stochastic estimation of the  $Q$ -function. Nonetheless, our analysis shows that the total error does not accumulate across iterations, and the convergence guarantee only incurs an  $\mathcal{O}(\delta + \tau + \epsilon + \sigma)$  term in the final bound, rather than growing with the number of iterations.

## 5 EXPERIMENTS

Our empirical study consists of two parts: comparative evaluation and ablation analysis. The comparative evaluation focuses on benchmarking our methods against representative baselines to assess overall performance. The ablation analysis investigates three key questions: (i) the effect of  $\tau$  on WPPG, (ii) the impact of latent variable dimension on WPPG-I, and (iii) the role of double- $Q$  learning in WPPG, which is postponed to the appendix.

### 5.1 COMPARATIVE EVALUATION

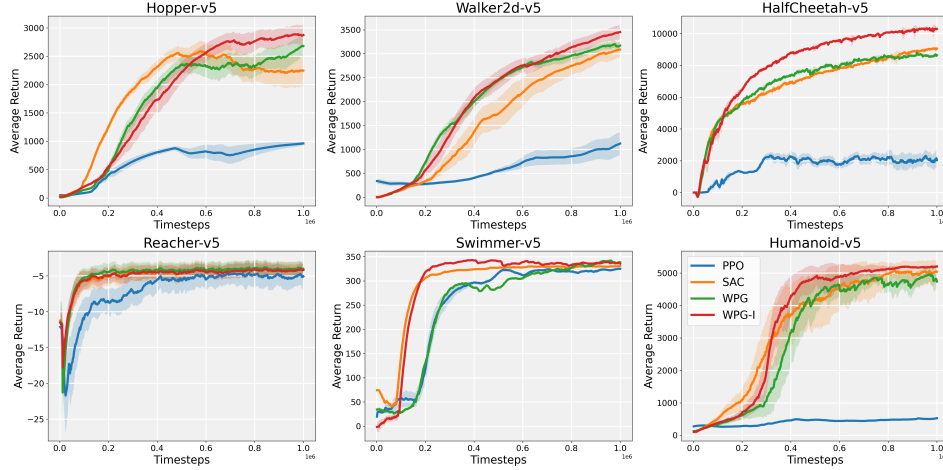


Figure 1: Training curves on MuJoCo continuous control benchmarks: Solid lines denote the mean episodic return, while shaded areas represent the 95% confidence interval computed over 10 independent evaluation runs with different random seeds.

**Evaluation Tasks** We evaluate our approach on a set of standard continuous control benchmarks from the MuJoCo suite<sup>1</sup>, including Hopper-v5, Walker2d-v5, HalfCheetah-v5, Reacher-v5, Swimmer-v5, and Humanoid-v5. These tasks cover a wide range of difficulties: from relatively low-dimensional and easy-to-learn tasks such as Swimmer and Hopper, to high-dimensional and challenging tasks such as Humanoid.

**Baseline Models** We compare against three representative baselines: (i) PPO, a KL-proximal policy optimization method that employs clipped surrogate objectives to constrain successive policy updates and improve training stability (Schulman et al., 2017b); (ii) SAC, a stochastic actor-critic algorithm formulated as an entropy-regularized policy optimization method, which augments the reward with a maximum-entropy term to encourage exploration and improve stability (Haarnoja et al.,

<sup>1</sup><https://gymnasium.farama.org/environments/mujoco/>



2018). (iii) WPO, a Wasserstein-proximal actor-critic algorithm that replaces the KL divergence commonly used in proximal methods with the Wasserstein distance, thereby constraining successive policy updates under the geometry of optimal transport. (Pfau et al., 2025). Our proposed methods include WPPG, with a Gaussian MLP policy actor, and WPPG-I, with an implicit MLP policy actor. Hyperparameters for PPO are taken from the RL Zoo project<sup>2</sup>, while those for SAC and WPO follow Pfau et al. (2025).

**Experiment Setup** For fair comparison, WPPG and WPPG-I adopt the double- $Q$  technique, taking the minimum of two  $Q$ -functions for both critic targets and action gradients, while WPO follows the original single- $Q$  formulation. We also evaluate a single- $Q$  variant of WPPG for consistency, with its comparison to WPO revisited in the ablation study. To further ensure fairness, SAC is evaluated with entropy coefficient self-tuning disabled, since  $\tau$  in WPPG and WPPG-I is also fixed rather than adaptively tuned while training. All off-policy methods share the same replay buffer structure. Additional implementation details for all methods and pseudo code of our algorithm are provided in the Appendix B.

**Results and Discussion** The learning curves for all tasks are shown in Figure 5. Across the six MuJoCo benchmarks, WPPG demonstrates performance comparable to SAC, which can be attributed to the fact that both methods employ Gaussian MLP policies. This observation suggests that Wasserstein geometry can match, and in some cases even surpass, the effectiveness of KL-based geometry in policy optimization. More importantly, WPPG-I consistently outperforms all baselines, achieving superior convergence speed and higher returns across nearly all tasks. The success of WPPG-I further indicates that WPPG can be naturally extended to implicit policy classes: although we use only a simple MLP-based implicit policy here, the framework readily accommodates richer architectures. We refer readers to the Appendix B.1 and B.2 for a detailed comparison between the two algorithms. In contrast, PPO lags behind due to slower learning and lower asymptotic performance, while WPO suffers from unstable convergence on challenging environments such as Humanoid and Swimmer and even fails to learn in Reacher. Overall, these results highlight that WPPG preserves the sample efficiency of off-policy actor-critic methods, while WPPG-I not only inherits these advantages but also demonstrates a consistent and significant margin over all baselines.

## 5.2 ABLATION STUDY

**Ablation on  $\tau$ .** In the preceding analysis, we showed that the parameter  $\tau$  originates from entropy regularization of the policy. Unlike SAC and related methods that explicitly add an entropy penalty term into the  $Q$ -function fitting objective, WPPG does not require such a penalty. Instead, Gaussian noise is injected when computing the movement direction of action samples, where the scale of the injected noise  $\tau$  corresponds to the magnitude of the entropy penalty.

To study the impact of  $\tau$ , we conducted ablation experiments on Humanoid environment. As illustrated in Figure 2, on Humanoid we observe that injecting noise with  $\tau$  in the range  $[0, 0.01]$  significantly accelerates convergence, while larger values 0.1 slows it down. This reflects a clear exploration-exploitation trade-off: noise injection encourages the policy to maintain entropy, thereby enabling exploration of richer reward information, but excessive noise hampers the ability of  $\nabla_a Q(s, a)$  to provide useful guidance for policy updates.

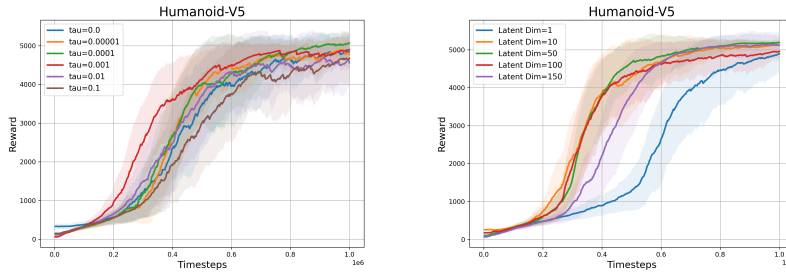


Figure 2: Ablation study on  $\tau$  (left) and *Latent Dimension* (right).

<sup>2</sup><https://github.com/DLR-RM/rl-baselines3-zoo>



**Ablation on Latent Dimension** We further evaluate the effect of the dimension of latent variable  $z$  in our implicit generative model on WPPG-I in the Humanoid environment. When the latent dimension is as small as 1, the model learns slowly due to insufficient stochasticity. With moderate dimensions (e.g., 10, 50, 100), learning is significantly accelerated, indicating that a reasonable amount of latent variables enhances exploration without overwhelming the policy. However, when the latent dimension becomes too large (e.g., 150), excessive non-informative variables begin to dominate the input and degrade learning speed. Empirically, we find that setting the latent dimension about one-third of the state dimension provides a good balance between exploration and stability.

**Ablation on Double Q Trick** We also evaluated the single- $Q$  variant of WPPG across all environments, and found that it outperforms WPO on nearly every task, with the corresponding results provided in the Appendix B.4. In addition, consistent with prior findings, adopting double- $Q$  further improves WPPG by both stabilizing training and enhancing overall performance.

## 6 LIMITATIONS AND CONCLUSION

**Limitations** Our empirical study focuses on standard continuous-control benchmarks, while the performance of WPPG and WPPG-I on ultra high-dimensional and more complex tasks remains unexplored. Moreover, the implicit policy formulation admits natural extensions to richer and more expressive classes, such as diffusion-based policies, which we have not yet investigated. Finally, whether Wasserstein geometry can serve as a foundation for reinforcement learning fine-tuning of large language models (LLMs) is an open and promising direction. We leave these broader applications to future work.

**Conclusion** In this work, we proposed Wasserstein Proximal Policy Gradient (WPPG), a novel framework for policy optimization that leverages Wasserstein geometry to design proximal updates directly in distribution space. Our method eliminates the need for policy densities or score functions, making it naturally applicable to implicit policies. Theoretically, we established linear convergence guarantees under entropy regularization and [a transport-entropy condition](#), covering both exact and approximate value function settings. To the best of our knowledge, this work is an early pioneering attempt to employ Wasserstein geometry for establishing global convergence guarantees. Empirically, WPPG demonstrates competitive performance against strong baselines, while its implicit extension WPPG-I consistently outperforms them across challenging continuous-control benchmarks. These results highlight the potential of Wasserstein geometry as a principled alternative to KL-based methods in reinforcement learning.

**Reproducibility Statement.** Our entire implementation is built on top of the Stable Baselines3 (SB3) interface, ensuring compatibility and transparency. We provide detailed pseudocode for both WPPG and WPPG-I, together with comprehensive hyperparameter specifications, in the Appendix B.1. These materials cover nearly all technical details of our approach, and the pseudocode alone suffices to fully reproduce the proposed models. In addition, our experiments follow standardized training and evaluation protocols across tasks, which further supports reproducibility and comparability with prior baselines.

## REFERENCES

- Mohammed Amin Abdullah, Hang Ren, Haitham Bou Ammar, Vladimir Milenkovic, Rui Luo, Mingtian Zhang, and Jun Wang. Wasserstein robust reinforcement learning. *arXiv preprint arXiv:1907.13196*, 2019.
- Alekh Agarwal, Sham M. Kakade, Jason D. Lee, and Gaurav Mahajan. On the Theory of Policy Gradient Methods: Optimality, Approximation, and Distribution Shift, October 2020. URL <http://arxiv.org/abs/1908.00261>. arXiv:1908.00261 [cs, stat].
- Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient Flows*. Birkhäuser, Basel, 2008. ISBN 978-3-7643-8721-1 978-3-7643-8722-8. doi: 10.1007/978-3-7643-8722-8. URL <http://link.springer.com/10.1007/978-3-7643-8722-8>.
- Jalaj Bhandari and Daniel Russo. Global optimality guarantees for policy gradient methods. *Operations Research*, 72(5):1906–1927, 2024.
- Shicong Cen, Chen Cheng, Yuxin Chen, Yuting Wei, and Yuejie Chi. Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 70(4):2563–2578, 2022.
- Yifan Chen and Wuchen Li. Optimal transport natural gradient for statistical manifolds with continuous sample space. *Information Geometry*, 3(1):1–32, 2020.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, 2025.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. Pmlr, 2018.
- Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the fokker-planck equation. *SIAM journal on mathematical analysis*, 29(1):1–17, 1998.
- Sham M Kakade. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- Guanghui Lan. Policy Mirror Descent for Reinforcement Learning: Linear Convergence, New Sampling Complexity, and Generalized Problem Classes, January 2021. URL <https://arxiv.org/abs/2102.00135v6>.
- Guanghui Lan. Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. *Mathematical programming*, 198(1):1059–1106, 2023.
- Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2001.
- Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(39):1–40, 2016.
- Ted Moskovitz, Michael Arbel, Ferenc Huszar, and Arthur Gretton. Efficient wasserstein natural gradients for reinforcement learning. In *International Conference on Learning Representations*, 2021.
- Aldo Pacchiano, Jack Parker-Holder, Yunhao Tang, Krzysztof Choromanski, Anna Choromanska, and Michael Jordan. Learning to score behaviors for guided policy optimization. In *International Conference on Machine Learning*, pp. 7445–7454. PMLR, 2020.
- David Pfau, Ian Davies, Diana Borsa, Joao G. M. Araujo, Brendan Tracey, and Hado van Hasselt. Wasserstein Policy Optimization, May 2025. URL <http://arxiv.org/abs/2505.00663>. arXiv:2505.00663 [cs].

- Pierre H Richemond and Brendan Maginnis. Diffusing policies: Towards wasserstein policy gradient flows. 2018.
- Filippo Santambrogio. *Optimal Transport for Applied Mathematicians: Calculus of Variations, PDEs, and Modeling*, volume 87 of *Progress in Nonlinear Differential Equations and Their Applications*. Springer International Publishing, Cham, 2015. ISBN 978-3-319-20827-5 978-3-319-20828-2. doi: 10.1007/978-3-319-20828-2. URL <https://link.springer.com/10.1007/978-3-319-20828-2>.
- John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust Region Policy Optimization, April 2017a. URL <http://arxiv.org/abs/1502.05477>. arXiv:1502.05477 [cs].
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms, August 2017b. URL <http://arxiv.org/abs/1707.06347>. arXiv:1707.06347 [cs].
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- Jun Song, Niao He, Lijun Ding, and Chaoyue Zhao. Provably Convergent Policy Optimization via Metric-aware Trust Region Methods, June 2023. URL <http://arxiv.org/abs/2306.14133>. arXiv:2306.14133 [cs].
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- Yunhao Tang and Shipra Agrawal. Implicit policy for reinforcement learning. *arXiv preprint arXiv:1806.06798*, 2018.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- Lin Xiao. On the convergence rates of policy gradient methods. *Journal of Machine Learning Research*, 23(282):1–36, 2022.
- Ruiyi Zhang, Changyou Chen, Chunyuan Li, and Lawrence Carin. Policy optimization as wasserstein gradient flows. In *International Conference on machine learning*, pp. 5737–5746. PMLR, 2018.
- Hanna Ziesche and Leonel Rozo. Wasserstein gradient flows for optimizing gaussian mixture policies. *Advances in Neural Information Processing Systems*, 36:21058–21080, 2023.

**LLM Usage Statement** We used large language models (LLMs) only for language editing and polishing of the manuscript, and for occasional assistance in writing boilerplate code such as plotting scripts, LaTeX formatting, and Stable Baselines3 (SB3) wrapper templates. All research ideas, algorithmic designs, theoretical analyses, and experimental implementations are original and conducted by the authors.

## A THEORETICAL DERIVATIONS AND PROOFS

We begin by showing a fundamental computation of the functional gradient of the value function  $V_\tau^\pi(s_0)$  with respect to  $\pi$ .

**Lemma 1** (Log-derivative of the trajectory law). *Consider a zero-mass perturbation  $\delta\pi(\cdot|s)$  (i.e.  $\int_{\mathcal{A}} \delta\pi(da|s) = 0$  for all  $s$ ) and the policy path  $\pi_\epsilon = \pi + \epsilon\delta\pi$ , with  $|\epsilon|$  small so that  $\pi_\epsilon \geq 0$  and  $\text{supp}(\delta\pi) \subseteq \{\pi > 0\}$ . Let  $\mathbb{P}_{\pi_\epsilon}$  be the path measure of the Markov chain induced by  $\pi_\epsilon$  and the transition kernel  $P(s'|s, a)$  from an initial law  $\rho_0$ . For any finite horizon  $T$  and any integrable test functional  $F$  of  $\tau_{0:T} = (s_0, a_0, \dots, s_T)$ ,*

$$\left. \frac{d}{d\epsilon} \right|_{\epsilon=0} \int F(\tau_{0:T}) d\mathbb{P}_{\pi_\epsilon}(\tau_{0:T}) = \int F(\tau_{0:T}) \left( \sum_{u=0}^{T-1} \frac{\delta\pi(a_u | s_u)}{\pi(a_u | s_u)} \right) d\mathbb{P}_\pi(\tau_{0:T}). \quad (16)$$

*If  $F$  is dominated by an integrable function uniformly, then letting  $T \rightarrow \infty$  and applying dominated convergence yields*

$$\left. \frac{d}{d\epsilon} \right|_{\epsilon=0} \int F(\tau) d\mathbb{P}_{\pi_\epsilon}(\tau) = \int F(\tau) \left( \sum_{u=0}^{\infty} \frac{\delta\pi(a_u | s_u)}{\pi(a_u | s_u)} \right) d\mathbb{P}_\pi(\tau). \quad (17)$$

*Proof.* Write the truncated path density under  $\pi_\epsilon$  as

$$d\mathbb{P}_{\pi_\epsilon}(\tau_{0:T}) = \rho_0(s_0) \prod_{t=0}^{T-1} [\pi_\epsilon(a_t | s_t) P(s_{t+1} | s_t, a_t)].$$

Only  $\pi_\epsilon$  depends on  $\epsilon$ , hence by the product rule,

$$\left. \frac{d}{d\epsilon} \right|_{\epsilon=0} \prod_{t=0}^{T-1} \pi_\epsilon(a_t | s_t) = \left( \prod_{t=0}^{T-1} \pi_\epsilon(a_t | s_t) \right) \sum_{u=0}^{T-1} \frac{\delta\pi(a_u | s_u)}{\pi_\epsilon(a_u | s_u)}.$$

Evaluating at  $\epsilon = 0$  gives

$$\left. \frac{d}{d\epsilon} \right|_{\epsilon=0} \pi_\epsilon(a_u | s_u) = \delta\pi(a_u | s_u),$$

hence (16) follows by dominated convergence under the stated integrability.  $\square$

**Lemma 2** (Functional policy gradient of Entropy). *Fix a state  $s \in \mathcal{S}$ . Let the policy at  $s$  admit a density  $\pi(\cdot|s)$  w.r.t. a reference measure  $\mathfrak{a}$  (and write  $h^\pi(s) := h(\pi(\cdot|s))$ ) with*

$$H^\pi(s) = \int_{\mathcal{A}} \pi(a|s) \log \pi(a|s) \mathfrak{a}.$$

*For any zero-mass direction  $\delta\pi(\cdot|s)$  (i.e.  $\int_{\mathcal{A}} \delta\pi(a|s) \mathfrak{a} = 0$ ) and the perturbation  $\pi_\epsilon(\cdot|s) = \pi(\cdot|s) + \epsilon\delta\pi(\cdot|s)$ , provided  $\pi_\epsilon(\cdot|s)$  stays nonnegative for  $|\epsilon|$  small and  $\pi(a|s) > 0$  on its support, the Gâteaux derivative of  $h^\pi(s)$  in the direction  $\delta\pi(\cdot|s)$  is*

$$\frac{\delta H^\pi(s)}{\delta\pi} := \left. \frac{d}{d\epsilon} \right|_{\epsilon=0} H^{\pi_\epsilon}(s) = \int_{\mathcal{A}} 1 + \log \pi(a|s) \delta\pi(da|s). \quad (18)$$

*Proof.* Let  $\pi_\epsilon(a) = \pi(a|s) + \epsilon\delta\pi(a|s)$ . Then

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \frac{H^{\pi_\epsilon}(s) - H^\pi(s)}{\epsilon} &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \left( \int \pi_\epsilon(a|s) \log \pi_\epsilon(a|s) \mathfrak{a} - \int \pi(a|s) \log \pi(a|s) \mathfrak{a} \right) \\ &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \left( \int (\pi(a|s) + \epsilon\delta\pi(a|s)) \log(\pi(a|s) + \epsilon\delta\pi(a|s)) \mathfrak{a} - \int \pi(a|s) \log \pi(a|s) \mathfrak{a} \right) \\ &= \int \lim_{\epsilon \rightarrow 0} \frac{\pi(a|s) (\ln(\pi(a|s) + \epsilon\delta\pi(a|s)) - \ln \pi(a|s))}{\epsilon} \mathfrak{a} + \delta\pi(a|s) \ln(\pi(a|s) + \epsilon\delta\pi(a|s)) \mathfrak{a} \\ &= \int (1 + \ln \pi) d\delta\pi \end{aligned}$$

which implies that  $\frac{\delta H^\pi(s)}{\delta\pi}(a) = 1 + \ln \pi(a|s)$ .  $\square$

**Lemma 3** (Functional policy gradient under entropy regularization). *Then for any zero-mass direction  $\delta\pi$ , the Gâteaux derivative (functional gradient)*

$$\frac{\delta V_\tau^\pi(s_0)}{\delta\pi}(s, a) = \frac{1}{1-\gamma} d_{s_0}^\pi(s) \left( Q_\tau^\pi(s, a) - \tau(1 + \ln \pi(a|s)) \right), \quad (19)$$

*Proof.* We should be very careful here: when computing the functional gradient of  $V_\tau^\pi(s)$  with respect to  $\pi$ , there are two contributing parts - one from the MDP dynamics and one from the entropy regularization. Let  $V_\tau^\pi(s) := \mathbb{E}_{\pi, s_0=s} \left[ \sum_{t \geq 0} \gamma^t (r(s_t, a_t) - \tau H^\pi(s_t)) \right]$  By Lemma 1 with  $F(\tau) = \sum_{t \geq 0} \gamma^t (r(s_t, a_t) - \tau H^\pi(s_t))$ ,

$$\frac{d}{d\epsilon} \Big|_{\epsilon=0} V_\tau^{\pi_\epsilon} = \sum_{t \geq 0} \gamma^t \mathbb{E}_\pi \left[ (r(s_t, a_t) - \tau H^\pi(s_t)) \sum_{u=0}^t \frac{\delta\pi(a_u|s_u)}{\pi(a_u|s_u)} \right] - \sum_{t \geq 0} \gamma^t \mathbb{E}_\pi \left[ \frac{d}{d\epsilon} \Big|_{\epsilon=0} \tau H^{\pi_\epsilon}(s) \right]$$

(I) *Trajectory-law part.*

$$(I) = \sum_{u \geq 0} \mathbb{E}_\pi \left[ \frac{\delta\pi(a_u|s_u)}{\pi(a_u|s_u)} \sum_{t \geq u} \gamma^t (r(s_t, a_t) - \tau H^\pi(s_t)) \right] = \sum_{u \geq 0} \gamma^u \mathbb{E}_\pi \left[ \frac{\delta\pi(a_u|s_u)}{\pi(a_u|s_u)} Q_\tau^\pi(s_u, a_u) \right].$$

Condition on  $s_u$  and expand over  $a$  to cancel  $\pi$ , then sum over time and use  $\sum_{u \geq 0} \gamma^u \mathbb{P}_\pi(s_u = s | s_0) = \frac{1}{1-\gamma} d_{s_0}^\pi(s)$  to get

$$(I) = \iint_{S \times \mathcal{A}} \frac{1}{1-\gamma} d_{s_0}^\pi(s) Q_\tau^\pi(s, a) \delta\pi(a|s) ds da.$$

(II) *Explicit entropy part.* By Lemma 2 for each state  $s$ ,  $H^\pi(s) = - \int \pi(a|s) \log \pi(a|s) da$ , we have

$$\frac{d}{d\epsilon} \Big|_{\epsilon=0} H^{\pi_\epsilon}(s) = \int_{\mathcal{A}} [1 + \log \pi(a|s)] \delta\pi(a|s) da,$$

hence

$$(II) = \sum_{t \geq 0} \gamma^t \mathbb{E}_\pi \left[ \tau \frac{d}{d\epsilon} \Big|_{\epsilon=0} H^{\pi_\epsilon}(s) \right] = \iint_{S \times \mathcal{A}} \frac{1}{1-\gamma} d_{s_0}^\pi(s) \tau [1 + \log \pi(a|s)] \delta\pi(a|s) ds da.$$

Finally, we have

$$\frac{\delta V_\tau^\pi(s_0)}{\delta\pi}(s, a) = \frac{1}{1-\gamma} d_{s_0}^\pi(s) (Q_\tau^\pi(s, a) - \tau(1 + \ln \pi(a|s)))$$

□

**Lemma 4** (Entropy upper bound from second moment). *Let  $X \in \mathbb{R}^d$  be absolutely continuous with differential entropy  $h(X)$ . Assume*

$$\mathbb{E}\|X\|^2 \leq R^2.$$

*Then*

$$h(X) \leq \frac{d}{2} \log \left( 2\pi e \frac{R^2}{d} \right).$$

*Proof.* Let  $\mu = \mathbb{E}X$  and  $\Sigma = \text{Cov}(X)$ . Among all distributions with covariance  $\Sigma$ , the Gaussian maximizes differential entropy, hence

$$h(X) \leq \frac{1}{2} \log((2\pi e)^d \det \Sigma).$$

Also,

$$\text{tr}(\Sigma) = \mathbb{E}\|X - \mu\|^2 \leq \mathbb{E}\|X\|^2 \leq R^2.$$

Let  $\lambda_1, \dots, \lambda_d$  be eigenvalues of  $\Sigma$ . Then  $\sum_i \lambda_i = \text{tr}(\Sigma) \leq R^2$  and  $\det \Sigma = \prod_i \lambda_i$ . By AM-GM,

$$\det \Sigma \leq \left( \frac{1}{d} \sum_{i=1}^d \lambda_i \right)^d \leq \left( \frac{R^2}{d} \right)^d.$$

Substituting gives

$$h(X) \leq \frac{1}{2} \log \left( (2\pi e)^d \left( \frac{R^2}{d} \right)^d \right) = \frac{d}{2} \log \left( 2\pi e \frac{R^2}{d} \right).$$

□

With this lemma, if the second moments of policies are uniformly upper bounded, then their entropies are also uniformly upper bounded, so are the soft-Q functions  $Q_\tau^\pi$  defined in (1).

Then we present the performance difference lemma which is the most important thing in reinforcement learning.

**Lemma 5** (Entropy Regularized Performance Difference Lemma). (*Lan, 2023, Lemma2*) For any two feasible policies  $\pi$  and  $\pi'$ , we have

$$V_{\tau}^{\pi'}(s) - V_{\tau}^{\pi}(s) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\pi'}} \left[ \langle A_{\tau}^{\pi}(s', \cdot), \pi'(\cdot|s') \rangle - \tau H^{\pi'}(s') + \tau H^{\pi}(s') \right],$$

where  $A_{\tau}^{\pi}(s', a) := Q_{\tau}^{\pi}(s', a) - V_{\tau}^{\pi}(s')$ .

For completeness, we provide the proof here and adapt it to our notation.

*Proof.* For simplicity, let us denote  $\xi^{\pi'}(s_0)$  the random process  $(s_t, a_t, s_{t+1})$ ,  $t \geq 0$ , generated by following the policy  $\pi'$  starting with the initial state  $s_0$ . It then follows from the definition of  $V_{\tau}^{\pi'}$  that

$$\begin{aligned} V_{\tau}^{\pi'}(s) - V_{\tau}^{\pi}(s) &= \mathbb{E}_{\xi^{\pi'}(s)} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) - \tau H^{\pi'}(s_t)) \right] - V_{\tau}^{\pi}(s) \\ &= \mathbb{E}_{\xi^{\pi'}(s)} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) - \tau H^{\pi'}(s_t) + V_{\tau}^{\pi}(s_t) - V_{\tau}^{\pi}(s_t)) \right] - V_{\tau}^{\pi}(s) \\ &\stackrel{(1)}{=} \mathbb{E}_{\xi^{\pi'}(s)} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) - \tau H^{\pi'}(s_t) + \gamma V_{\tau}^{\pi}(s_{t+1}) - V_{\tau}^{\pi}(s_t)) \right] \\ &\quad + \mathbb{E}_{\xi^{\pi'}(s)} [V_{\tau}^{\pi}(s_0)] - V_{\tau}^{\pi}(s) \\ &\stackrel{(2)}{=} \mathbb{E}_{\xi^{\pi'}(s)} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) - \tau H^{\pi'}(s_t) + \gamma V_{\tau}^{\pi}(s_{t+1}) - V_{\tau}^{\pi}(s_t)) \right] \\ &= \mathbb{E}_{\xi^{\pi'}(s)} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) - \tau H^{\pi}(s_t) + \gamma V_{\tau}^{\pi}(s_{t+1}) - V_{\tau}^{\pi}(s_t) - \tau H^{\pi'}(s_t) + \tau H^{\pi}(s_t)) \right] \\ &\stackrel{(3)}{=} \mathbb{E}_{\xi^{\pi'}(s)} \left[ \sum_{t=0}^{\infty} \gamma^t (Q_{\tau}^{\pi}(s_t, a_t) - V_{\tau}^{\pi}(s_t) - \tau H^{\pi'}(s_t) + \tau H^{\pi}(s_t)) \right]. \end{aligned}$$

where (1) follows by taking the term  $V_{\tau}^{\pi}(s_0)$  outside the summation, (2) follows from the fact that  $\mathbb{E}_{\xi^{\pi'}(s)} [V_{\tau}^{\pi}(s_0)] = V_{\tau}^{\pi}(s)$  since the random process starts with  $s_0 = s$ , and (3) follows from 1. The previous conclusion then imply that

$$\begin{aligned} V_{\tau}^{\pi'}(s) - V_{\tau}^{\pi}(s) &= \frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} d_{\gamma}^{\pi'}(s') \pi'(a'|s') [A_{\tau}^{\pi}(s', a') + \tau H^{\pi'}(s') - \tau H^{\pi}(s')] \\ &= \frac{1}{1-\gamma} \sum_{s' \in \mathcal{S}} d_{\gamma}^{\pi'}(s') [A_{\tau}^{\pi}(s', \pi'(\cdot|s')) - \tau H^{\pi'}(s') + \tau H^{\pi}(s')], \end{aligned}$$

which immediately implies the result.  $\square$

Remember our definition of  $\nu^*$  as the steady state distribution induced by  $\pi^*$ .

**Lemma 6.** (*Lan, 2023, Lemma3*)

$$\mathbb{E}_{s \sim \nu^*} [Q_{\tau}^{\pi}(s, \cdot), \pi^*(\cdot|s) - \pi(\cdot|s) - \tau H^{\pi^*}(s) + \tau H^{\pi}(s)] = \mathbb{E}_{s \sim \nu^*} [(1-\gamma) (V_{\tau}^{\pi^*}(s) - V_{\tau}^{\pi}(s))]. \quad (20)$$

For completeness, we provide the proof here and adapt it to our notation.

*Proof.* It follows from Lemma 5 (with  $\pi' = \pi^*$ ) that

$$(1-\gamma) [V_{\tau}^{\pi^*}(s) - V_{\tau}^{\pi}(s)] = \mathbb{E}_{s' \sim d_s^{\pi^*}} [A_{\tau}^{\pi}(s', \cdot), \pi^*(\cdot|s') + \tau H^{\pi}(s') - \tau H^{\pi^*}(s')].$$

Noting that:

$$\langle A_\tau^\pi(s', \cdot), \pi^*(\cdot|s') \rangle = \langle Q_\tau^\pi(s', \cdot), \pi^*(\cdot|s') \rangle - V_\tau^\pi(s') \quad (21)$$

$$= \langle Q_\tau^\pi(s', \cdot), \pi^*(\cdot|s') \rangle - \langle Q_\tau^\pi(s', \cdot), \pi(\cdot|s') \rangle \quad (22)$$

$$= \langle Q_\tau^\pi(s', \cdot), \pi^*(\cdot|s') - \pi(\cdot|s') \rangle, \quad (23)$$

Combining the above two relations and taking expectation w.r.t.  $\nu^*$ , we obtain

$$\begin{aligned} (1 - \gamma) \mathbb{E}_{s \sim \nu^*} [V_\tau^{\pi^*}(s) - V_\tau^\pi(s)] &= \mathbb{E}_{s \sim \nu^*, s' \sim d_{s^*}^\pi} [\langle Q_\tau^\pi(s', \cdot), \pi^*(\cdot|s') - \pi(\cdot|s') \rangle + \tau H^\pi(s') - \tau H^{\pi^*}(s')] \\ &= \mathbb{E}_{s \sim \nu^*} [\langle Q_\tau^\pi(s, \cdot), \pi^*(\cdot|s) - \pi(\cdot|s) - \tau H^{\pi^*}(s) \rangle + \tau H^\pi(s)] \end{aligned}$$

where the second identity is due to  $\nu^*$  is the steady state distribution induced by  $\pi^*$ .  $\square$

Next we will show a geometry property of the squared  $W_2$  distance that we repeatedly leverage in our convergence analysis.

**Lemma 7.** *Let  $\nu \in \mathcal{P}_2(\mathbb{R}^d)$ . For any  $\rho, \mu \in \mathcal{P}_2(\mathbb{R}^d)$ , let  $(\varphi^{\rho \rightarrow \nu}, \psi^{\rho \rightarrow \nu})$  be an optimal Kantorovich dual pair for the cost  $c(x, y) = \frac{1}{2} \|x - y\|^2$  between  $(\rho, \nu)$ , i.e.*

$$\varphi(x) + \psi(y) \leq \frac{1}{2} \|x - y\|^2 \quad \text{for all } x, y \in \mathbb{R}^d \quad (24)$$

$$\int \varphi^{\rho \rightarrow \nu} d\rho + \int \psi^{\rho \rightarrow \nu} d\nu = \frac{1}{2} W_2^2(\rho, \nu). \quad (25)$$

Then, for every  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ ,

$$\frac{1}{2} W_2^2(\mu, \nu) \geq \frac{1}{2} W_2^2(\rho, \nu) + \int_{\mathbb{R}^d} \varphi^{\rho \rightarrow \nu}(x) (\mu - \rho)(dx). \quad (26)$$

*Proof.* Recall the dual formulation (valid for any  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ ):

$$\frac{1}{2} W_2^2(\mu, \nu) = \sup_{\varphi, \psi} \left\{ \int \varphi d\mu + \int \psi d\nu : \varphi(x) + \psi(y) \leq \frac{1}{2} \|x - y\|^2 \right\}. \quad (27)$$

Since the feasibility constraint (24) is pointwise in  $(x, y)$ , the optimal pair  $(\varphi^{\rho \rightarrow \nu}, \psi^{\rho \rightarrow \nu})$  for  $(\rho, \nu)$  is also a feasible pair in the supremum (27) for  $(\mu, \nu)$ . Therefore,

$$\frac{1}{2} W_2^2(\mu, \nu) \geq \int \varphi^{\rho \rightarrow \nu} d\mu + \int \psi^{\rho \rightarrow \nu} d\nu. \quad (28)$$

By optimality for  $(\rho, \nu)$ , we have (25). Subtracting  $\int \varphi^{\rho \rightarrow \nu} d\rho$  on the right-hand side of (28) yields

$$\frac{1}{2} W_2^2(\mu, \nu) \geq \left( \int \varphi^{\rho \rightarrow \nu} d\rho + \int \psi^{\rho \rightarrow \nu} d\nu \right) + \int \varphi^{\rho \rightarrow \nu} d(\mu - \rho) = \frac{1}{2} W_2^2(\rho, \nu) + \int \varphi^{\rho \rightarrow \nu} d(\mu - \rho),$$

which is exactly (26).  $\square$

**Lemma 8** (Wasserstein proximal one step inequality). *Fix a state  $s \in \mathcal{S}$  and let the per-state JKO/proximal update be*

$$\bar{\pi}_{k+1}(\cdot|s) \in \operatorname{argmax}_{q \in \mathcal{P}_A} \left\{ \langle Q^{\pi_k}(s, \cdot), q \rangle - \tau H^q(s) - \frac{1}{2\eta_k} W_2^2(q, \pi_k(\cdot|s)) \right\}, \quad (29)$$

where  $\langle f, q \rangle := \int_{a \in \mathcal{A}} f(a) q(a) da$ . Then for any competitor  $p \in \Delta_A$ ,

$$\begin{aligned} &\eta_k \left( \langle Q^{\pi_k}(s, \cdot), p - \bar{\pi}_{k+1}(\cdot|s) \rangle - \tau H^p(s) + \tau H^{\bar{\pi}_{k+1}}(s) \right) + \frac{1}{2} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s)) \\ &\leq \frac{1}{2} W_2^2(p, \pi_k(\cdot|s)) - \frac{\eta_k \lambda \tau}{2} W_2^2(p, \bar{\pi}_{k+1}(\cdot|s)). \end{aligned}$$

*Proof.* Let  $\varphi^{\bar{\pi}_{k+1} \rightarrow \pi_k}(s, \cdot)$  be a Kantorovich potential for the pair  $(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s))$  under cost  $c(a, a') = \frac{1}{2} \|a - a'\|^2$ . Note that  $\langle Q^{\pi_k}(s, \cdot), q \rangle$  is a linear functional of  $q$ ,  $-\tau H^q(s)$  is strongly concave in  $q$ , and  $-\frac{1}{2\eta_k} W_2^2(q, \pi_k(\cdot|s))$  is concave in  $q$ . The first-order optimality condition of the concave program (29) states that

$$\left\langle \eta_k (Q^{\pi_k}(s, \cdot) - \tau(1 + \ln \bar{\pi}_{k+1}(\cdot|s))) - \varphi^{\bar{\pi}_{k+1} \rightarrow \pi_k}(\cdot, a), p(\cdot|s) - \bar{\pi}_{k+1}(\cdot|s) \right\rangle \leq 0, \quad \forall p(\cdot|s) \in \mathcal{P}_A. \quad (30)$$

Next, apply 7 for  $F(p) = \frac{1}{2} W_2^2(p, \pi_k(\cdot|s))$  and arbitrary  $p(\cdot|s)$ :

$$\frac{1}{2} W_2^2(p(\cdot|s), \pi_k(\cdot|s)) \geq \frac{1}{2} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s)) + \langle \varphi^{\bar{\pi}_{k+1} \rightarrow \pi_k}(s, \cdot), p(\cdot|s) - \bar{\pi}_{k+1}(\cdot|s) \rangle. \quad (31)$$

Rearranging (31) gives

$$\langle \varphi^{\bar{\pi}_{k+1} \rightarrow \pi_k}, p - \bar{\pi}_{k+1} \rangle \leq \frac{1}{2} W_2^2(p, \pi_k) - \frac{1}{2} W_2^2(\bar{\pi}_{k+1}, \pi_k),$$



where the arguments  $(\cdot|s)$  are omitted for readability. Plug this bound into optimal condition to obtain

$$\left\langle \eta_k(Q^{\pi_k}(s, \cdot) - \tau(1 + \ln \bar{\pi}_{k+1})), p - \bar{\pi}_{k+1} \right\rangle + \frac{1}{2}W_2^2(\bar{\pi}_{k+1}, \pi_k) \leq \frac{1}{2}W_2^2(p, \pi_k), \quad \forall p(\cdot|s) \in \mathcal{P}_{\mathcal{A}}. \quad (32)$$

By noting the two facts that

$$\langle 1, p - \bar{\pi}_{k+1}(\cdot|s) \rangle = 0,$$

and

$$\langle \ln \bar{\pi}_{k+1}(\cdot|s), p \rangle = \langle \ln \bar{\pi}_{k+1}(\cdot|s), p \rangle - \langle \ln p, p \rangle + \langle \ln p, p \rangle \quad (33)$$

$$= -\text{KL}(p \| \bar{\pi}_{k+1}(\cdot|s)) + H^p(s), \quad (34)$$

we have

$$\begin{aligned} \eta_k \left( \langle Q^{\pi_k}(s, \cdot), p - \bar{\pi}_{k+1}(\cdot|s) \rangle - \tau H^p(s) + \tau H^{\bar{\pi}_{k+1}}(s) \right) + \frac{1}{2}W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s)) \\ \leq \frac{1}{2}W_2^2(p, \pi_k(\cdot|s)) - \frac{\eta_k \lambda \tau}{2}W_2^2(p, \bar{\pi}_{k+1}(\cdot|s)). \end{aligned} \quad (35)$$

□

**Lemma 9.** For any  $s \in \mathcal{S}$ , we have

$$V_{\tau}^{\bar{\pi}_{k+1}}(s) - V_{\tau}^{\pi_k}(s) \geq \langle Q_{\tau}^{\pi_k}(s, \cdot), \bar{\pi}_{k+1}(\cdot|s) - \pi_k(\cdot|s) \rangle - \tau H^{\bar{\pi}_{k+1}}(s) + \tau H^{\pi_k}(s). \quad (36)$$

*Proof.* It follows from Lemma 5 (with  $\pi' = \pi^{k+1}$ ,  $\pi = \pi^k$ ) that

$$V_{\tau}^{\bar{\pi}_{k+1}}(s) - V_{\tau}^{\pi_k}(s) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\pi_k}} \left[ \langle A_{\tau}^{\pi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') \rangle - \tau H^{\bar{\pi}_{k+1}}(s') + \tau H^{\pi_k}(s') \right]. \quad (37)$$

And

$$\begin{aligned} \langle A_{\tau}^{\pi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') \rangle &= \langle Q_{\tau}^{\pi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') \rangle - V_{\tau}^{\pi_k}(s') \\ &= \langle Q_{\tau}^{\pi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') \rangle - \langle Q_{\tau}^{\pi_k}(s', \cdot), \pi_k(\cdot|s') \rangle \\ &= \langle Q_{\tau}^{\pi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_k(\cdot|s') \rangle. \end{aligned}$$

Combining the two identities above, we obtain

$$V_{\tau}^{\bar{\pi}_{k+1}}(s) - V_{\tau}^{\pi_k}(s) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\bar{\pi}_{k+1}}} \left[ \langle Q_{\tau}^{\pi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_k(\cdot|s') \rangle - \tau H^{\bar{\pi}_{k+1}}(s') + \tau H^{\pi_k}(s') \right]. \quad (38)$$

Now we conclude from Lemma 8 with  $p = \pi_k(\cdot|s')$  for any  $s'$  that

$$\langle Q_{\tau}^{\pi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_k(\cdot|s') \rangle - \tau H^{\bar{\pi}_{k+1}}(s') + \tau H^{\pi_k}(s') \geq \frac{\eta_k \lambda \tau}{2} W_2^2(\pi_k(\cdot|s'), \bar{\pi}_{k+1}(\cdot|s')). \quad (39)$$

The previous two conclusions then clearly imply the result in (36).

It also follows from (39) that

$$\begin{aligned} \mathbb{E}_{s' \sim d_s^{\pi_k}} \left[ \langle Q_{\tau}^{\pi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_k(\cdot|s') \rangle - \tau H^{\bar{\pi}_{k+1}}(s') + \tau H^{\pi_k}(s') \right] \\ \geq d_s^{\pi_k}(s) \left[ \langle Q_{\tau}^{\pi_k}(s, \cdot), \bar{\pi}_{k+1}(\cdot|s) - \pi_k(\cdot|s) \rangle - \tau H^{\bar{\pi}_{k+1}}(s) + \tau H^{\pi_k}(s) \right] \\ \geq (1-\gamma) \left[ \langle Q_{\tau}^{\pi_k}(s, \cdot), \bar{\pi}_{k+1}(\cdot|s) - \pi_k(\cdot|s) \rangle - \tau H^{\bar{\pi}_{k+1}}(s) + \tau H^{\pi_k}(s) \right], \end{aligned}$$

where the last inequality follows from the fact that  $d_s^{\pi_k}(s) \geq (1-\gamma)$  due to the definition of  $d_s^{\pi_k}$  and  $s_0 = s$  with probability one. Then by (38) and the above inequality, the claim follows. □

**Lemma 10** (Function Approximation Error). Under Assumption 1, we can bound the error of the value function induced by the function approximation step, i.e., for any  $s$ ,

$$|V_{\tau}^{\bar{\pi}_{k+1}}(s) - V_{\tau}^{\pi_{k+1}}(s)| \leq \mathcal{O}(\delta + \tau). \quad (40)$$

*Proof.* It follows from Lemma 5 (with  $\pi' = \bar{\pi}_{k+1}$ ,  $\pi = \pi_{k+1}$ ) that

$$V_{\tau}^{\bar{\pi}_{k+1}}(s) - V_{\tau}^{\pi_{k+1}}(s) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\bar{\pi}_{k+1}}} \left[ \langle Q_{\tau}^{\pi_{k+1}}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_{k+1} \rangle - \tau H^{\bar{\pi}_{k+1}}(s') + \tau H^{\pi_{k+1}}(s') \right]. \quad (41)$$

For any  $s'$ , the first term is bounded by

$$\begin{aligned} \langle Q_{\tau}^{\pi_{k+1}}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_{k+1} \rangle &\leq \|Q_{\tau}^{\pi_{k+1}}(s', \cdot)\|_{\infty} \text{TV}(\bar{\pi}_{k+1}(\cdot|s'), \pi_{k+1}(\cdot|s')) \\ &\leq B\delta, \end{aligned} \quad (42)$$

where the first inequality is due to Hölder's inequality, and the second follows from Pinsker's inequality, which says

$$\text{TV}(\bar{\pi}_{k+1}(\cdot|s'), \pi_{k+1}(\cdot|s')) \leq \sqrt{\frac{\text{KL}(\bar{\pi}_{k+1}(\cdot|s') \parallel \pi_{k+1}(\cdot|s'))}{2}} \leq \frac{\delta}{\sqrt{2}}.$$

By Lemma 4, we have  $\mathbf{H}^{\pi_{k+1}}(s') \leq \frac{d}{2} \log\left(2\pi e \frac{R^2}{d}\right)$ ,  $\mathbf{H}^{\bar{\pi}_{k+1}}(s') \leq \frac{d}{2} \log\left(2\pi e \frac{R^2}{d}\right)$

For ease of notation, let  $C := d \log\left(2\pi e \frac{R^2}{d}\right)$

Finally, we obtain

$$\begin{aligned} |V_{\tau}^{\bar{\pi}_{k+1}}(s) - V_{\tau}^{\pi_{k+1}}(s)| &= \left| \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\bar{\pi}_{k+1}}} [\langle Q_{\tau}^{\pi_{k+1}}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_{k+1} \rangle - \tau \mathbf{H}^{\bar{\pi}_{k+1}}(s') + \tau \mathbf{H}^{\pi_{k+1}}(s')] \right| \\ &\leq \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\bar{\pi}_{k+1}}} [|\langle Q_{\tau}^{\pi_{k+1}}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_{k+1} \rangle| + |\tau \mathbf{H}^{\bar{\pi}_{k+1}}(s') - \tau \mathbf{H}^{\pi_{k+1}}(s')|] \\ &\leq \frac{B}{1-\gamma} \delta + \frac{C}{1-\gamma} \tau. \end{aligned} \tag{43}$$

□

**Theorem 1 (Linear Convergence).** Suppose Assumption 1 holds and set the step size  $\eta_k = \eta \geq \frac{1}{\gamma\lambda\tau}$ . Then for any  $k \geq 0$ , the iterates of (WPPG) satisfy

$$J(\pi_*) - J(\pi_k) + \lambda\tau\mathcal{D}(\pi_k, \pi^*) \leq \gamma^k [J(\pi^*) - J(\pi_0) + \lambda\tau\mathcal{D}(\pi_0, \pi^*)] + \mathcal{O}(\delta + \tau)$$

where  $J$  is defined in (9), and  $\mathcal{D}(\pi_k, \pi^*) := \mathbb{E}_{s \sim \nu^*} [\frac{1}{2} \mathbf{W}_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s))]$ .

Consequently, in order to achieve an error of  $\mathcal{O}(\varepsilon + \delta + \tau)$ , the required iteration complexity is

$$\mathcal{O}\left(\frac{1}{1-\gamma} \log \frac{J(\pi^*) - J(\pi_0) + \lambda\tau\mathcal{D}(\pi_0, \pi^*)}{\varepsilon}\right).$$

*Proof.* By Lemma 8 with  $p = \pi^*$ , we have

$$\begin{aligned} &\langle Q_{\tau}^{\pi_k}(s, \cdot), \pi^*(\cdot|s) - \pi_{k+1}(\cdot|s) \rangle - \tau \mathbf{H}^{\pi^*}(s) + \tau \mathbf{H}^{\pi_{k+1}}(s) + \frac{1}{2\eta_k} \mathbf{W}_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s)) \\ &\leq \frac{1}{2\eta_k} \mathbf{W}_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) - \frac{\lambda\tau}{2\eta_k} \mathbf{W}_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s)). \end{aligned}$$

Combining with (36), we obtain

$$\begin{aligned} &[\langle Q_{\tau}^{\pi_k}(s, \cdot), \pi^*(\cdot|s) - \pi_k(\cdot|s) \rangle - \tau \mathbf{H}^{\pi^*}(s) + \tau \mathbf{H}^{\pi_k}(s)] \\ &\quad + [V_{\tau}^{\pi_k}(s) - V_{\tau}^{\bar{\pi}_{k+1}}(s)] \\ &\quad + \frac{1}{2\eta_k} \mathbf{W}_2^2(\pi_k(\cdot|s), \bar{\pi}_{k+1}(\cdot|s)) \\ &\leq [\langle Q_{\tau}^{\pi_k}(s, \cdot), \pi^*(\cdot|s) - \pi_k(\cdot|s) \rangle - \tau \mathbf{H}^{\pi^*}(s) + \tau \mathbf{H}^{\pi_k}(s)] \\ &\quad - [\langle Q_{\tau}^{\pi_k}(s, \cdot), \bar{\pi}_{k+1}(\cdot|s) - \pi_k(\cdot|s) \rangle - \tau \mathbf{H}^{\bar{\pi}_{k+1}}(s) + \tau \mathbf{H}^{\pi_k}(s)] \\ &\quad + \frac{1}{2\eta_k} \mathbf{W}_2^2(\pi_k(\cdot|s), \bar{\pi}_{k+1}(\cdot|s)) \\ &= (\langle Q_{\tau}^{\pi_k}(s, \cdot), \pi^*(\cdot|s) - \bar{\pi}_{k+1}(\cdot|s) \rangle - \tau \mathbf{H}^{\pi^*}(s) + \tau \mathbf{H}^{\bar{\pi}_{k+1}}(s)) \\ &\quad + \frac{1}{2\eta_k} \mathbf{W}_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s)) \\ &\leq \frac{1}{2\eta_k} \mathbf{W}_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) - \frac{\lambda\tau}{2} \mathbf{W}_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s)). \end{aligned}$$

Taking expectation with respect to  $\nu^*$  on both sides of the inequality, we obtain

$$\begin{aligned} \mathbb{E}_{s \sim \nu^*} \left[ (1 - \gamma) (V_{\tau}^{\pi^*}(s) - V_{\tau}^{\pi_k}(s)) \right] &+ \mathbb{E}_{s \sim \nu^*} \left[ V_{\tau}^{\pi_k}(s) - V_{\tau}^{\bar{\pi}_{k+1}}(s) \right] \\ &+ \mathbb{E}_{s \sim \nu^*} \left[ \frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s), \bar{\pi}_{k+1}(\cdot|s)) \right] \\ &\leq \mathbb{E}_{s \sim \nu^*} \left[ \frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) - \frac{\lambda\tau}{2} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s)) \right]. \end{aligned}$$

Then by Lemma 10, we have

$$\begin{aligned} \mathbb{E}_{s \sim \nu^*} \left[ V_{\tau}^{\pi_k}(s) - V_{\tau}^{\bar{\pi}_{k+1}}(s) \right] &= \mathbb{E}_{s \sim \nu^*} \left[ V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi_{k+1}}(s) + V_{\tau}^{\pi_{k+1}}(s) - V_{\tau}^{\bar{\pi}_{k+1}}(s) \right] \\ &\geq \mathbb{E}_{s \sim \nu^*} \left[ V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi_{k+1}}(s) \right] - \left( \frac{B}{1-\gamma} \delta + \frac{C}{1-\gamma} \tau \right). \end{aligned} \quad (44)$$

Note that we have assume the action space is bounded by  $R$  and  $\pi_{k+1}$  satisfy  $T_2$

$$\begin{aligned} &|W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s)) - W_2^2(\pi_{k+1}(\cdot|s), \pi^*(\cdot|s))| \\ &= \left| (W_2(\bar{\pi}_{k+1}, \pi^*) - W_2(\pi_{k+1}, \pi^*)) (W_2(\bar{\pi}_{k+1}, \pi^*) + W_2(\pi_{k+1}, \pi^*)) \right| \\ &\leq \sqrt{2}R W_2(\bar{\pi}_{k+1}(\cdot|s), \pi_{k+1}(\cdot|s)) \\ &\leq 4R\sqrt{\frac{1}{\lambda}}\delta. \end{aligned} \quad (45)$$

Hence,

$$W_2^2(\pi_{k+1}(\cdot|s), \pi^*(\cdot|s)) \geq W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s)) - 4R\sqrt{\frac{1}{\lambda}}\delta. \quad (46)$$

Similarly,

$$|W_2^2(\pi_{k+1}(\cdot|s), \pi_k(\cdot|s)) - W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s))| \leq 4R\sqrt{\frac{1}{\lambda}}\delta. \quad (47)$$

Combining (44), (46), and (47), we obtain

$$\begin{aligned} \mathbb{E}_{s \sim \nu^*} \left[ (1 - \gamma) (V_{\tau}^{\pi^*}(s) - V_{\tau}^{\pi_k}(s)) \right] &+ \mathbb{E}_{s \sim \nu^*} \left[ V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi_{k+1}}(s) \right] \\ &+ \mathbb{E}_{s \sim \nu^*} \left[ \frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s), \pi_{k+1}(\cdot|s)) \right] \\ &\leq \mathbb{E}_{s \sim \nu^*} \left[ \frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) - \frac{\lambda\tau}{2} W_2^2(\pi_{k+1}(\cdot|s), \pi^*(\cdot|s)) \right] \\ &+ \left( \frac{B}{1-\gamma} + \frac{2R}{\eta_k} \sqrt{\frac{1}{\lambda}} + 2R\tau\sqrt{\lambda} \right) \delta + \frac{C}{1-\gamma} \tau. \end{aligned} \quad (48)$$

By rewriting

$$V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi_{k+1}}(s) = V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi^*}(s) + V_{\tau}^{\pi^*}(s) - V_{\tau}^{\pi_{k+1}}(s),$$

and rearranging the inequality, we have

$$\begin{aligned} \mathbb{E}_{s \sim \nu^*} \left[ V_{\tau}^{\pi^*}(s) - V_{\tau}^{\pi_{k+1}}(s) \right] &+ \lambda\tau \mathbb{E}_{s \sim \nu^*} \left[ \frac{1}{2} W_2^2(\pi_{k+1}(\cdot|s), \pi^*(\cdot|s)) \right] \\ &+ \mathbb{E}_{s \sim \nu^*} \left[ \frac{1}{2} W_2^2(\pi_k(\cdot|s), \pi_{k+1}(\cdot|s)) \right] \\ &\leq \gamma \mathbb{E}_{s \sim \nu^*} \left[ V_{\tau}^{\pi^*}(s) - V_{\tau}^{\pi_k}(s) + \frac{1}{2\eta_k\gamma} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) \right] \\ &+ \left( \frac{B}{1-\gamma} + \frac{2R}{\eta_k} \sqrt{\frac{1}{\lambda}} + 2R\tau\sqrt{\lambda} \right) \delta + \frac{C}{1-\gamma} \tau. \end{aligned}$$

Thus,

$$\begin{aligned} \mathbb{E}_{s \sim \nu^*} \left[ V_{\tau}^{\pi^*}(s) - V_{\tau}^{\pi_{k+1}}(s) + \frac{\lambda\tau}{2} W_2^2(\pi_{k+1}(\cdot|s), \pi^*(\cdot|s)) \right] \\ \leq \gamma \mathbb{E}_{s \sim \nu^*} \left[ V_{\tau}^{\pi^*}(s) - V_{\tau}^{\pi_k}(s) + \frac{1}{2\eta_k\gamma} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) \right] \\ + \left( \frac{B}{1-\gamma} + \frac{R}{\eta_k} \sqrt{\frac{2}{\lambda}} + 2R\tau\sqrt{2\lambda} \right) \delta + \frac{C}{1-\gamma} \tau. \end{aligned}$$

Recalling the definitions of  $J$  (9) and  $\mathcal{D}$ , we obtain

$$J(\pi^*) - J(\pi_{k+1}) + \lambda\tau\mathcal{D}(\pi_{k+1}, \pi^*) \leq \gamma \left[ J(\pi^*) - J(\pi_k) + \frac{1}{\eta_k\gamma}\mathcal{D}(\pi_k, \pi^*) \right] + \left( \frac{B}{1-\gamma} + \frac{2R}{\eta_k} \sqrt{\frac{1}{\lambda}} + 2R\tau\sqrt{\lambda} \right) \delta + \frac{C}{1-\gamma}\tau. \quad (49)$$

Choosing  $\eta_k = \eta \geq \frac{1}{\gamma\lambda\tau}$  in the JKO scheme, we obtain

$$J(\pi^*) - J(\pi_{k+1}) + \lambda\tau\mathcal{D}(\pi_{k+1}, \pi^*) \leq \gamma \left[ J(\pi^*) - J(\pi_k) + \lambda\tau\mathcal{D}(\pi_k, \pi^*) \right] + \left( \frac{B}{1-\gamma} + 3R\tau\sqrt{2\lambda} \right) \delta + \frac{C}{1-\gamma}\tau,$$

which implies

$$J(\pi^*) - J(\pi_k) + \lambda\tau\mathcal{D}(\pi_k, \pi^*) \leq \gamma^k \left[ J(\pi^*) - J(\pi_0) + \lambda\tau\mathcal{D}(\pi_0, \pi^*) \right] + \mathcal{O}(\delta + \tau).$$

□

For the ease of presentation, we denote  $\Delta_k = Q^{\pi_k, \xi_k} - Q^{\pi_k}$  and  $\xi_{0:k} = \{\xi_0, \xi_1, \dots, \xi_k\}$  in the following paper.

**Lemma 11.** *Under Assumption 2, for any state  $s$  we have:*

$$\mathbb{E}_{\xi_{0:k}} [\langle \Delta_k(\cdot, s), \pi_{k+1}(\cdot|s) - \pi_k(\cdot|s) \rangle] \leq 2\eta_k\sigma_k^2 + \frac{1}{2\eta_k} \mathbb{E}_{\xi_{0:k}} W_2^2(\pi_k(\cdot|s), \pi_{k+1}(\cdot|s)) \quad (50)$$

*Proof.* For any  $s$ , let  $\gamma(a, a'|s)$  be the optimal couple of the two distribution  $\pi_k(\cdot|s)$  and  $\pi_{k+1}(\cdot|s)$  in  $W_2$ .

$$\begin{aligned} & \mathbb{E}_{\xi_{0:k}} [\langle \Delta_k(a, s), \pi_{k+1}(a|s) - \pi_k(a|s) \rangle | \xi_{0:k-1}] \\ &= \mathbb{E}_{\xi_k} \left[ \int_{\mathcal{A}} \Delta_k(\cdot, s) d(\pi_{k+1}(\cdot|s) - \pi_k(\cdot|s)) | \xi_{0:k-1} \right] \\ &= \mathbb{E}_{\xi_{0:k}} \left[ \int \int_{\mathcal{A} \times \mathcal{A}} \Delta_k(a, s) - \Delta_k(a', s) d\gamma(a, a'|s) | \xi_{0:k-1} \right] \\ &= \mathbb{E}_{\xi_{0:k}} \left[ \int \int_{\mathcal{A} \times \mathcal{A}} \int \langle \nabla_a \Delta_k((1-t)a' + ta, s), a - a' \rangle dt d\gamma(a, a'|s) | \xi_{0:k-1} \right] \\ &= \int \int_{\mathcal{A} \times \mathcal{A}} \int \mathbb{E}_{\xi_{0:k}} [\langle \nabla_a \Delta_k((1-t)a' + ta, s), a - a' \rangle | \xi_{0:k-1}] dt d\gamma(a, a'|s) \\ &\leq \int \int_{\mathcal{A} \times \mathcal{A}} \int \mathbb{E}_{\xi_{0:k}} [2\eta_k \|\nabla_a \Delta_k((1-t)a' + ta, s)\|_2^2 + \frac{1}{2\eta_k} \|a - a'\|_2^2 | \xi_{0:k-1}] dt d\gamma(a, a'|s) \\ &\leq 2\eta_k\sigma_k^2 + \frac{1}{2\eta_k} \mathbb{E}_{\xi_{0:k}} [W_2^2(\pi_k(\cdot|s), \pi_{k+1}(\cdot|s)) | \xi_{0:k-1}] \end{aligned}$$

The second equality applies the definition of an optimal coupling  $\gamma(\cdot, \cdot|s) \in \Gamma(\pi_{k+1}(\cdot|s), \pi_k(\cdot|s))$ , which means has the same marginal distribution as  $\pi_k$  and  $\pi_{k+1}$ . The second inequality uses Young's inequality  $\langle u, v \rangle \leq \frac{1}{2\eta_k} \|u\|^2 + \frac{\eta_k}{2} \|v\|^2$  to separate the two terms. The last inequality bounds the variance term of the stochastic gradient by  $\sigma_k^2$  yields the last inequality, where the quadratic term recovers the squared Wasserstein distance between  $\pi_k(\cdot|s)$  and  $\pi_{k+1}(\cdot|s)$ .

Taking expectation with respect to  $\xi_{0:k-1}$  on both sides, we have the final result:

$$\mathbb{E}_{\xi_{0:k}} [\langle \Delta_k(\cdot, s), \pi_{k+1}(\cdot|s) - \pi_k(\cdot|s) \rangle] \leq 2\eta_k\sigma_k^2 + \frac{1}{2\eta_k} \mathbb{E}_{\xi_{0:k}} W_2^2(\pi_k(\cdot|s), \pi_{k+1}(\cdot|s))$$

□

**Theorem 2 (Linear Convergence).** *Suppose Assumptions 1 and 2 hold, and for all  $k \geq 0$ ,  $\epsilon_k \leq \epsilon$ ,  $\sigma_k \leq \sigma$  and  $\|Q^{\pi_k, \xi_k}\|_\infty \leq B$ . Then the iterates of (11) using step size  $\eta_k = \eta \geq \frac{1}{\gamma\lambda\tau}$  satisfies*

$$\mathbb{E}_{\xi_{0:k-1}} [J(\pi^*) - J(\pi_k) + \lambda\tau\mathcal{D}(\pi_k, \pi^*)] \leq \gamma^k [J(\pi^*) - J(\pi_0) + \lambda\tau\mathcal{D}(\pi_0, \pi^*)] + \mathcal{O}(\delta + \tau + \epsilon + \sigma) \quad (15)$$

where  $J$  is defined in (9), and  $\mathcal{D}(\pi_k, \pi^*) := \mathbb{E}_{s \sim \nu^*} [\frac{1}{2} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s))]$ . Consequently, in order to achieve an error of  $\mathcal{O}(\epsilon + \tau + \delta + \epsilon + \sigma)$  in expectation, the required iteration complexity

is

$$\mathcal{O}\left(\frac{1}{1-\gamma} \log \frac{J(\pi^*) - J(\pi_0) + \lambda\tau\mathcal{D}(\pi_0, \pi^*)}{\varepsilon}\right).$$

*Proof.* By Lemma 8 applied to 11 with  $p = \pi^*$ , we have

$$\begin{aligned} & \langle Q_{\tau}^{\pi_k, \xi_k}(s, \cdot), \pi^*(\cdot|s) - \bar{\pi}_{k+1}(\cdot|s) \rangle - \tau H^{\pi^*}(s) + \tau H^{\bar{\pi}_{k+1}}(s) + \frac{1}{2\eta_k} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s)) \\ & \leq \frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) - \frac{\lambda\tau}{2} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s)). \end{aligned} \quad (51)$$

By Lemma 8 applied to 11 with  $p = \pi_k$ , we have

$$\begin{aligned} & \left( \langle Q_{\tau}^{\pi_k, \xi_k}(s, \cdot), \pi_k(\cdot|s) - \bar{\pi}_{k+1}(\cdot|s) \rangle - \tau H^{\pi_k}(s) + \tau H^{\bar{\pi}_{k+1}}(s) \right) + \frac{1}{2\eta_k} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s)) \\ & \leq -\frac{\lambda\tau}{2} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s)) \leq 0. \end{aligned}$$

Which implies that

$$\begin{aligned} & \mathbb{E}_{s' \sim d_s^{\pi_k}} [\langle Q_{\tau}^{\pi_k, \xi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_k(\cdot|s') \rangle - \tau H^{\bar{\pi}_{k+1}}(s') + \tau H^{\pi_k}(s') + \frac{1}{2\eta_k} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s))] \\ & \leq d_s^{\pi_k}(s) [\langle Q_{\tau}^{\pi_k, \xi_k}(s, \cdot), \bar{\pi}_{k+1}(\cdot|s) - \pi_k(\cdot|s) \rangle - \tau H^{\bar{\pi}_{k+1}}(s) + \tau H^{\pi_k}(s) + \frac{1}{2\eta_k} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s))] \\ & \leq (1-\gamma) [\langle Q_{\tau}^{\pi_k, \xi_k}(s, \cdot), \bar{\pi}_{k+1}(\cdot|s) - \pi_k(\cdot|s) \rangle - \tau H^{\bar{\pi}_{k+1}}(s) + \tau H^{\pi_k}(s) + \frac{1}{2\eta_k} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s))] \end{aligned} \quad (52)$$

where the last inequality follows from the fact that  $d_s^{\pi_k}(s) \geq (1-\gamma)$  due to the definition of  $d_s^{\pi_k}$  and  $s_0 = s$  with probability one.

Note that we can still use the performance difference identity 38

$$\begin{aligned} V_{\tau}^{\bar{\pi}_{k+1}}(s) - V_{\tau}^{\pi_k}(s) &= \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\pi_k}} [\langle Q_{\tau}^{\pi_k, \xi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_k(\cdot|s') \rangle - \tau H^{\bar{\pi}_{k+1}}(s') + \tau H^{\pi_k}(s')] \\ &= \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\pi_k}} [\langle Q_{\tau}^{\pi_k, \xi_k}(s', \cdot), \bar{\pi}_{k+1}(\cdot|s') - \pi_k(\cdot|s') \rangle - \tau H^{\bar{\pi}_{k+1}}(s') + \tau H^{\pi_k}(s') \\ & \quad - \langle \Delta_k(\cdot, s'), \bar{\pi}_{k+1}(\cdot|s') - \pi_k(\cdot|s') \rangle] \end{aligned} \quad (53)$$

By multiplying both sides by -1 and taking expectation with respect to  $\xi_{0:k}$  gives

$$\begin{aligned} & \mathbb{E}_{\xi_{0:k}} [V_{\tau}^{\pi_k}(s) - V_{\tau}^{\bar{\pi}_{k+1}}(s)] \\ & \leq \frac{1}{1-\gamma} \mathbb{E}_{\xi_{0:k}} \mathbb{E}_{s' \sim d_s^{\pi_k}} [\langle Q_{\tau}^{\pi_k, \xi_k}(s', \cdot), \pi_k(\cdot|s') - \bar{\pi}_{k+1}(\cdot|s') \rangle - \tau H^{\pi_k}(s') + \tau H^{\bar{\pi}_{k+1}}(s') \\ & \quad + \frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s'), \bar{\pi}_{k+1}(\cdot|s'))] + 2\eta_k \sigma_k^2 \\ & \leq \mathbb{E}_{\xi_{0:k}} [\langle Q_{\tau}^{\pi_k, \xi_k}(s, \cdot), \pi_k(\cdot|s) - \bar{\pi}_{k+1}(\cdot|s) \rangle - \tau H^{\pi_k}(s) + \tau H^{\bar{\pi}_{k+1}}(s) \\ & \quad + \frac{1}{2\eta_k} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi_k(\cdot|s))] + \frac{2\eta_k \sigma_k^2}{1-\gamma}. \end{aligned} \quad (54)$$

Taking expectation with  $\xi_{0:k}$  on 51 and combine with 54, we have:

$$\begin{aligned} & \mathbb{E}_{\xi_{0:k}} [\langle Q_{\tau}^{\pi_k, \xi_k}(s, \cdot), \pi_k(\cdot|s) - \pi^*(\cdot|s) \rangle + \tau H^{\pi_k}(s) - \tau H^{\pi^*}(s) + V_{\tau}^{\bar{\pi}_{k+1}}(s) - V_{\tau}^{\pi_k}(s)] \\ & \leq \mathbb{E}_{\xi_{0:k}} [\frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) - \frac{\lambda\tau}{2} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s))] + \frac{2\eta_k \sigma_k^2}{1-\gamma}. \end{aligned}$$

Finally, averaging over the distribution  $s \sim \nu^*$  and noting that  $s$  and  $\xi_{0:k}$  are independent, we have

$$\begin{aligned} & \mathbb{E}_{s \sim \nu^*, \xi_{0:k}} [\langle Q_{\tau}^{\pi_k, \xi_k}(s, \cdot), \pi^*(\cdot|s) - \pi_k(\cdot|s) \rangle - \tau H^{\pi^*}(s) + \tau H^{\pi_k}(s) + V_{\tau}^{\bar{\pi}_{k+1}}(s) - V_{\tau}^{\pi_k}(s)] \\ & \leq \mathbb{E}_{s \sim \nu^*, \xi_{0:k}} [\frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) - \frac{\lambda\tau}{2} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s))] + \frac{2\eta_k \sigma_k^2}{1-\gamma}. \end{aligned} \quad (55)$$

Noting that

$$\begin{aligned}
& \mathbb{E}_{\xi_k} [\langle Q_{\tau}^{\pi_k, \xi_k}(s, \cdot), \pi^*(\cdot|s) - \pi_k(\cdot|s) \rangle \mid \xi_{0:k-1}] \\
&= \mathbb{E}_{\xi_k} [\langle Q_{\tau}^{\pi_k}(s, \cdot), \pi^*(\cdot|s) - \pi_k(\cdot|s) \rangle + \langle \bar{Q}_{\tau}^{\pi_k}(s, \cdot) - Q_{\tau}^{\pi_k}(s, \cdot), \pi^*(\cdot|s) - \pi_k(\cdot|s) \rangle \\
&+ \langle Q_{\tau}^{\pi_k, \xi_k}(s, \cdot) - \bar{Q}_{\tau}^{\pi_k}(s, \cdot), \pi^*(\cdot|s) - \pi_k(\cdot|s) \rangle \mid \xi_{0:k-1}] \\
&\geq \langle Q_{\tau}^{\pi_k}(s, \cdot), \pi^*(\cdot|s) - \pi_k(\cdot|s) \rangle - 2\epsilon_k
\end{aligned} \tag{56}$$

The first equality expands  $Q_{\tau}^{\pi_k, \xi_k}$  into its expectation  $Q_{\tau}^{\pi_k}$  plus two error terms, namely the bias  $\bar{Q}_{\tau}^{\pi_k} - Q_{\tau}^{\pi_k}$  and the stochastic fluctuation  $Q_{\tau}^{\pi_k, \xi_k} - \bar{Q}_{\tau}^{\pi_k}$ . Taking conditional expectation w.r.t.  $\xi_k$  eliminates the mean of the fluctuation term. Finally, using the uniform error bound  $\|\bar{Q}_{\tau}^{\pi_k} - Q_{\tau}^{\pi_k}\|_{\infty} \leq \epsilon_k$  and noting that both  $\pi_k(\cdot|s)$  and  $\pi^*(\cdot|s)$  are probability measures (which implies  $\|\pi^*(\cdot|s) - \pi_k(\cdot|s)\|_1 \leq 2$ ), Hölder's inequality yields the desired bound.

Combining 55 and 56 and using Lemma 6,

$$\begin{aligned}
& \mathbb{E}_{s \sim \nu^*, \xi_{0:k}} [(1 - \gamma)(V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi^*}(s)) + V_{\tau}^{\pi_{k+1}}(s) - V_{\tau}^{\pi_k}(s)] \\
&\leq \mathbb{E}_{s \sim \nu^*, \xi_{0:k}} \left[ \frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) - \frac{\lambda\tau}{2} W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s)) \right] + 2\epsilon_k + \frac{2\eta_k\sigma_k^2}{1 - \gamma}.
\end{aligned} \tag{57}$$

Then by 10, we have

$$|V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi^*}(s)| \leq \mathcal{O}(\delta + \tau) \tag{58}$$

And Note that the action space is bounded by  $R$  and  $\pi_{k+1}$  satisfy  $T_2$

$$\begin{aligned}
& |W_2^2(\bar{\pi}_{k+1}(\cdot|s), \pi^*(\cdot|s)) - W_2^2(\pi_{k+1}(\cdot|s), \pi^*(\cdot|s))| \\
&= \left| (W_2(\bar{\pi}_{k+1}, \pi^*) - W_2(\pi_{k+1}, \pi^*)) (W_2(\bar{\pi}_{k+1}, \pi^*) + W_2(\pi_{k+1}, \pi^*)) \right| \\
&\leq \sqrt{2}R W_2(\bar{\pi}_{k+1}(\cdot|s), \pi_{k+1}(\cdot|s)) \\
&\leq \mathcal{O}(\delta).
\end{aligned} \tag{59}$$

combining (57), (58), (59), we have

$$\begin{aligned}
& \mathbb{E}_{s \sim \nu^*, \xi_{0:k}} [(1 - \gamma)(V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi^*}(s)) + V_{\tau}^{\pi_{k+1}}(s) - V_{\tau}^{\pi_k}(s)] \\
&\leq \mathbb{E}_{s \sim \nu^*, \xi_{0:k}} \left[ \frac{1}{2\eta_k} W_2^2(\pi_k(\cdot|s), \pi^*(\cdot|s)) - \frac{\lambda\tau}{2} W_2^2(\pi_{k+1}(\cdot|s), \pi^*(\cdot|s)) \right] \\
&+ \mathcal{O}(\delta + \tau) + 2\epsilon_k + \frac{2\eta_k\sigma_k^2}{1 - \gamma}.
\end{aligned} \tag{60}$$

Decomposing  $V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi_{k+1}}(s)$  into  $V_{\tau}^{\pi_k}(s) - V_{\tau}^{\pi^*}(s) - (V_{\tau}^{\pi_{k+1}}(s) - V_{\tau}^{\pi^*}(s))$ , recalling our definition of  $J$  (9) and rearranging the terms in the above inequality, we get

$$\begin{aligned}
& \mathbb{E}_{\xi_{0:k}} [J(\pi^*) - J(\pi_{k+1}) + \lambda\tau\mathcal{D}(\pi_{k+1}, \pi^*)] \\
&\leq \mathbb{E}_{\xi_{0:k-1}} [\gamma(J(\pi^*) - J(\pi_k)) + \frac{1}{\eta_k}\mathcal{D}(\pi^*, \pi_k)] + \mathcal{O}(\delta + \tau) + 2\epsilon_k + \frac{\eta_k\sigma_k^2}{2(1 - \gamma)}.
\end{aligned} \tag{61}$$

By choosing  $\eta_k = \eta \geq \frac{1}{\gamma\lambda\tau}$ , and for all  $k \geq 0$ ,  $\epsilon_k \leq \epsilon$ ,  $\sigma_k \leq \sigma$ , we get

$$\mathbb{E}_{\xi_{0:k-1}} [J(\pi^*) - J(\pi_k) + \lambda\tau\mathcal{D}(\pi_k, \pi^*)] \leq \gamma^k [J(\pi^*) - J(\pi_0) + \lambda\tau\mathcal{D}(\pi_0, \pi^*)] + \mathcal{O}(\delta + \tau + \epsilon + \sigma) \tag{62}$$

□

## B NUMERICAL

### B.1 OVERALL ALGORITHM

**WPPG vs. WPPG-I: commonalities and differences** Both WPPG and WPPG-I are off-policy actor-critic methods built on the same backbone: (i) replay-based training with 1-step TD targets; (ii) Double- $Q$  critics with target networks and Polyak averaging; (iii) multi-sample bootstrap for target construction (average over  $K$  next-action samples and take  $\min(Q_1, Q_2)$ ); (iv) actor updates driven by action-gradient matching, i.e., aligning the policy’s action increment with a noisy target direction, and (v)  $\tanh$  squashing that maps actions to box constraints.

*Key difference.* WPPG employs an explicit Tanh–Gaussian policy  $a = \text{Affine}(\tanh(\mu_\theta(s) + \sigma_\theta(s) \odot \varepsilon))$  with a closed-form density (useful if one wishes to incorporate entropy/KL terms). WPPG-I uses a latent-conditioned *implicit* policy  $a = \text{Affine}(\tanh(f_\theta([s, z])))$ , where  $z \sim \mathcal{N}(0, I)$  is concatenated with the state; the policy distribution is implicit (no closed-form  $\log \pi$ ), and learning relies purely on pathwise gradients through  $\nabla_a Q$  and the shared latent variable reforwarding trick. Operationally, WPPG controls stochasticity via the Gaussian actor’s output scale, whereas WPPG-I controls it via the *input* latent variables (its scale and dimensionality), enabling richer, state-conditional exploration.

---

#### Algorithm 1 WPPG with Replay and Double- $Q$ Critics (Gaussian policy)

---

**Require:** Initialize actor  $\pi_\theta(a|s) = \mathcal{N}(\mu_\theta(s), \Sigma_\theta(s))$  with Tanh squash; twin critics  $Q_{w_1}, Q_{w_2}$  and targets  $\bar{Q}_{w_1}, \bar{Q}_{w_2}$ ; target actor  $\bar{\pi}_\theta$ ; replay buffer  $\mathcal{D}$ ; step size  $\eta$ , noise scale  $\tau$ , samples per state  $K$ , discount  $\gamma$ , Polyak  $\sigma$ .

```

1: for each episode do
2:   Initialize  $s_0$ 
3:   for  $t = 0$  to  $T - 1$  do
4:     Sample  $a_t \sim \pi_\theta(\cdot|s_t)$ , execute, observe  $r_t, s_{t+1}$  and store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
5:     if  $\text{len}(\mathcal{D}) \geq \text{batch\_size}$  then
6:       Sample a minibatch  $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^B$  from  $\mathcal{D}$ 
7:       Compute 1-step TD targets (multi-sample bootstrap using target nets):
8:       For each  $s'_i$ , draw  $\epsilon_{i,1:K} \sim \mathcal{N}(0, I)$  and set  $a'_{i,k} \leftarrow \bar{\pi}_\theta(s'_i; \epsilon_{i,k})$ 
9:        $\hat{Q}_i \leftarrow \frac{1}{K} \sum_{k=1}^K \min(\bar{Q}_{w_1}(s'_i, a'_{i,k}), \bar{Q}_{w_2}(s'_i, a'_{i,k}))$ 
10:       $y_i \leftarrow r_i + \gamma \hat{Q}_i$ 
11:      Critic update (train both critics):
12:       $w_j \leftarrow w_j - \beta_Q \nabla_{w_j} \frac{1}{B} \sum_i (Q_{w_j}(s_i, a_i) - y_i)^2, \quad j \in \{1, 2\}$ 
13:      Actor update (WPPG step with action-sample direction):
14:      For each  $s_i$ , draw shared  $\epsilon_{i,1:K}$  and form  $a_{i,k} = \pi_\theta(s_i; \epsilon_{i,k})$ ; let  $A_i = [a_{i,1:K}]$ 
15:      Compute  $q_{i,k} = \min(Q_{w_1}(s_i, a_{i,k}), Q_{w_2}(s_i, a_{i,k}))$ 
16:      Obtain  $\nabla_a Q$  at samples:  $G_i = [\nabla_a q_{i,k}]_{k=1}^K$ 
17:      Form noisy target direction:  $\Delta_i^* \leftarrow \eta G_i + \xi_i$ , where  $\xi_i \sim \mathcal{N}(0, 2\tau\eta I)$ 
18:      Re-sample  $A'_i = \pi_\theta(s_i; \epsilon_{i,1:K})$  with the same  $\epsilon$  and define  $\Delta_i \leftarrow A'_i - A_i$ 
19:      Update actor by matching directions:  $\theta \leftarrow \theta - \beta_\pi \nabla_\theta \frac{1}{BK} \sum_{i,k} \|\Delta_{i,k} - \Delta_{i,k}^*\|_2^2$ 
20:      Target updates (Polyak):  $\bar{w}_j \leftarrow \sigma w_j + (1 - \sigma)\bar{w}_j, \bar{\theta} \leftarrow \sigma\theta + (1 - \sigma)\bar{\theta}, j \in \{1, 2\}$ 
21:    end if
22:     $s_{t+1} \leftarrow s'$ 
23:  end for
24: end for

```

---



**Algorithm 2** WPPG-I with Replay and Double- $Q$  Critics (Implicit Policy)

**Require:** Implicit actor  $a = g_\theta(s, z)$  with Tanh squash ( $z \sim \mathcal{N}(0, I_M)$ ); twin critics  $Q_{w_1}, Q_{w_2}$  and targets  $\bar{Q}_{w_1}, \bar{Q}_{w_2}$ ; target actor  $\bar{g}_\theta$ ; replay buffer  $\mathcal{D}$ ; step size  $\eta$ , noise scale  $\tau$ , samples per state  $K$ , discount  $\gamma$ , Polyak  $\sigma$ .

```

1: for each episode do
2:   Initialize  $s_0$ 
3:   for  $t = 0$  to  $T - 1$  do
4:     Sample  $z_t \sim \mathcal{N}(0, I_M)$ , set  $a_t = g_\theta(s_t, z_t)$ , step env, observe  $(r_t, s_{t+1}, d_t)$ 
5:     Store  $(s_t, a_t, r_t, s_{t+1}, d_t)$  into  $\mathcal{D}$ 
6:     if  $\text{len}(\mathcal{D}) \geq \text{batch\_size}$  then
7:       Sample a minibatch  $\{(s_i, a_i, r_i, s'_i, d_i)\}_{i=1}^B$  from  $\mathcal{D}$ 

8:       Compute 1-step TD targets (multi-sample bootstrap):
9:       For each  $s'_i$ , draw  $z'_{i,1:K} \sim \mathcal{N}(0, I_M)$  and set  $a'_{i,k} \leftarrow \bar{g}_\theta(s'_i, z'_{i,k})$ 
10:       $\hat{Q}_i \leftarrow \frac{1}{K} \sum_{k=1}^K \min(\bar{Q}_{w_1}(s'_i, a'_{i,k}), \bar{Q}_{w_2}(s'_i, a'_{i,k}))$ 
11:       $y_i \leftarrow r_i + \gamma(1 - d_i) \hat{Q}_i$ 

12:      Critic update (train both critics):
13:       $w_j \leftarrow w_j - \beta_Q \nabla_{w_j} \frac{1}{B} \sum_i (Q_{w_j}(s_i, a_i) - y_i)^2, \quad j \in \{1, 2\}$ 

14:      Actor update (direction matching with shared noise):
15:      For each  $s_i$ , draw shared  $z_{i,1:K} \sim \mathcal{N}(0, I_M)$  and set  $a_{i,k}^{(0)} \leftarrow g_\theta(s_i, z_{i,k})$ 
16:      Compute  $q_{i,k} = \min(Q_{w_1}(s_i, a_{i,k}^{(0)}), Q_{w_2}(s_i, a_{i,k}^{(0)}))$ 
17:      Obtain  $G_i = [\nabla_a q_{i,k}]_{k=1}^K$  at  $a^{(0)}$  (stop grad to critics)
18:      Form target direction:  $\Delta_i^* \leftarrow \eta G_i + \xi_i$ , where  $\xi_i \sim \mathcal{N}(0, 2\tau\eta I)$ 
19:      Reforward with the same  $z$ :  $a_{i,k}^{(1)} \leftarrow g_\theta(s_i, z_{i,k})$ , define  $\Delta_{i,k} \leftarrow a_{i,k}^{(1)} - a_{i,k}^{(0)}$ 
20:      Update actor:  $\theta \leftarrow \theta - \beta_\pi \nabla_\theta \frac{1}{BK} \sum_{i,k} \|\Delta_{i,k} - \Delta_i^*\|_2^2$ 

21:      Target updates (Polyak):  $\bar{w}_j \leftarrow \sigma w_j + (1 - \sigma)\bar{w}_j, \bar{\theta} \leftarrow \sigma\theta + (1 - \sigma)\bar{\theta}, j \in \{1, 2\}$ 
22:    end if
23:     $s_{t+1} \leftarrow s'$ 
24:  end for
25: end for

```

## B.2 IMPLEMENTATIONS

## B.2.1 ACTOR AND POLICY

All our models and baselines are implemented under the standard actor-critic framework. Below we provide the implementation details for the actor and critic components separately.

**Action Squashing.** For consistency across methods, we apply a *tanh* squashing function to map sampled actions into the valid box  $[a_{\min}, a_{\max}]$  for all algorithms. This squashing is crucial for the implicit policy: without it, when the injected-latent dimension is high, many actions are hard-clipped at the bounds, which prevents meaningful exploration and gradients, often leading to training failure. Below we present concise formulations of the two actors used.

**Tanh-Gaussian MLP Policy (used in WPPG/SAC/WPO).** Given state  $s \in \mathbb{R}^S$ , the actor outputs  $\mu_\theta(s), \log \sigma_\theta(s) \in \mathbb{R}^A$  and samples

$$a = \frac{a_{\max} - a_{\min}}{2} \odot \tanh(\mu_\theta(s) + \sigma_\theta(s) \odot \varepsilon) + \frac{a_{\max} + a_{\min}}{2}, \quad \varepsilon \sim \mathcal{N}(0, I_A),$$

i.e., a tanh-squashed Gaussian mapped to  $[a_{\min}, a_{\max}]$  via an MLP producing  $(\mu_\theta, \log \sigma_\theta)$ .

**Noise-Conditioned Deterministic Policy (used in WPPG-I).** Given state  $s \in \mathbb{R}^S$  and latent variables  $z \in \mathbb{R}^M$ ,

$$a = \frac{a_{\max} - a_{\min}}{2} \odot \tanh(f_\theta([s, z])) + \frac{a_{\max} + a_{\min}}{2}, \quad z \sim \mathcal{N}(0, I_M),$$

where  $f_\theta$  is an MLP taking the concatenated input  $[s, z]$ . This defines an implicit policy (no closed-form density) with tanh-squashed outputs mapped to  $[a_{\min}, a_{\max}]$ .

### B.2.2 CRITIC

**Critic Learning Target** For all off-policy algorithms (SAC, WPPG, WPPG-I, WPO), the critic is trained with 1-step TD targets that average over  $K$  bootstrap action samples and use Double- $Q$  when available:

$$y_t = r_t + \gamma(1-d_t) \frac{1}{K} \sum_{k=1}^K \min_{j \in \{1,2\}} Q_{\bar{w}_j}(s_{t+1}, a'_{t+1,k}), \quad a'_{t+1,k} = g_{\bar{\theta}}(s_{t+1}, \varepsilon_k), \quad \varepsilon_k \sim \mathcal{N}(0, I).$$

Here,  $j \in \{1, 2\}$  indexes the two target critics used by Double- $Q$  (the per-sample minimum is taken), and the outer average is over  $K$  target actions drawn from the target actor  $g_{\bar{\theta}}$ . For single- $Q$  methods (e.g., WPO uses a single  $Q$ , or the single- $Q$  WPPG ablation), replace  $\min_{j \in \{1,2\}} Q_{\bar{w}_j}$  by  $Q_{\bar{w}}$ . In contrast, PPO retains its on-policy generalized advantage estimation (GAE) for actor updates.

### B.3 HYPERPARAMETERS

**Neural Network Architecture** All experiments are based on two standard network configurations: a larger network with hidden sizes (256, 256) and ReLU activation, and a smaller network with hidden sizes (64, 64) and Tanh activation. We use the larger network for Hopper, Humanoid, and HalfCheetah, and the smaller network for all other tasks. The same choice is applied uniformly across all models, and the actor and critic share the same network architecture.

**Replay Buffer** For consistency, all off-policy algorithms (SAC, WPPG, WPPG-I, WPO) use the same replay buffer configuration as summarized in Table 1, ensuring identical storage capacity, sampling scheme, and update frequency across methods.

**Training Setup** All off-policy models share the same basic training setup: each is trained for  $1 \times 10^6$  timesteps, with evaluation performed every 2000 steps. Target networks are updated via Polyak averaging to stabilize critic training. Both actor and critic use a learning rate of  $3 \times 10^{-4}$ . The discount factor is set to  $\gamma = 0.99$  for all tasks, except Swimmer where  $\gamma = 0.9999$ . The configuration is summarized in Table 2. For PPO, hyperparameters are environment-specific and detailed below.

Table 1: Replay buffer.

Hyperparameter	Value
Buffer size	1,000,000
Batch size	256
Learning starts	10,000
Train frequency	1 (step)
Gradient steps per update	1
Number of environments	1

Table 2: Training Setup.

Hyperparameter	Value
Discount factor $\gamma$	0.99
Polyak coefficient	0.005
Learning rate (actor/critic)	$3 \times 10^{-4}$
Target update interval	1
Total timesteps	1,000,000
Optimizer	Adam

**Model Specific Hyperparameters** The hyperparameters of all models are summarized in the tables below. For WPO and SAC, we follow the settings reported in the WPO paper, while the hyperparameters of PPO are taken from RL Zoo. Although we find that tuning hyperparameters for each environment can often improve performance, for simplicity and fairness of comparison we adopt a single unified set of hyperparameters for all off-policy methods across all tasks, which yields stable and competitive learning performance.

### B.4 ADDITIONAL EXPERIMENT RESULTS

**Multi-Run Evaluation** To more comprehensively demonstrate the behavior of WPPG and WPPG-I, we further evaluate both methods with multiple training runs. Specifically, each algorithm is trained 5 times with different random seeds. In the corresponding plots, the solid line denotes the mean return across the 5 runs, while the shaded area indicates the range between the minimum and maximum returns over these runs (see Fig.3).

Table 3: WPPG-I hyperparameters.

Hyperparameter	Default Value
Action samples	32
Step size $\eta$	10
Noise scale $\tau$	0.00001
Actor Latent Dimension	$\frac{1}{3} \times \text{State Dimension}$

Table 4: WPPG hyperparameters.

Hyperparameter	Default Value
Action samples	32
Step size $\eta$	0.01
Noise scale $\tau$	0.00001

Table 5: SAC hyperparameters.

Hyperparameter	Default Value
Entropy Coefficient $\alpha$	0.001
Maximum Policy Variance	$\exp(4)$
Minimum Policy Variance	$\exp(-10)$

Table 6: WPO hyperparameters.

Hyperparameter	Default Value
KL Mean Penalty $\alpha_\mu$	0.001
KL Variance Penalty $\alpha_\Sigma$	0.001
Action samples	32

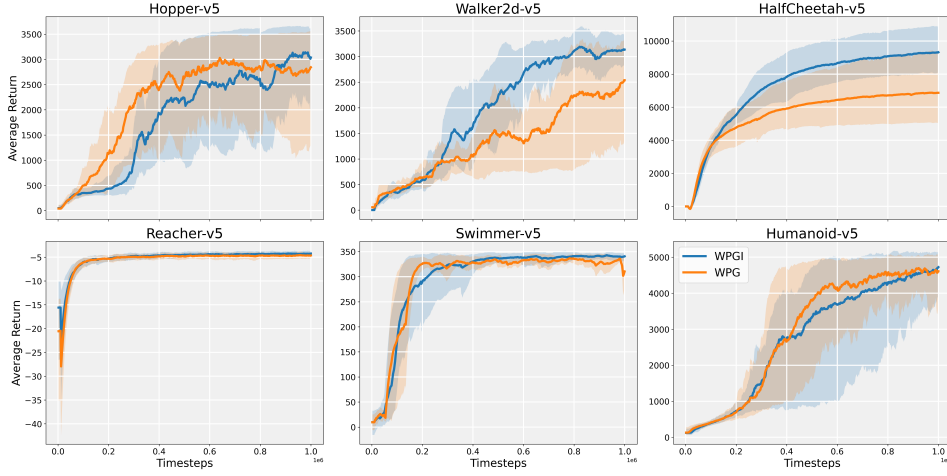


Figure 3: Multi-Run Evaluation

**Combined Humanoid Task** Our method shares some similarities with SAC in that both are based on entropy regularization and use the action gradient of the Q-function for policy updates. However, the key advantage of WPPG lies in its ability to train an implicit policy. To better showcase this benefit, we follow the construction in WPO Pfau et al. (2025) and create a combined task that increases the action dimensionality: multiple Humanoid environments are run in parallel, their states are concatenated and fed into a single agent, which outputs the concatenated actions jointly. As shown in the combined Humanoid task, WPPG-I converges to consistently higher returns than SAC, indicating that the implicit policy is able to discover action distributions that achieve higher rewards. (see Fig.4).

**Ablation on Double Q Function** Double-Q plays a crucial role for WPPG. As shown in the figures, although the single-Q variant of WPPG outperforms WPO on most environments, it fails to achieve fast and stable learning on challenging tasks such as Humanoid. Beyond stability, the use of double-Q also opens up interesting directions for further exploration; for example, one could choose the Q-function with the smaller gradient magnitude to provide the action-sample update direction. We leave such extensions for future work.

**Additional Ablation Study** Beyond the main results, we also conduct additional ablation studies on Humanoid-v5 and HalfCheetah-v5, systematically varying key hyperparameters (e.g., the Wasserstein step size  $\eta$ , the number of sampled actions, and the latent dimension of the implicit policy). These experiments, reported in the supplementary material, further validate the robustness of our method and illustrate how performance and stability depend on these design choices. (see Fig.6 and Fig.7).

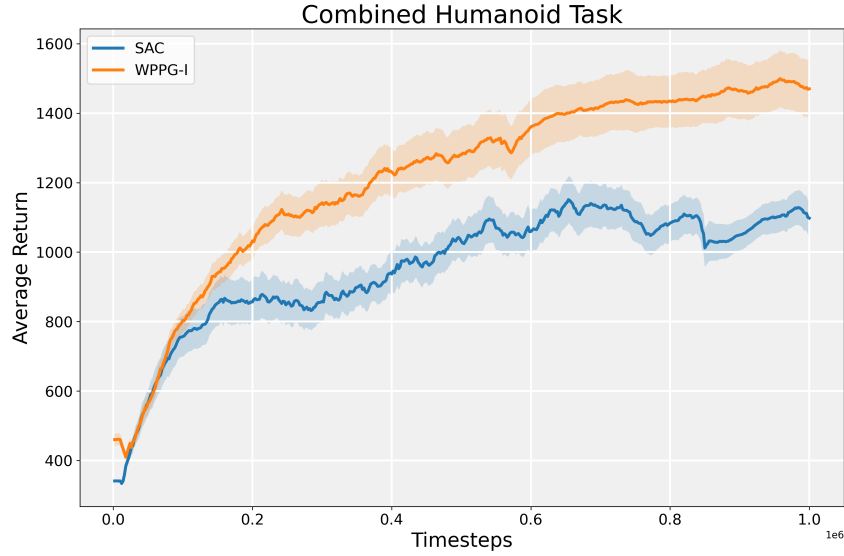


Figure 4: Combined Humanoid Task

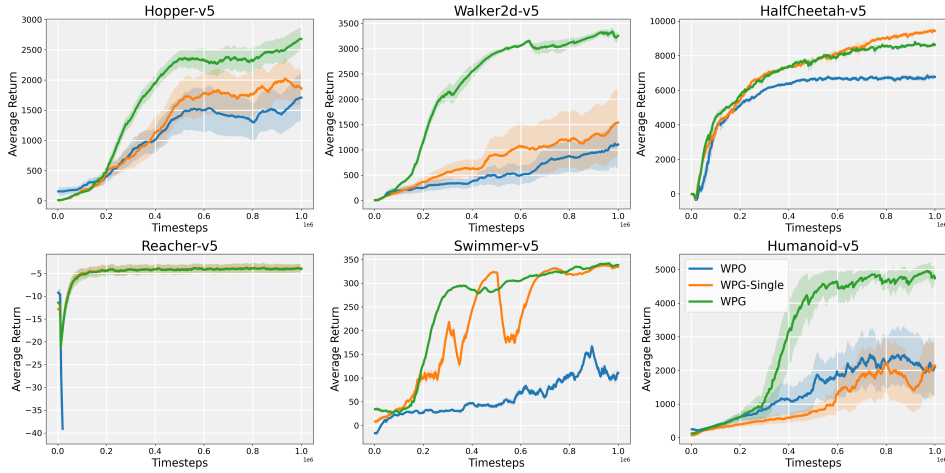


Figure 5: Ablation on Double Q Function

**Training Speed Comparison** All experiments are conducted on an NVIDIA RTX 4090 GPU. The training times on the high-dimensional Humanoid-v5 task are reported in Table 7. As can be seen, WPPG, WPPG-I, and WPO achieve comparable training speeds, all of which are faster than SAC.

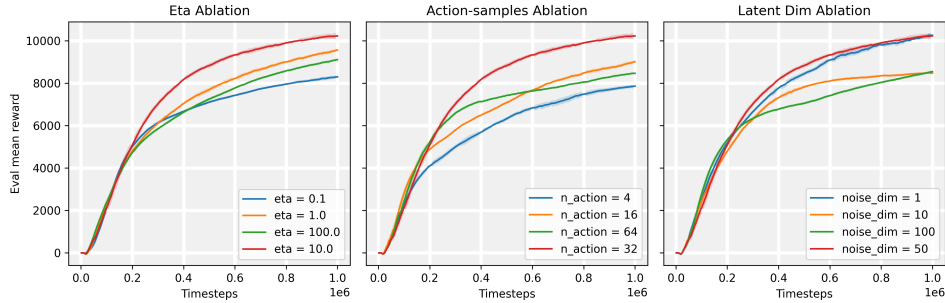


Figure 6: Additional Ablation on Eta, Action Samples and Latent Dim with HalfCheetah-v5 Task

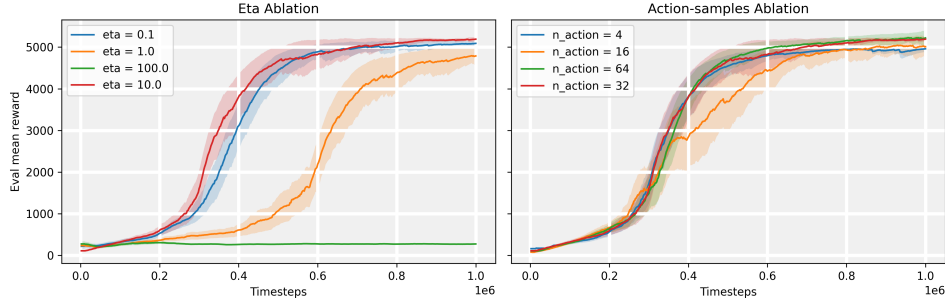


Figure 7: Additional Ablation on Eta, Action Samples and Latent Dim with Humanoid-v5 Task

Table 7: Training time on Humanoid-v5 (in seconds).

Model	Time (s)
SAC	19451
WPPG-I	12113
WPPG	13135
WPO	10649
PPO	2104

## C ADDITIONAL DISCUSSION

We now discuss a practically relevant case where the action space  $\mathcal{A}$  is convex and bounded. To ease the presentation, we fix a state  $s$  in the discussion below, and all constants should be interpreted as a uniform constant for all states.

Consider the policy update

$$\rho_{k+1} \in \operatorname{argmin}_{\rho \in \mathcal{P}(\mathcal{A})} \left\{ \int_{\mathcal{A}} Q_k(s, a) \rho(a) da + \tau \int_{\mathcal{A}} \rho \log \rho da + \frac{1}{2\eta} W_2^2(\rho, \rho_k) \right\}.$$

Assume  $Q_k$  has a uniformly bounded gradient. Assume the initial policy distribution  $\rho_0$  satisfying

$$0 < m_0 \leq \rho_0(a) \leq M_0 < \infty \quad \text{for all } a \in \mathcal{A}. \quad (63)$$

Fix a continuous distribution  $\nu$  on  $\mathcal{A}$  with density  $f$  such that

$$0 < \tilde{m} \leq f(a) \leq \tilde{M} < \infty \quad \text{for all } a \in \mathcal{A}. \quad (64)$$

The optimality condition yields

$$\frac{\rho_{k+1} - \rho_k}{\eta} = \tau \Delta \rho_{k+1} + \nabla \cdot (\rho_{k+1} \nabla Q) \quad \text{in } A, \quad (65)$$

with Neumann boundary condition

$$(\tau \nabla \rho_{k+1} + \rho_{k+1} \nabla Q) \cdot n = 0 \quad \text{on } \partial A. \quad (66)$$

By the maximum principle, if  $\rho_k$  has uniformly positive lower and upper bounds, so does  $\rho_{k+1}$ . For each  $k \geq 0$ , let  $T_k : \mathcal{A} \rightarrow \mathcal{A}$  be the Brenier map pushing  $\nu$  to  $\rho_k$ . Then there exists a convex potential  $\phi_k : \mathcal{A} \rightarrow \mathbb{R}$  such that

$$T_k(a) = \nabla \phi_k(a).$$

The convex gradient map  $T_k = \nabla \phi_k$  solves the Monge-Ampère equation

$$\det D^2 \phi_k(x) = \frac{f(x)}{g_k(\nabla \phi_k(x))} \quad \text{for } x \in A.$$

Then by Caffarelli’s regularity theory,  $\phi_k$  has uniformly bounded Hessian. Therefore

$$\operatorname{Lip}(T_k) = \sup_{x \in A} \|DT_k(x)\| = \sup_{x \in A} \|D^2 \phi_k(x)\| \leq C \quad \text{for all } k \geq 0.$$

Thus, the family  $\{T_k\}_{k \geq 0}$  of optimal transport maps is uniformly Lipschitz. As a direct consequence, if  $\nu$  satisfies a  $T_2$  inequality, then so do the JKO iterates  $\rho_k$  with a uniform constant.