
Model-Based Reparameterization Policy Gradient Methods: Theory and Practical Algorithms

Shenao Zhang¹ Boyi Liu² Zhaoran Wang² Tuo Zhao¹

Abstract

ReParameterization (RP) Policy Gradient Methods (PGMs) have been widely adopted for continuous control tasks in robotics and computer graphics. However, recent studies have revealed that, when applied to long-term reinforcement learning problems, model-based RP PGMs may experience chaotic and non-smooth optimization landscapes with exploding gradient variance, which leads to slow convergence. This is in contrast to the conventional belief that reparameterization methods have low gradient estimation variance in problems such as training deep generative models. To comprehend this phenomenon, we conduct a theoretical examination of model-based RP PGMs and search for solutions to the optimization difficulties. Specifically, we analyze the convergence of the model-based RP PGMs and pinpoint the smoothness of function approximators as a major factor that affects the quality of gradient estimation. Based on our analysis, we propose a spectral normalization method to mitigate the exploding variance issue caused by long model unrolls. Our experimental results demonstrate that proper normalization significantly reduces the gradient variance of model-based RP PGMs. As a result, the performance of the proposed method is comparable or superior to other gradient estimators, such as the Likelihood Ratio (LR) gradient estimator.

1 Introduction

Reinforcement Learning (RL) has seen tremendous success in a variety of sequential decision-making applications, such as strategy games (Silver et al., 2017; Vinyals et al., 2019) and robotics (Duan et al., 2016; Wang et al., 2019b), by identifying actions that maximize long-term accumulated rewards. As one of the most popular methodologies, the

policy gradient methods (PGM) (Sutton et al., 1999; Kakade, 2001; Silver et al., 2014) seek to search for the optimal policy by iteratively computing and following a stochastic gradient direction with respect to the policy parameters. Therefore, the quality of the stochastic gradient estimation is essential for the effectiveness of PGMs.

Two main categories have emerged in the realm of stochastic gradient estimation: (1) Likelihood Ratio (LR) estimators, which perform zeroth-order estimation through the sampling of function evaluations (Williams, 1992; Konda & Tsitsiklis, 1999; Kakade, 2001), and (2) ReParameterization (RP) gradient estimators, which harness the differentiability of the function approximation (Figurnov et al., 2018; Ruiz et al., 2016; Clavera et al., 2020; Suh et al., 2022a).

Despite the wide adoption of both LR and RP PGMs in practice, the majority of the literature on the theoretical properties of PGMs focuses on LR PGMs. The optimality and approximation error of LR PGMs have been heavily investigated under various settings (Agarwal et al., 2021; Wang et al., 2019a; Bhandari & Russo, 2019). Conversely, the theoretical underpinnings of RP PGMs remain to be fully explored, with a dearth of research on the quality of RP gradient estimators and the convergence of RP PGMs.

RP gradient estimators have established themselves as a reliable technique for training deep generative models such as variational autoencoders (Figurnov et al., 2018). From a stochastic optimization perspective, previous studies (Ruiz et al., 2016; Mohamed et al., 2020) have shown that RP gradient methods enjoy small variance, which leads to better convergence and performance. However, recent research (Parmas et al., 2018; Metz et al., 2021) has reported an opposite observation: When applied to long-horizon reinforcement learning problems, model-based RP PGMs tend to encounter chaotic optimization procedures and highly non-smooth optimization landscapes with exploding gradient variance, causing slow convergence.

Such an intriguing phenomenon inspires us to delve deeper into the theoretical properties of RP gradient estimators in search of a remedy for the issue of exploding gradient variance in model-based RP PGMs. To this end, we present a unified theoretical framework for the examination of model-

¹Georgia Tech, Atlanta, GA, USA ²Northwestern University, Evanston, IL, USA.

based RP PGMs and establish their convergence results. Our analysis implies that the smoothness and accuracy of the learned model are crucial determinants of the exploding variance of RP gradients: (1) Both the gradient variance and bias exhibit a polynomial dependence on the Lipschitz continuity of the learned model and policy w.r.t. the input state, with degrees that increase linearly with the steps of model value expansion, and (2) the bias also depends on the error of the estimated model and value.

Our findings suggest that imposing smoothness on the model and policy can greatly decrease the variance of RP gradient estimators. To put this discovery into practice, we propose a spectral normalization method to enforce the smoothness of the learned model and policy. It's worth noting that this method can enhance the algorithm's efficiency without substantially compromising accuracy when the underlying transition kernel is smooth. However, if the transition kernel is not smooth, enforcing smoothness may lead to increased error in the learned model and introduce bias. In such cases, a balance should be struck between model bias and gradient variance. Nonetheless, our empirical study demonstrates that the reduced gradient variance when applying spectral normalization leads to a significant performance boost, even with the cost of a higher bias. Furthermore, our results highlight the potential of investigating model-based RP PGMs, as they demonstrate superiority over other model-based and Likelihood Ratio (LR) gradient estimator alternatives.

2 Background

Reinforcement Learning. We consider learning to optimize an infinite-horizon γ -discounted Markov Decision Process (MDP) over repeated episodes of interaction. We denote by $\mathcal{S} \subseteq \mathbb{R}^{d_s}$ and $\mathcal{A} \subseteq \mathbb{R}^{d_a}$ the state space and the action space, respectively. When taking an action $a \in \mathcal{A}$ at a state $s \in \mathcal{S}$, the agent receives a reward $r(s, a)$ and the MDP transits to a new state s' according to probability $s' \sim f(\cdot | s, a)$.

We aim to find a policy π that maps a state to an action distribution to maximize the expected cumulative reward. We denote by $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ and $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ the state value function and the state-action value function associated with π , respectively, which are defined as follows,

$$V^\pi(s) = (1 - \gamma) \cdot \mathbb{E}_{\pi, f} \left[\sum_{i=0}^{\infty} \gamma^i \cdot r(s_i, a_i) \mid s_0 = s \right],$$

$$Q^\pi(s, a) = (1 - \gamma) \cdot \mathbb{E}_{\pi, f} \left[\sum_{i=0}^{\infty} \gamma^i \cdot r(s_i, a_i) \mid s_0 = s, a_0 = a \right].$$

Here $s \in \mathcal{S}$, $a \in \mathcal{A}$, and the expectation $\mathbb{E}_{\pi, f}[\cdot]$ is taken with respect to the dynamic induced by the policy π and the transition probability f .

We denote by ζ the initial state distribution. Under policy π , the state and state-action visitation measure $\nu_\pi(s)$ over

\mathcal{S} and $\sigma_\pi(s, a)$ over $\mathcal{S} \times \mathcal{A}$ are defined as follows,

$$\nu_\pi(s) = (1 - \gamma) \cdot \sum_{i=0}^{\infty} \gamma^i \cdot \mathbb{P}(s_i = s),$$

$$\sigma_\pi(s, a) = (1 - \gamma) \cdot \sum_{i=0}^{\infty} \gamma^i \cdot \mathbb{P}(s_i = s, a_i = a),$$

respectively. Here the summations are taken with respect to the trajectory induced by $s_0 \sim \zeta$, $a_i \sim \pi(\cdot | s_i)$, and $s_{i+1} \sim f(\cdot | s_i, a_i)$. The objective $J(\pi)$ of RL is defined as the expected policy value as follows,

$$J(\pi) = \mathbb{E}_{s_0 \sim \zeta} [V^\pi(s_0)] = \mathbb{E}_{(s, a) \sim \sigma_\pi} [r(s, a)]. \quad (2.1)$$

Stochastic Gradient Estimation. The underlying problem of policy gradient, i.e., computing the gradient of an expectation with respect to the parameters of the sampling distribution, takes the form $\nabla_\theta \mathbb{E}_{p(x; \theta)} [y(x)]$. To restore the RL objective, we can set $p(x; \theta)$ as the trajectory distribution conditioned on the policy parameter θ and $y(x)$ as the cumulative reward. In the sequel, we introduce two commonly used gradient estimators.

Likelihood Ratio (LR) Gradient (Zeroth-Order): By leveraging the *score function*, LR gradients only require samples of the function values. Since $\nabla_\theta \log p(x; \theta) = \nabla_\theta p(x; \theta) / p(x; \theta)$, the LR gradient takes the form of

$$\nabla_\theta \mathbb{E}_{p(x; \theta)} [y(x)] = \mathbb{E}_{p(x; \theta)} [y(x) \nabla_\theta \log p(x; \theta)]. \quad (2.2)$$

ReParameterization (RP) Gradient (First-Order): RP gradient benefits from the structural characteristics of the objective, i.e., how the overall objective is affected by the operations applied to the sources of randomness as they pass through the measure and into the cost function (Mohamed et al., 2020). From the simulation property of continuous distribution, we have the following equivalence between direct and indirect ways of drawing samples,

$$\hat{x} \sim p(x; \theta) \iff \hat{x} = g(\epsilon; \theta), \epsilon \sim p. \quad (2.3)$$

Derived from the *law of the unconscious statistician* (LOTUS) (Grimmett & Stirzaker, 2020), i.e., $\mathbb{E}_{p(x; \theta)} [y(x)] = \mathbb{E}_{p(\epsilon)} [y(g(\epsilon; \theta))]$, the RP gradient can be formulated as

$$\nabla_\theta \mathbb{E}_{p(x; \theta)} [y(x)] = \mathbb{E}_{p(\epsilon)} [\nabla_\theta y(g(\epsilon; \theta))].$$

3 Analytic Reparameterization Gradient in Reinforcement Learning

In this section, we present two fundamental *analytic* forms of the RP gradient in RL. We first consider the Policy-Value Gradient (PVG) method, which is model-free and can be expanded sequentially to obtain the Analytic Policy Gradient (APG) method. Then we discuss potential obstacles that may arise when developing practical algorithms.

We consider a policy $\pi_\theta(s, \varsigma)$ with noise ς in continuous action spaces. To ensure that the first-order gradient through the value function is well-defined, we make the following assumption on the continuity of the MDP.

Assumption 3.1 (Continuous MDP). We assume that $f(s' | s, a)$, $\pi_\theta(s, \varsigma)$, $r(s, a)$, and $\nabla_a r(s, a)$ are continuous in all parameters and variables s, a, s' .

Policy-Value Gradient. The reparameterization Policy-Value Gradient (PVG) takes the following general form,

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{s \sim \zeta, \varsigma \sim p} [\nabla_\theta Q^{\pi_\theta}(s, \pi_\theta(s, \varsigma))]. \quad (3.1)$$

By performing sequential decision-making, any immediate action could lead to changes in all future states and rewards. Therefore, the value gradient $\nabla_\theta Q^{\pi_\theta}$ possesses a recursive structure. Adapted from the deterministic policy gradient theorem (Silver et al., 2014; Lillicrap et al., 2015) by taking stochasticity into consideration, we rewrite (3.1) as

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{s \sim \nu_\pi, \varsigma} [\nabla_\theta \pi_\theta(s, \varsigma) \cdot \nabla_a Q^{\pi_\theta}(s, a)|_{a=\pi_\theta(s, \varsigma)}].$$

Here, $\nabla_a Q^{\pi_\theta}$ can be estimated using a critic, which leads to model-free frameworks (Heess et al., 2015; Amos et al., 2021). Notably, as a result of the recursive structure of $\nabla_\theta Q^{\pi_\theta}$, the expectation is taken over the state visitation ν_π instead of the initial distribution ζ .

By sequentially expanding PVG, we obtain the following dynamical representation of the policy gradient.

Analytic Policy Gradient. Due to the simulation property of continuous distributions in (2.3), we interchangeably write $a \sim \pi_\theta(\cdot | s)$ with $a = \pi_\theta(s, \varsigma)$ and $s' \sim f(\cdot | s, a)$ with $s' = f(s, a, \xi^*)$, where ξ^* is sampled from the unknown distribution p . From the Bellman equation $V^{\pi_\theta}(s) = \mathbb{E}_\varsigma[(1-\gamma) \cdot r(s, \pi_\theta(s, \varsigma)) + \gamma \cdot \mathbb{E}_{\xi^*}[V^{\pi_\theta}(f(s, \pi_\theta(s, \varsigma), \xi^*))]]$, we obtain the following backward recursions of gradient,

$$\nabla_\theta V^{\pi_\theta}(s) = \mathbb{E}_\varsigma [(1-\gamma) \nabla_a r \nabla_\theta \pi_\theta \quad (3.2)$$

$$+ \gamma \mathbb{E}_{\xi^*} [\nabla_{s'} V^{\pi_\theta}(s') \nabla_a f \nabla_\theta \pi_\theta + \nabla_\theta V^{\pi_\theta}(s')]],$$

$$\nabla_s V^{\pi_\theta}(s) = \mathbb{E}_\varsigma [(1-\gamma)(\nabla_s r + \nabla_a r \nabla_s \pi_\theta) \quad (3.3)$$

$$+ \gamma \mathbb{E}_{\xi^*} [\nabla_{s'} V^{\pi_\theta}(s') (\nabla_s f + \nabla_a f \nabla_s \pi_\theta)]].$$

See §A for detailed derivations of (3.2) and (3.3). Now we have the RP gradient backpropagated through the transition path starting at s . By taking an expectation over the initial state distribution, we obtain the Analytic Policy Gradient (APG) that takes the following form,

$$\text{APG:} \quad \nabla_\theta J(\pi_\theta) = \mathbb{E}_{s \sim \zeta} [\nabla_\theta V^{\pi_\theta}(s)].$$

There remain challenges when developing practical algorithms: (1) The above formulas require the gradient information of the transition function f . In this work, however, we consider a common RL setting where f and its derivatives are unknown and need to be fitted by a model. It is thus natural to ask how the properties of the model (e.g., prediction accuracy and model smoothness) affect the gradient estimation and the convergence of the resulting algorithms, and (2) even if we have access to an accurate model, unrolling

it over full sequences faces practical difficulties. The memory and computational cost scale linearly with the unroll length. Long chains of nonlinear mappings can also lead to exploding or vanishing gradients and even worse, chaotic phenomena (Bollt, 2000) and difficulty in optimization (Pascanu et al., 2013; Maclaurin et al., 2015; Vicol et al., 2021; Metz et al., 2019). These difficulties demand some form of truncation when performing RP PGMs.

4 Model-Based RP Policy Gradient Methods

Through the application of Model Value Expansion (MVE) for model truncation, this section unveils two RP policy gradient frameworks constructed upon MVE.

4.1 h -Step Model Value Expansion

To handle the difficulties inherent in full unrolls, many algorithms employ direct truncation, where the long sequence is broken down into short sub-sequences and backpropagation is applied accordingly, e.g., Truncated BPTT (Werbos, 1990). However, such an approach over-prioritizes short-term dependencies, which leads to biased gradient estimates.

In model-based RL (MBRL), one viable solution is to adopt the h -step Model Value Expansion (Feinberg et al., 2018), which decomposes the value estimation $\hat{V}^\pi(s)$ into the rewards gleaned from the learned model and a residual estimated by a critic function \hat{Q}_ω , that is,

$$\hat{V}^{\pi_\theta}(s) = (1-\gamma) \cdot \left(\sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_i, \hat{a}_i) + \gamma^h \cdot \hat{Q}_\omega(\hat{s}_h, \hat{a}_h) \right),$$

where $\hat{s}_0 = s$, $\hat{a}_i = \pi_\theta(\hat{s}_i, \varsigma)$, and $\hat{s}_{i+1} = \hat{f}_\psi(\hat{s}_i, \hat{a}_i, \xi)$. Here, the noise variables ς and ξ can be sampled from the fixed distributions or inferred from the real samples, which we will discuss in the following section.

4.2 Model-Based RP Gradient Estimation

Utilizing the pathwise gradient with respect to θ , we present the following two frameworks.

Model Derivatives on Predictions (DP). A straightforward way to compute the first-order gradient is to link the reward, model, policy, and critic together and backpropagate through them. Specifically, the differentiation is carried out on the trajectories simulated by the model \hat{f}_ψ , which serves as a tool for *both* the prediction of states and the evaluation of derivatives. The corresponding RP-DP estimator of gradient $\nabla_\theta J(\pi_\theta)$ is denoted as $\hat{\nabla}_\theta^{\text{DP}} J(\pi_\theta)$, which takes the form of

$$\hat{\nabla}_\theta^{\text{DP}} J(\pi_\theta) = \frac{1}{N} \sum_{n=1}^N \nabla_\theta \left(\sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) + \gamma^h \cdot \hat{Q}_\omega(\hat{s}_{h,n}, \hat{a}_{h,n}) \right), \quad (4.1)$$

where $\hat{s}_{0,n} \sim \mu_{\pi_\theta}$, $\hat{a}_{i,n} = \pi_\theta(\hat{s}_{i,n}, \varsigma_n)$, and $\hat{s}_{i+1,n} = \hat{f}_\psi(\hat{s}_{i,n}, \hat{a}_{i,n}, \xi_n)$ with noises $\varsigma_n \sim p(\varsigma)$ and $\xi_n \sim p(\xi)$.

Here, μ_{π_θ} is the distribution where the initial states of the simulated trajectories are sampled. In Section 5, we study a general form of μ_{π_θ} that is a mixture of the initial state distribution ζ and the state visitation ν_{π_θ} .

Various algorithms can be instantiated from (4.1) with different choices of h . When $h = 0$, the framework reduces to model-free policy gradients, such as RP(0) (Amos et al., 2021) and the variants of DDPG (Lillicrap et al., 2015), e.g., SAC (Haarnoja et al., 2018). When $h \rightarrow \infty$, the resulting algorithm is BPTT (Grzeszczuk et al., 1998; Degraeve et al., 2019; Bastani, 2020) where only the model is learned. Recent model-based approaches, such as MAAC (Clavera et al., 2020) and related algorithms (Parmas et al., 2018; Amos et al., 2021; Li et al., 2021), require a carefully selected h .

Model Derivatives on Real Samples (DR). An alternative approach is to use the learned differentiable model solely for the calculation of derivatives, with the aid of Monte-Carlo estimates obtained from *real* samples. By replacing $\nabla_a f, \nabla_s f$ in (3.2)-(3.3) with $\nabla_a \hat{f}_\psi, \nabla_s \hat{f}_\psi$ and setting the termination of backpropagation at the h -th step as $\hat{\nabla} V^{\pi_\theta}(\hat{s}_{h,n}) = \nabla \hat{V}_\omega(\hat{s}_{h,n})$, we are able to derive a dynamic representation of $\hat{\nabla}_\theta V^{\pi_\theta}$, which we defer to §A.

The corresponding RP-DR gradient estimator takes the following form,

$$\hat{\nabla}_\theta^{\text{DR}} J(\pi_\theta) = \frac{1}{N} \sum_{n=1}^N \hat{\nabla}_\theta V^{\pi_\theta}(\hat{s}_{0,n}), \quad (4.2)$$

where $\hat{s}_{0,n} \sim \mu_{\pi_\theta}$. Equation (4.2) can be specified as (A.9), which is in the same format as (4.1), but with the noise variables ς_n, ξ_n inferred from the real data sample (s_i, a_i, s_{i+1}) via the relation $a_i = \pi_\theta(s_i, \varsigma_n)$ and $s_{i+1} = \hat{f}_\psi(s_i, a_i, \xi_n)$ (see §A for details). Algorithms such as SVG (Heess et al., 2015) and its variants (Abbeel et al., 2006; Atkeson, 2012) are examples of this method.

4.3 Algorithmic Framework

The pseudocode of model-based RP PGMs is presented in Algorithm 1, where three update procedures are performed iteratively. In other words, the policy, model, and critic are updated at each iteration $t \in [T]$, generating sequences of $\{\pi_{\theta_t}\}_{t \in [T+1]}$, $\{\hat{f}_{\psi_t}\}_{t \in [T]}$, and $\{\hat{Q}_{\omega_t}\}_{t \in [T]}$, respectively.

Policy Update. The update rule for policy parameter θ with learning rate η is as follows,

$$\theta_{t+1} \leftarrow \theta_t + \eta \cdot \hat{\nabla}_\theta J(\pi_{\theta_t}), \quad (4.3)$$

where $\hat{\nabla}_\theta J(\pi_{\theta_t})$ is specified as $\hat{\nabla}_\theta^{\text{DP}} J(\pi_{\theta_t})$ or $\hat{\nabla}_\theta^{\text{DR}} J(\pi_{\theta_t})$.

Model Update. By predicting the mean of transition with minimized mean squared error (MSE) or fitting a probabilistic function with maximum likelihood estimation (MLE), canonical MBRL methods learn forward models that predict how the system evolves when an action a is taken at state s .

Algorithm 1 Model-Based RP Policy Gradient Method

Input: Number of iterations T , learning rate η , batch size N , state distribution $\mu(\cdot)$

- 1: **for** iteration $t \in [T]$ **do**
- 2: Update the model parameter ψ_t by MSE or MLE
- 3: Update the critic parameter ω_t by performing TD
- 4: Sample states from μ_{π_t} and estimate $\hat{\nabla}_\theta J(\pi_{\theta_t}) = \hat{\nabla}_\theta^{\text{DP}} J(\pi_{\theta_t})$ (4.1) or $\hat{\nabla}_\theta^{\text{DR}} J(\pi_{\theta_t})$ (4.2)
- 5: Update the policy parameter θ_t by $\theta_{t+1} \leftarrow \theta_t + \eta \cdot \hat{\nabla}_\theta J(\pi_{\theta_t})$ and execute $\pi_{\theta_{t+1}}$
- 6: **end for**
- 7: **Output:** $\{\pi_{\theta_t}\}_{t \in [T]}$

However, accurate state predictions do not imply accurate RP gradient estimation. Thus, we introduce $\epsilon_f(t)$ to represent the model (gradient) error at iteration t , defined as

$$\epsilon_f(t) = \max_{i \in [h]} \mathbb{E}_{\mathbb{P}(s_i, a_i), \mathbb{P}(\hat{s}_i, \hat{a}_i)} \left[\left\| \frac{\partial s_i}{\partial s_{i-1}} - \frac{\partial \hat{s}_i}{\partial \hat{s}_{i-1}} \right\|_2 + \left\| \frac{\partial s_i}{\partial a_{i-1}} - \frac{\partial \hat{s}_i}{\partial \hat{a}_{i-1}} \right\|_2 \right], \quad (4.4)$$

where $\mathbb{P}(s_i, a_i)$ is the true state-action distribution at the i -th timestep by following $s_0 \sim \nu_{\pi_{\theta_t}}, a_j \sim \pi_{\theta_t}(\cdot | s_j), s_{j+1} \sim f(\cdot | s_j, a_j)$, with policy and transition noise sampled from a fixed distribution. Similarly, $\mathbb{P}(\hat{s}_i, \hat{a}_i)$ is the model rollout distribution at the i -th timestep by following $\hat{s}_0 \sim \nu_{\pi_{\theta_t}}, \hat{a}_j \sim \pi_{\theta_t}(\cdot | \hat{s}_j), \hat{s}_{j+1} \sim \hat{f}_{\psi_t}(\cdot | \hat{s}_j, \hat{a}_j)$, where the noise is sampled when we use RP-DP gradient estimator and is inferred from real samples when we use RP-DR gradient estimator (in this case $\mathbb{P}(\hat{s}_i, \hat{a}_i) = \mathbb{P}(s_i, a_i)$).

In MBRL, it is common to learn a state-predictive model that can make multi-step predictions. However, this presents a challenge in reconciling the discrepancy between minimizing state prediction error and the gradient error of the model. Although it is natural to consider regularizing the models' directional derivatives to be consistent with the samples (Li et al., 2021), we contend that the use of state-predictive models does *not* cripple our analysis of gradient bias based on ϵ_f : For learned models that extrapolate beyond the visited regions, the gradient error can still be bounded via finite difference. In other words, ϵ_f can be expressed as the mean squared training error with an additional measure of the model class complexity to capture its generalizability. This same argument can also be applied to the case of learning a critic through temporal difference.

Critic Update. For any policy π , its value function Q^π satisfies the Bellman equation, which has a unique solution. In other words, $Q = \mathcal{T}^\pi Q$ if and only if $Q = Q^\pi$. The Bellman operator \mathcal{T}^π is defined for any $(s, a) \in \mathcal{S} \times \mathcal{A}$ as

$$\mathcal{T}^\pi Q(s, a) = \mathbb{E}_{\pi, f} [(1 - \gamma) \cdot r(s, a) + \gamma \cdot Q(s', a')].$$

We aim to approximate the state-action value function Q^π with a critic \hat{Q}_ω . Due to the solution uniqueness of the Bell-

man equation, it can be achieved by minimizing the mean-squared Bellman error $\mathbb{E}[(\hat{Q}_\omega(s, a) - \mathcal{T}^\pi \hat{Q}_\omega(s, a))^2]$ via Temporal Difference (TD) (Sutton, 1988; Cai et al., 2019). We define the critic error at the t -th iteration as follows,

$$\epsilon_v(t) = \alpha^2 \cdot \mathbb{E}_{\mathbb{P}(s_h, a_h), \mathbb{P}(\hat{s}_h, \hat{a}_h)} \left[\left\| \frac{\partial Q^{\pi_{\theta_t}}}{\partial s} - \frac{\partial \hat{Q}_{\omega_t}}{\partial \hat{s}} \right\|_2 + \left\| \frac{\partial Q^{\pi_{\theta_t}}}{\partial a} - \frac{\partial \hat{Q}_{\omega_t}}{\partial \hat{a}} \right\|_2 \right], \quad (4.5)$$

where $\alpha = (1 - \gamma)/\gamma^h$ and $\mathbb{P}(s_h, a_h), \mathbb{P}(\hat{s}_h, \hat{a}_h)$ are distributions at timestep h with the same definition as in (4.4). The inclusion of α^2 ensures that the critic error remains in alignment with the single-step model error ϵ_f : (1) The critic estimates the tail terms that occur after h steps in the model expansion, therefore the step-average critic error should be inversely proportional to the tail discount summation $\sum_{i=h}^\infty \gamma^i = 1/\alpha$, and (2) the quadratic form shares similarities with the canonical MBRL analysis – the cumulative error of the model trajectories scales linearly with the single-step prediction error and quadratically with the considered horizon (i.e., tail after the h -th step). This is because the cumulative error is linear in the considered horizon and the maximum state discrepancy, which is linear in the single-step error and, again, the horizon (Janner et al., 2019).

5 Main Results

In what follows, we present our main theoretical results, whose detailed proofs are deferred to §C. Specifically, we analyze the convergence of model-based RP PGMs and, more importantly, study the correlation between the convergence rate, gradient bias, variance, smoothness of the model, and approximation error. Based on our theory, we propose various algorithmic designs for model-based RP PGMs.

To begin with, we impose a common regularity condition on the policy functions following previous works (Xu et al., 2019; Pirota et al., 2015; Zhang et al., 2020; Agarwal et al., 2021). The assumption below essentially ensures the smoothness of the objective $J(\pi_\theta)$, which is required by most existing analyses of policy gradient methods (Wang et al., 2019a; Bastani, 2020; Agarwal et al., 2020).

Assumption 5.1 (Lipschitz and Bounded Score Function). We assume that the score function of policy π_θ is Lipschitz continuous and has bounded norm $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$, that is,

$$\begin{aligned} \|\log \pi_{\theta_1}(a | s) - \log \pi_{\theta_2}(a | s)\|_2 &\leq L_1 \cdot \|\theta_1 - \theta_2\|_2, \\ \|\log \pi_\theta(a | s)\|_2 &\leq B_\theta. \end{aligned}$$

We characterize the convergence of RP PGMs by first providing the following proposition.

Proposition 5.2 (Convergence to Stationary Point). We define the gradient bias b_t and variance v_t as

$$\begin{aligned} b_t &= \|\nabla_\theta J(\pi_{\theta_t}) - \mathbb{E}[\hat{\nabla}_\theta J(\pi_{\theta_t})]\|_2, \\ v_t &= \mathbb{E}[\|\hat{\nabla}_\theta J(\pi_{\theta_t}) - \mathbb{E}[\hat{\nabla}_\theta J(\pi_{\theta_t})]\|_2^2]. \end{aligned}$$

Suppose the absolute value of the reward $r(s, a)$ is bounded by $|r(s, a)| \leq r_m$ for $(s, a) \in \mathcal{S} \times \mathcal{A}$. Let $\delta = \sup \|\theta\|_2$, $L = r_m \cdot L_1/(1 - \gamma)^2 + (1 + \gamma) \cdot r_m \cdot B_\theta^2/(1 - \gamma)^3$, and $c = (\eta - L\eta^2)^{-1}$. It then holds for $T \geq 4L^2$ that

$$\begin{aligned} \min_{t \in [T]} \mathbb{E}[\|\nabla_\theta J(\pi_{\theta_t})\|_2^2] &\leq \frac{4c}{T} \cdot \mathbb{E}[J(\pi_{\theta_T}) - J(\pi_{\theta_1})] \\ &\quad + \frac{4}{T} \left(\sum_{t=0}^{T-1} c(2\delta \cdot b_t + \frac{\eta}{2} \cdot v_t) + b_t^2 + v_t \right). \end{aligned}$$

Proposition 5.2 illustrates the interdependence between the convergence and the variance, bias of the gradient estimators. In order for model-based RP PGMs to converge, it is imperative to maintain both the variance and bias at sublinear growth rates. Prior to examining the upper bound of b_t and v_t , we make the following Lipschitz assumption, which has been implemented in a plethora of preceding studies (Pirota et al., 2015; Clavera et al., 2020; Li et al., 2021).

Assumption 5.3 (Lipschitz Continuity). We assume that $r(s, a)$, $f(s, a, \xi^*)$, $\hat{f}_\psi(s, a, \xi)$, $\pi_\theta(s, \varsigma)$, $\hat{Q}_\omega(s, a)$ are $L_r, L_f, L_{\hat{f}}, L_\pi, L_{\hat{Q}}$ Lipschitz continuous (details in §B).

Let $\tilde{L}_g = \max\{L_g, 1\}$, where L_g is the Lipschitz of function g . We have the following result for gradient variance.

Proposition 5.4 (Gradient Variance). Under Assumption 5.3, for any $t \in [T]$, the gradient variance of the estimator $\hat{\nabla}_\theta J(\pi_\theta)$, which can be specified as $\hat{\nabla}_\theta^{\text{DP}} J(\pi_\theta)$ or $\hat{\nabla}_\theta^{\text{DR}} J(\pi_\theta)$, can be bounded by

$$v_t = O\left(h^4 \left(\frac{1 - \gamma^h}{1 - \gamma}\right)^2 \tilde{L}_{\hat{f}}^{4h} \tilde{L}_\pi^{4h} / N + \gamma^{2h} h^4 \tilde{L}_{\hat{f}}^{4h} \tilde{L}_\pi^{4h} / N\right).$$

We observe that the variance upper bound exhibits a polynomial dependence on the Lipschitz continuity of the model and policy, where the degrees are linear in the model unroll length. This makes sense intuitively, as the transition can be highly chaotic when $L_{\hat{f}} > 1$ and $L_\pi > 1$. This can result in diverging trajectories and variable gradient directions during training, leading to significant variance in the gradients.

Remark 5.5. Model-based RP PGMs with non-smooth models and policies can suffer from large variance and highly non-smooth loss landscapes, which can lead to slow convergence or failure during training even in simple toy examples (Parmas et al., 2018; Metz et al., 2021; Suh et al., 2022a). Proposition 5.4 suggests that one can add smoothness regularization to avoid exploding gradient variance. See our discussion at the end of this section for more details.

Model-based RP PGMs possess unique advantages by utilizing proxy models for variance reduction. By enforcing the smoothness of the model, the gradient variance is reduced without a burden when the underlying transition is smooth. However, in cases of non-smooth dynamics, such as contact-rich tasks (Suh et al., 2022a; Pang et al., 2022), doing so may introduce additional bias due to increased model

estimation error. This necessitates a trade-off between the model error and gradient variance. Nevertheless, our empirical study demonstrates that smoothness regularization improves performance, despite the cost of increased bias.

Next, we study the gradient bias. We consider the case where the state distribution μ_π , which is used for estimating the RP gradient, is a mixture of the initial distribution ζ of the MDP and the state visitation ν_π . In other words, we consider $\mu_\pi = \beta \cdot \nu_\pi + (1 - \beta) \cdot \zeta$, where $\beta \in [0, 1]$. This form is of particular interest as it encompasses various state sampling schemes that can be employed, such as when $h = 0$ and $h \rightarrow \infty$: When not utilizing a model, such as in SVG(0) (Heess et al., 2015; Amos et al., 2021) and DDPG (Lillicrap et al., 2015), states are sampled from ν_π ; while when unrolling the model over full sequences, as in BPTT, states are sampled from the initial distribution.

Given that the effects of policy actions extend to all future states and rewards, unless we know the exact policy value function, its gradient $\nabla_\theta Q^{\pi_\theta}$ cannot be simply represented by quantities in any finite timescale. Hence, differentiating through a single critic function requires extra attention, as the true value gradient has recursive structures. To tackle this issue, we provide the gradient bias bound below that is based on the measure of discrepancy between the initial distribution ζ and the state visitation ν_π .

Proposition 5.6 (Gradient Bias). We denote $\kappa = \sup_\pi \mathbb{E}_{\nu_\pi}[(d\zeta/d\nu_\pi(s))^2]^{1/2}$, where $d\zeta/d\nu_\pi$ is the Radon-Nikodym derivative of ζ with respect to ν_π . Let $\kappa' = \beta + \kappa \cdot (1 - \beta)$. Under Assumption 5.3, for any $t \in [T]$, the gradient bias is bounded by

$$b_t = O\left(\kappa\kappa'h^2(1 - \gamma^h)\tilde{L}_f^h\tilde{L}_f^h\tilde{L}_\pi^{2h}\epsilon_{f,t}/(1 - \gamma) + \kappa'h\gamma^{3h}\tilde{L}_f^h\tilde{L}_\pi^h\epsilon_{v,t}/(1 - \gamma)^2\right),$$

where $\epsilon_{f,t}$ and $\epsilon_{v,t}$ are the shorthand notations of $\epsilon_f(t)$ defined in (4.4) and $\epsilon_v(t)$ defined in (4.5), respectively.

The analysis above yields the identification of an optimal model expansion step h^* that achieves the best convergence rate, whose form is presented by the following proposition.

Proposition 5.7 (Optimal Model Expansion Step). Given $L_f \leq 1$, if we regularize the model and policy so that $L_{\hat{f}} \leq 1$ and $L_\pi \leq 1$, then when $\gamma \approx 1$, the optimal model expansion step h^* at iteration t that minimizes the convergence rate upper bound satisfies $h^* = \max\{h^*, 0\}$, where $h^* = O(\epsilon_{v,t}/((1 - \gamma)(\epsilon_{f,t} + \epsilon_{v,t})))$ scales linearly with $\epsilon_{v,t}/(\epsilon_{f,t} + \epsilon_{v,t})$ and the effective task horizon $1/(1 - \gamma)$.

In Proposition 5.7, the Lipschitz condition of the underlying dynamics, i.e., $L_f \leq 1$, ensures the stability of the system. This can be seen in the linear system example, where the transitions are determined by the eigenspectrum of the family of transformations, leading to exponential divergence of trajectories w.r.t. the largest eigenvalue. In cases where this

condition is not met in practical control systems, finding the best model unroll length may require trial and error. Fortunately, we have observed through experimentation that enforcing smoothness offers a much wider range of unrolling lengths that still provide satisfactory results.

Remark 5.8. Our analysis reveals that as the error scale $\epsilon_{v,t}/(\epsilon_{f,t} + \epsilon_{v,t})$ increases, so too does the value of h^* . This finding can inform the practical algorithms to rely more on the model by performing longer unrolls when the model error $\epsilon_{f,t}$ is small, while avoiding long unrolls when the critic error $\epsilon_{v,t}$ is small.

Finally, we characterize the algorithm convergence rate.

Corollary 5.9 (Convergence Rate). Let $\varepsilon(T) = \sum_{t=0}^{T-1} b_t$. We have for $T \geq 4L^2$ that

$$\min_{t \in [T]} \mathbb{E} \left[\|\nabla_\theta J(\pi_{\theta_t})\|_2^2 \right] \leq 16\delta \cdot \varepsilon(T)/\sqrt{T} + 4\varepsilon^2(T)/T + O(1/\sqrt{T}).$$

The convergence rate can be further clarified by determining how quickly the errors of model and critic approach zero, i.e., $\sum_{t=0}^{T-1} \epsilon_f(t) + \epsilon_v(t)$. Such results can be accomplished by conducting a more fine-grained investigation of the model and critic function classes, such as utilizing overparameterized neural nets with width scaling with T to bound the training error, as done in (Cai et al., 2019; Liu et al., 2019), and incorporating complexity measures of the model and critic function classes to bound $\epsilon_f(t)$ and $\epsilon_v(t)$. Their forms, however, are beyond the scope of this paper.

A Spectral Normalization Method. To ensure a smooth transition and faster convergence, we propose using a Spectral Normalization (SN) (Miyato et al., 2018) model-based RP PGM that applies SN to all layers of the deep model network and policy network. While other techniques, such as adversarial regularization (Shen et al., 2020), exist, we focus primarily on SN as it directly regulates the Lipschitz constant of the function. Specifically, the Lipschitz constant L_g of a function g satisfies $L_g = \sup_x \sigma_{\max}(\nabla g(x))$, where $\sigma_{\max}(W)$ denotes the largest singular value of the matrix W , defined as $\sigma_{\max}(W) = \max_{\|x\|_2 \leq 1} \|Wx\|_2$. For neural network f with linear layers $g(x) = W_i x$ and 1-Lipschitz activation (e.g., ReLU and leaky ReLU), we have $L_g = \sigma_{\max}(W_i)$ and $L_{\hat{f}} \leq \prod_i \sigma_{\max}(W_i)$. By normalizing the spectral norm of W_i with $W_i^{\text{SN}} = W_i/\sigma_{\max}(W_i)$, SN guarantees that the Lipschitz of f is upper-bounded by 1.

6 Related Work

Policy Gradient Methods. Within the RL field, the LR estimator is the basis of most policy gradient algorithms, e.g., REINFORCE (Williams, 1992) and actor-critic methods (Sutton et al., 1999; Kakade, 2001; Kakade & Langford, 2002; Degris et al., 2012). Recent works (Agarwal et al., 2021; Wang et al., 2019a; Bhandari & Russo, 2019; Liu

et al., 2019) have shown the global convergence of LR policy gradient under certain conditions, while less attention has been focused on RP PGMs. Remarkably, the analysis in (Li et al., 2021) is based on the strong assumptions on the *chained* gradient and ignores the impact of value approximation, which oversimplifies the problem by reducing the h -step model value expansion to single-step model unrolls. Besides, Clavera et al. (2020) only focused on the gradient bias while still neglecting the necessary visitation analysis. Despite the utilization in our method of Spectral Normalization on the learned model to control the gradient variance, SN has also been applied in deep RL to *value* functions in order to enable deeper neural nets (Bjorck et al., 2021) or regulate the value-aware model error (Zheng et al., 2022).

Differentiable Simulation. This paper delves into the model-based setting, where a model that captures the transition of an MDP is employed to train a control policy. Recent approaches (Mora et al., 2021; Suh et al., 2022a;b; Xu et al., 2022) based on differentiable simulators (Freeman et al., 2021; Heiden et al., 2021b) assume that gradients of simulation outcomes w.r.t. actions are explicitly given. To deal with the discontinuities and empirical bias phenomenon in the differentiable simulation caused by contact dynamics, previous works proposed using penalty-based contact formulations (Geilinger et al., 2020; Xu et al., 2021) or adopting randomized smoothing for hard-contact dynamics (Suh et al., 2022a;b). However, these are not in direct comparison to our analysis, which relies on model function approximators.

7 Experiments

7.1 Instantiations and Comparisons of RP PGMs

We begin by evaluating several algorithms originating from the RP policy gradient methods in several MuJoCo (Todorov et al., 2012) locomotion tasks in Figure 1.

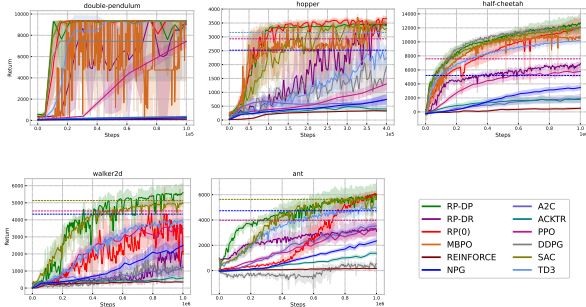


Figure 1. Evaluation of model-based RP PGMs in MuJoCo tasks. The dashed lines represent the return at convergence.

We use RP-DP and RP-DR to distinguish whether the model derivatives are calculated on predictions (4.1) or on real samples (4.2). Specifically, RP-DP is implemented as MAAC (Clavera et al., 2020) with entropy regularization, as suggested by (Amos et al., 2021); RP-DR is implemented as

SVG (Heess et al., 2015). We also evaluate model-free PGMs, including RP(0) (Amos et al., 2021), DDPG (Lillicrap et al., 2015), and its variants such as SAC (Haarnoja et al., 2018) and TD3 (Fujimoto et al., 2018).

The results indicate that RP-DP consistently outperforms or matches the performance of existing methods such as MBPO (Janner et al., 2019) and LR PGMs, including REINFORCE (Sutton et al., 1999), NPG (Kakade, 2001), ACKTR (Wu et al., 2017), and PPO (Schulman et al., 2017). This highlights the significance and potential of model-based RP PGMs. Due to space limitations, we refer readers to §D.4 for larger versions of the figures. Further implementation details and discussions can be found in §D.1 and §D.2.

7.2 Gradient Variance and Loss Landscape

Our prior investigations have revealed that vanilla model-based RP PGMs tend to have highly non-smooth landscapes due to the significant increase in gradient variance. We now conduct experiments to validate this phenomenon in practice. In Figure 2, we plot the mean gradient variance of the vanilla RP-DP algorithm during training. To visualize the loss landscapes, we plot in Figure 3 the negative value estimate along two directions that are randomly selected in the policy parameter space of a training policy.

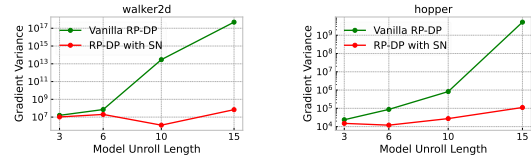


Figure 2. The gradient variance of the vanilla RP-DP explodes while adding spectral normalization solves this issue.

We can observe that for vanilla RP policy gradient algorithms, the gradient variance explodes in exponential rate with respect to the model unroll length. This results in a loss landscape that is highly non-smooth for larger unrolling steps. This renders the importance of smoothness regularization. Specifically, incorporating Spectral Normalization (SN) (Miyato et al., 2018) in the model and policy neural nets leads to a marked reduction in mean gradient variance for all unroll length settings, resulting in a much smoother loss surface compared to the vanilla implementation.

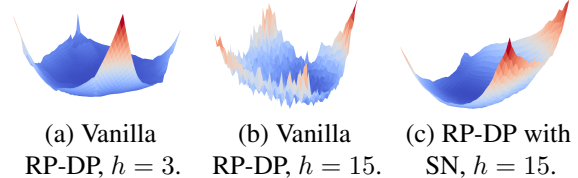


Figure 3. 2D projection of the loss surface in hopper.

7.3 Benefit of Smoothness Regularization

In this section, we investigate the effect of smoothness regularization to support our claim: The gradient variance has polynomial dependence on the Lipschitz continuity of the

model and policy, which is a contributing factor to training. Our results in Figure 4 show that SN-based RP PGMs achieve equivalent or superior performance compared to the vanilla implementation. Importantly, for longer model unrolls (e.g., 10 in walker2d and 15 in hopper), vanilla RP PGMs fail to produce reliable performance. SN-based methods, on the other hand, significantly boost training.

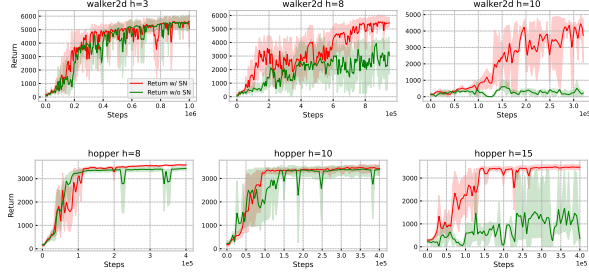


Figure 4. Performance of vanilla and SN model-based RP PGMs. Additionally, we explore different choices of model unroll lengths and examine the impact of spectral normalization, with results shown in Figure 5. We find that by utilizing SN, the curse of chaos can be mitigated, allowing for longer model unrolls. This is crucial for practical algorithmic designs: The most popular model-based RP PGMs such as (Clavera et al., 2020; Amos et al., 2021) often rely on a carefully chosen (small) h (e.g., $h = 3$). When the model is good enough, a small h may not fully leverage the accurate gradient information. As evidence, approaches (Xu et al., 2022; Mora et al., 2021) based on differentiable simulators typically adopt longer unrolls compared to model-based approaches. Therefore, with SN, more accurate multi-step predictions should enable more efficient learning without making the underlying optimization process harder. SN-based approaches also provide more robustness since the return is insensitive to h and the variance of return is smaller compared to the vanilla implementation when h is large.

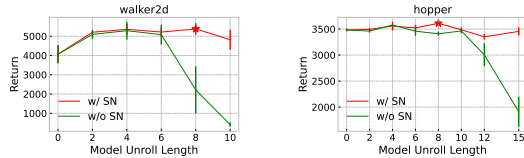


Figure 5. Spectral normalization with different model unrolls.

Ablation on Variance. By plotting the gradient variance of RP-DP during training in Figure 6, we can discern that for walker $h = 10$ and hopper $h = 15$ a key contributor to the failure of vanilla RP-DP is the exploding gradient variances. On the contrary, the SN-based approach excels in training performance as a result of the drastically reduced variance.

Ablation on Bias. When the underlying MDP is itself contact-rich and has non-smooth or even discontinuous dynamics, explicitly regularizing the Lipschitz of the transition model may lead to large error ϵ_f and thus large gradient bias. Therefore, it is also important to study if SN causes such a negative effect and if it does, how to trade off between the

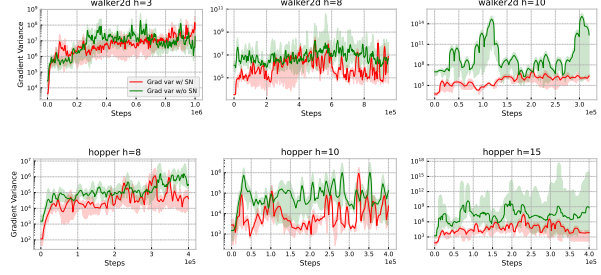


Figure 6. Gradient variance of RP PGMs during training.

model bias and gradient variance. To efficiently obtain an accurate first-order gradient (instead of via finite difference in MuJoCo), we conduct ablation based on the *differentiable* simulator dFlex (Heiden et al., 2021a; Xu et al., 2022), where Analytic Policy Gradient (APG) can be implemented.

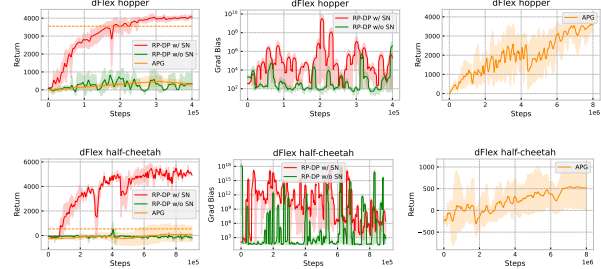


Figure 7. Performance and gradient bias in differentiable simulation. The last column is the full training curves of APG, which needs 20 times more steps in the hopper task to reach a comparable return with RP-DP-SN in the first column.

Figure 7 illustrates the crucial role that SN plays in dFlex locomotion tasks. It is noteworthy that the higher bias of the SN method does *not* impede performance, but rather improves it, indicating that the primary obstacle in training RP PGMs is the large variance in gradients. This suggests that even in differentiable simulations, one may still use a smooth proxy model when the dynamics have bumps or discontinuous jumps, sharing similarities with the gradient smoothing techniques (Suh et al., 2022a;b) applied to APG.

8 Conclusion & Future Work

In this work, we study the convergence of model-based reparameterization policy gradient methods and identify the determining factors that affect the quality of gradient estimation. Based on our theory, we propose a spectral normalization (SN) method to mitigate the exploding gradient variance issue. Our experimental results also support the proposed theory and method. Since SN-based RP PGMs allow longer model unrolls without introducing additional optimization hardness, learning more accurate multi-step models to fully leverage their gradient information should be a fruitful future direction. It will also be interesting to explore different smoothness regularization designs and apply them to a broader range of algorithms, such as using proxy models in differentiable simulation to obtain smooth policy gradients, which we would like to leave as future work.

References

- Abbeel, P., Quigley, M., and Ng, A. Y. Using inaccurate models in reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pp. 1–8, 2006.
- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. Optimality and approximation with policy gradient methods in markov decision processes. In *Conference on Learning Theory*, pp. 64–66. PMLR, 2020.
- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021.
- Amos, B., Stanton, S., Yarats, D., and Wilson, A. G. On the model-based stochastic value gradient for continuous reinforcement learning. In *Learning for Dynamics and Control*, pp. 6–20. PMLR, 2021.
- Atkeson, C. G. Efficient robust policy optimization. In *2012 American Control Conference (ACC)*, pp. 5220–5227. IEEE, 2012.
- Bastani, O. Sample complexity of estimating the policy gradient for nearly deterministic dynamical systems. In *International Conference on Artificial Intelligence and Statistics*, pp. 3858–3869. PMLR, 2020.
- Bhandari, J. and Russo, D. Global optimality guarantees for policy gradient methods. *arXiv preprint arXiv:1906.01786*, 2019.
- Bjorck, J., Gomes, C. P., and Weinberger, K. Q. Towards deeper deep reinforcement learning. *arXiv preprint arXiv:2106.01151*, 2021.
- Bollt, E. M. Controlling chaos and the inverse frobenius–perron problem: global stabilization of arbitrary invariant measures. *International Journal of Bifurcation and Chaos*, 10(05):1033–1050, 2000.
- Cai, Q., Yang, Z., Lee, J. D., and Wang, Z. Neural temporal-difference learning converges to global optima. *Advances in Neural Information Processing Systems*, 32, 2019.
- Clavera, I., Fu, V., and Abbeel, P. Model-augmented actor-critic: Backpropagating through paths. *arXiv preprint arXiv:2005.08068*, 2020.
- Degrave, J., Hermans, M., Dambre, J., et al. A differentiable physics engine for deep learning in robotics. *Frontiers in neurorobotics*, pp. 6, 2019.
- Degrís, T., White, M., and Sutton, R. S. Off-policy actor-critic. *arXiv preprint arXiv:1205.4839*, 2012.
- Duan, Y., Chen, X., Houthoofd, R., Schulman, J., and Abbeel, P. Benchmarking deep reinforcement learning for continuous control. In *International conference on machine learning*, pp. 1329–1338. PMLR, 2016.
- Feinberg, V., Wan, A., Stoica, I., Jordan, M. I., Gonzalez, J. E., and Levine, S. Model-based value expansion for efficient model-free reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning (ICML 2018)*, 2018.
- Figurnov, M., Mohamed, S., and Mnih, A. Implicit reparameterization gradients. *Advances in neural information processing systems*, 31, 2018.
- Freeman, C. D., Frey, E., Raichuk, A., Girgin, S., Mordatch, I., and Bachem, O. Brax—a differentiable physics engine for large scale rigid body simulation. *arXiv preprint arXiv:2106.13281*, 2021.
- Fujimoto, S., Hoof, H., and Meger, D. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.
- Geilinger, M., Hahn, D., Zehnder, J., Bächer, M., Thomaszewski, B., and Coros, S. Add: Analytically differentiable dynamics for multi-body systems with frictional contact. *ACM Transactions on Graphics (TOG)*, 39(6):1–15, 2020.
- Grimmett, G. and Stirzaker, D. *Probability and random processes*. Oxford university press, 2020.
- Grzeszczuk, R., Terzopoulos, D., and Hinton, G. Neuroanimator: Fast neural network emulation and control of physics-based models. In *Proceedings of the 25th annual conference on Computer graphics and interactive techniques*, pp. 9–20, 1998.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. 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.
- Heess, N., Wayne, G., Silver, D., Lillicrap, T., Erez, T., and Tassa, Y. Learning continuous control policies by stochastic value gradients. *Advances in neural information processing systems*, 28, 2015.
- Heiden, E., Macklin, M., Narang, Y., Fox, D., Garg, A., and Ramos, F. Disect: A differentiable simulation engine for autonomous robotic cutting. *arXiv preprint arXiv:2105.12244*, 2021a.
- Heiden, E., Millard, D., Coumans, E., Sheng, Y., and Sukhatme, G. S. Neursim: Augmenting differentiable

- simulators with neural networks. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 9474–9481. IEEE, 2021b.
- Janner, M., Fu, J., Zhang, M., and Levine, S. When to trust your model: Model-based policy optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Kakade, S. and Langford, J. Approximately optimal approximate reinforcement learning. In *In Proc. 19th International Conference on Machine Learning*. Citeseer, 2002.
- Kakade, S. M. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001.
- Konda, V. and Tsitsiklis, J. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- Li, C., Wang, Y., Chen, W., Liu, Y., Ma, Z.-M., and Liu, T.-Y. Gradient information matters in policy optimization by back-propagating through model. In *International Conference on Learning Representations*, 2021.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Liu, B., Cai, Q., Yang, Z., and Wang, Z. Neural trust region/proximal policy optimization attains globally optimal policy. *Advances in neural information processing systems*, 32, 2019.
- Maclaurin, D., Duvenaud, D., and Adams, R. Gradient-based hyperparameter optimization through reversible learning. In *International conference on machine learning*, pp. 2113–2122. PMLR, 2015.
- Metz, L., Maheswaranathan, N., Nixon, J., Freeman, D., and Sohl-Dickstein, J. Understanding and correcting pathologies in the training of learned optimizers. In *International Conference on Machine Learning*, pp. 4556–4565. PMLR, 2019.
- Metz, L., Freeman, C. D., Schoenholz, S. S., and Kachman, T. Gradients are not all you need. *arXiv preprint arXiv:2111.05803*, 2021.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- Mohamed, S., Rosca, M., Figurnov, M., and Mnih, A. Monte carlo gradient estimation in machine learning. *J. Mach. Learn. Res.*, 21(132):1–62, 2020.
- Mora, M. A. Z., Peychev, M. P., Ha, S., Vechev, M., and Coros, S. Pods: Policy optimization via differentiable simulation. In *International Conference on Machine Learning*, pp. 7805–7817. PMLR, 2021.
- Pang, T., Suh, H., Yang, L., and Tedrake, R. Global planning for contact-rich manipulation via local smoothing of quasi-dynamic contact models. *arXiv preprint arXiv:2206.10787*, 2022.
- Parmas, P., Rasmussen, C. E., Peters, J., and Doya, K. Pippo: Flexible model-based policy search robust to the curse of chaos. In *International Conference on Machine Learning*, pp. 4065–4074. PMLR, 2018.
- Pascanu, R., Mikolov, T., and Bengio, Y. On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pp. 1310–1318. PMLR, 2013.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Pineda, L., Amos, B., Zhang, A., Lambert, N. O., and Calandra, R. Mbrl-lib: A modular library for model-based reinforcement learning. *arXiv preprint arXiv:2104.10159*, 2021.
- Pirotta, M., Restelli, M., and Bascetta, L. Policy gradient in lipschitz markov decision processes. *Machine Learning*, 100(2):255–283, 2015.
- Rachelson, E. and Lagoudakis, M. G. On the locality of action domination in sequential decision making. 2010.
- Ruiz, F. R., AUEB, T. R., Blei, D., et al. The generalized reparameterization gradient. *Advances in neural information processing systems*, 29, 2016.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Shen, Q., Li, Y., Jiang, H., Wang, Z., and Zhao, T. Deep reinforcement learning with robust and smooth policy. In *International Conference on Machine Learning*, pp. 8707–8718. PMLR, 2020.
- Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., and Riedmiller, M. Deterministic policy gradient algorithms. In *International conference on machine learning*, pp. 387–395. PMLR, 2014.
- Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel,

- T., et al. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*, 2017.
- Suh, H., Simchowicz, M., Zhang, K., and Tedrake, R. Do differentiable simulators give better policy gradients? *arXiv preprint arXiv:2202.00817*, 2022a.
- Suh, H. J. T., Pang, T., and Tedrake, R. Bundled gradients through contact via randomized smoothing. *IEEE Robotics and Automation Letters*, 7(2):4000–4007, 2022b.
- Sutton, R. S. Learning to predict by the methods of temporal differences. *Machine learning*, 3(1):9–44, 1988.
- Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- Todorov, E., Erez, T., and Tassa, Y. MuJoCo: A physics engine for model-based control. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 5026–5033. IEEE, 2012.
- Vicol, P., Metz, L., and Sohl-Dickstein, J. Unbiased gradient estimation in unrolled computation graphs with persistent evolution strategies. In *International Conference on Machine Learning*, pp. 10553–10563. PMLR, 2021.
- Vinyals, O., Babuschkin, I., Chung, J., Mathieu, M., Jaderberg, M., Czarnecki, W. M., Dudzik, A., Huang, A., Georgiev, P., Powell, R., et al. Alphastar: Mastering the real-time strategy game starcraft ii. *DeepMind blog*, 2, 2019.
- Wang, L., Cai, Q., Yang, Z., and Wang, Z. Neural policy gradient methods: Global optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*, 2019a.
- Wang, T., Bao, X., Clavera, I., Hoang, J., Wen, Y., Langlois, E., Zhang, S., Zhang, G., Abbeel, P., and Ba, J. Benchmarking model-based reinforcement learning. *arXiv preprint arXiv:1907.02057*, 2019b.
- Weng, J. et al. A highly modularized deep reinforcement learning library. *arXiv preprint arXiv:2107.14171*, 2021.
- Werbos, P. J. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560, 1990.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- Wu, Y., Mansimov, E., Grosse, R. B., Liao, S., and Ba, J. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. *Advances in neural information processing systems*, 30, 2017.
- Xu, J., Chen, T., Zlokapa, L., Foshey, M., Matusik, W., Sueda, S., and Agrawal, P. An end-to-end differentiable framework for contact-aware robot design. *arXiv preprint arXiv:2107.07501*, 2021.
- Xu, J., Makovychuk, V., Narang, Y., Ramos, F., Matusik, W., Garg, A., and Macklin, M. Accelerated policy learning with parallel differentiable simulation. *arXiv preprint arXiv:2204.07137*, 2022.
- Xu, P., Gao, F., and Gu, Q. Sample efficient policy gradient methods with recursive variance reduction. *arXiv preprint arXiv:1909.08610*, 2019.
- Zhang, K., Koppel, A., Zhu, H., and Basar, T. Global convergence of policy gradient methods to (almost) locally optimal policies. *SIAM Journal on Control and Optimization*, 58(6):3586–3612, 2020.
- Zheng, R., Wang, X., Xu, H., and Huang, F. Is model ensemble necessary? model-based rl via a single model with lipschitz regularized value function. In *Deep Reinforcement Learning Workshop NeurIPS 2022*, 2022.

A Recursive Expression of Analytic Policy Gradient

In what follows, we interchangeably write $\nabla_a x$ and dx/da as the gradient and denote by $\partial x/\partial a$ the partial derivative.

Derivation of Analytic Policy Gradient. First of all, we provide the derivation of (3.2) and (3.3), i.e., the backward recursions of the gradient in APG.

Following (Heess et al., 2015), we define the operator

$$\nabla_\theta^i = \sum_{j \geq i} \frac{da_j}{d\theta} \cdot \frac{\partial}{\partial a_j} + \sum_{j > i} \frac{ds_j}{d\theta} \cdot \frac{\partial}{\partial s_j}. \quad (\text{A.1})$$

We begin by expanding the total derivative operator by chain rule as

$$\begin{aligned} \frac{d}{d\theta} &= \sum_{i \geq 0} \frac{da_i}{d\theta} \cdot \frac{\partial}{\partial a_i} + \sum_{i > 0} \frac{ds_i}{d\theta} \cdot \frac{\partial}{\partial s_i} \\ &= \frac{da_0}{d\theta} \cdot \frac{\partial}{\partial a_0} + \frac{ds_1}{d\theta} \cdot \frac{\partial}{\partial s_1} + \sum_{i \geq 1} \frac{da_i}{d\theta} \cdot \frac{\partial}{\partial a_i} + \sum_{i > 1} \frac{ds_i}{d\theta} \cdot \frac{\partial}{\partial s_i}. \end{aligned} \quad (\text{A.2})$$

Here, the expansion holds when $d/d\theta$ operates on policies and models that are differentiable with respect to all states s_i and actions a_i .

Plugging (A.2) into (A.1), we obtain the following recursive formula for ∇_θ^i ,

$$\nabla_\theta^i = \frac{da_i}{d\theta} \cdot \frac{\partial}{\partial a_i} + \frac{da_t}{d\theta} \cdot \frac{ds_{i+1}}{da_t} \cdot \frac{\partial}{\partial s_{i+1}} + \nabla_\theta^{i+1}. \quad (\text{A.3})$$

By the Bellman equation, we have

$$V^{\pi_\theta}(s) = \mathbb{E}_\varsigma \left[(1 - \gamma) \cdot r(s, \pi_\theta(s, \varsigma)) + \gamma \cdot \mathbb{E}_{\xi^*} \left[V^{\pi_\theta} \left(f(s, \pi_\theta(s, \varsigma), \xi^*) \right) \right] \right]. \quad (\text{A.4})$$

Combining (A.3) and (A.4) gives

$$\begin{aligned} \nabla_\theta V^{\pi_\theta}(s) &= \frac{dV^{\pi_\theta}(s)}{d\theta} = \frac{d}{d\theta} \mathbb{E}_\varsigma \left[(1 - \gamma) \cdot r(s, \pi_\theta(s, \varsigma)) + \gamma \cdot \mathbb{E}_{\xi^*} \left[V^{\pi_\theta} \left(f(s, \pi_\theta(s, \varsigma), \xi^*) \right) \right] \right] \\ &= \mathbb{E}_\varsigma \left[(1 - \gamma) \cdot \frac{\partial r}{\partial a} \cdot \frac{da}{d\theta} + \gamma \cdot \mathbb{E}_{\xi^*} \left[\frac{da}{d\theta} \cdot \frac{\partial s'}{\partial a} \cdot \frac{dV^{\pi_\theta}(s')}{ds'} + \frac{dV^{\pi_\theta}(s')}{d\theta} \right] \right], \end{aligned} \quad (\text{A.5})$$

which corresponds to (3.2).

For the $dV^{\pi_\theta}(s)/ds$ term on the right-hand side of (A.5), we have the following recursion,

$$\begin{aligned} \nabla_s V^{\pi_\theta}(s) &= \frac{dV^{\pi_\theta}(s)}{ds} = \frac{d}{ds} \mathbb{E}_\varsigma \left[(1 - \gamma) \cdot r(s, \pi_\theta(s, \varsigma)) + \gamma \cdot \mathbb{E}_{\xi^*} \left[V^{\pi_\theta} \left(f(s, \pi_\theta(s, \varsigma), \xi^*) \right) \right] \right] \\ &= \mathbb{E}_\varsigma \left[(1 - \gamma) \cdot \left(\frac{\partial r}{\partial s} + \frac{\partial r}{\partial a} \cdot \frac{\partial a}{\partial s} \right) + \gamma \cdot \mathbb{E}_{\xi^*} \left[\frac{\partial s'}{\partial s} \cdot \frac{dV^{\pi_\theta}(s')}{ds'} + \frac{\partial s'}{\partial a} \cdot \frac{\partial a}{\partial s} \cdot \frac{dV^{\pi_\theta}(s')}{ds'} \right] \right], \end{aligned} \quad (\text{A.6})$$

which corresponds to (3.3).

Therefore, we complete the derivative of (3.2) and (3.3).

Derivation of RP-DR Policy Gradient. By the same arguments in (A.5) and (A.6) with $\nabla_a f$ (or $\partial s'/\partial a$) and $\nabla_s f$ (or $\partial s'/\partial s$) replaced with $\nabla_a \hat{f}_\psi$ and $\nabla_s \hat{f}_\psi$, we obtain

$$\begin{aligned} \widehat{\nabla}_\theta V^{\pi_\theta}(\hat{s}_{i,n}) &= (1 - \gamma) \nabla_a r(\hat{s}_{i,n}, \hat{a}_{i,n}) \nabla_\theta \pi_\theta(\hat{s}_{i,n}, \varsigma_n) \\ &\quad + \gamma \widehat{\nabla}_s V^{\pi_\theta}(\hat{s}_{i+1,n}) \nabla_a \hat{f}_\psi(\hat{s}_{i,n}, \hat{a}_{i,n}, \xi_n) \nabla_\theta \pi_\theta(\hat{s}_{i,n}, \varsigma_n) + \gamma \widehat{\nabla}_\theta V^{\pi_\theta}(\hat{s}_{i+1,n}), \end{aligned} \quad (\text{A.7})$$

$$\begin{aligned} \widehat{\nabla}_s V^{\pi_\theta}(\hat{s}_{i,n}) &= (1 - \gamma) (\nabla_s r(\hat{s}_{i,n}, \hat{a}_{i,n}) + \nabla_a r(\hat{s}_{i,n}, \hat{a}_{i,n}) \nabla_s \pi_\theta(\hat{s}_{i,n}, \varsigma_n)) \\ &\quad + \gamma \widehat{\nabla}_s V^{\pi_\theta}(\hat{s}_{i+1,n}) (\nabla_s \hat{f}_\psi(\hat{s}_{i,n}, \hat{a}_{i,n}, \xi_n) + \nabla_a \hat{f}_\psi(\hat{s}_{i,n}, \hat{a}_{i,n}, \xi_n) \nabla_s \pi_\theta(\hat{s}_{i,n}, \varsigma_n)), \end{aligned} \quad (\text{A.8})$$

where the termination of backpropagation at the h -th step is $\widehat{\nabla} V^{\pi_\theta}(\hat{s}_{h,n}) = \nabla \hat{V}_\omega(\hat{s}_{h,n})$.

Combining (A.7) and (A.8), we obtain the RP-DR policy gradient estimator as follows,

$$\widehat{\nabla}_{\theta}^{\text{DR}} J(\pi_{\theta}) = \frac{1}{N} \sum_{n=1}^N \widehat{\nabla}_{\theta} V^{\pi_{\theta}}(\widehat{s}_{0,n}) = \frac{1}{N} \sum_{n=1}^N \nabla_{\theta} \left(\sum_{i=0}^{h-1} \gamma^i r(\widehat{s}_{i,n}, \widehat{a}_{i,n}) + \gamma^h \widehat{Q}_{\omega}(\widehat{s}_{h,n}, \widehat{a}_{h,n}) \right), \quad (\text{A.9})$$

where $\widehat{s}_{0,n} \sim \mu_{\pi_{\theta}}$, $\widehat{a}_{i,n} = \pi_{\theta}(\widehat{s}_{i,n}, \varsigma_n)$, and $\widehat{s}_{i+1,n} = \widehat{f}_{\psi}(\widehat{s}_{i,n}, \widehat{a}_{i,n}, \xi_n)$. Here, ς_n and ξ_n are inferred by solving $a_i = \pi_{\theta}(s_i, \varsigma_n)$ and $s_{i+1} = \widehat{f}_{\psi}(s_i, a_i, \xi_n)$, respectively, where (s_i, a_i, s_{i+1}) is the real data sample. For example, for a state s_{i+1} sampled from a one-dimensional Gaussian transition model $s_{i+1} \sim \mathcal{N}(\phi(s_i, a_i), \sigma^2)$, where the variance is σ and the mean $\phi(s_i, a_i)$ is the output of some function parameterized by ϕ , the noise ξ_n can be inferred as $\xi_n = (s_{i+1} - \phi(s_i, a_i))/\sigma$.

B Complete Statement of Assumption 5.3

Assumption B.1 (Lipschitz Continuous Functions). We assume that $r(s, a)$, $f(s, a, \xi^*)$, $\widehat{f}_{\psi}(s, a, \xi)$, $\pi_{\theta}(s, \varsigma)$, $\widehat{Q}_{\omega}(s, a)$ are $L_r, L_f, L_{\widehat{f}}, L_{\pi}, L_{\widehat{Q}}$ -Lipschitz continuous, respectively, that is, for any ψ, θ, ω ,

$$\begin{aligned} |r(s_1, a_1) - r(s_2, a_2)| &\leq L_r \cdot \|(s_1 - s_2, a_1 - a_2)\|_2, \\ \|f(s_1, a_1, \xi_1^*) - f(s_2, a_2, \xi_2^*)\|_2 &\leq L_f \cdot \|(s_1 - s_2, a_1 - a_2, \xi_1^* - \xi_2^*)\|_2, \\ \|\widehat{f}_{\psi}(s_1, a_1, \xi_1) - \widehat{f}_{\psi}(s_2, a_2, \xi_2)\|_2 &\leq L_{\widehat{f}} \cdot \|(s_1 - s_2, a_1 - a_2, \xi_1 - \xi_2)\|_2, \\ \|\pi_{\theta}(s_1, \varsigma_1) - \pi_{\theta}(s_2, \varsigma_2)\|_2 &\leq L_{\pi} \cdot \|(s_1 - s_2, \varsigma_1 - \varsigma_2)\|_2, \\ |\widehat{Q}_{\omega}(s_1, a_1) - \widehat{Q}_{\omega}(s_2, a_2)| &\leq L_{\widehat{Q}} \cdot \|(s_1 - s_2, a_1 - a_2)\|_2. \end{aligned}$$

Additionally, we assume that the policy $\pi_{\theta}(s, \varsigma)$ is L_{θ} -Lipschitz continuous in θ , which implies $\|\nabla_{\theta} \pi_{\theta}(s, \varsigma)\|_2 \leq L_{\theta}$ for any state $s \in \mathcal{S}$.

C Proofs

C.1 Proof of Proposition 5.2

As a preparation before proving Proposition 5.2, we first present the following lemma stating that the objective in (2.1) is Lipschitz smooth under Assumption 5.1.

Lemma C.1 (Smooth Objective, (Zhang et al., 2020) Lemma 3.2). The objective $J(\pi_{\theta})$ is L -smooth in θ , such that $\|\nabla_{\theta} J(\pi_{\theta_1}) - \nabla_{\theta} J(\pi_{\theta_2})\|_2 \leq L \|\theta_1 - \theta_2\|_2$, where

$$L = \frac{r_m \cdot L_1}{(1 - \gamma)^2} + \frac{(1 + \gamma) \cdot r_m \cdot B_{\theta}^2}{(1 - \gamma)^3}.$$

Then we are ready to prove Proposition 5.2.

Proof of Proposition 5.2. From the policy update rule in (4.3), we have $\widehat{\nabla}_{\theta} J(\pi_{\theta_t}) = (\theta_{t+1} - \theta_t)/\eta$. By Assumption 5.3, we have

$$\begin{aligned} J(\pi_{\theta_{t+1}}) - J(\pi_{\theta_t}) &\geq \nabla_{\theta} J(\pi_{\theta_t})^{\top} (\theta_{t+1} - \theta_t) - \frac{L}{2} \|\theta_{t+1} - \theta_t\|_2^2 \\ &= \eta \nabla_{\theta} J(\pi_{\theta_t})^{\top} \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) - \frac{L\eta^2}{2} \|\widehat{\nabla}_{\theta} J(\pi_{\theta_t})\|_2^2. \end{aligned} \quad (\text{C.1})$$

By basic algebra, we have for $\nabla_{\theta} J(\pi_{\theta_t})^{\top} \widehat{\nabla}_{\theta} J(\pi_{\theta_t})$ that

$$\begin{aligned} &\nabla_{\theta} J(\pi_{\theta_t})^{\top} \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \\ &= \left(\nabla_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] \right)^{\top} \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) - \left(\widehat{\nabla}_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] \right)^{\top} \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) + \widehat{\nabla}_{\theta} J(\pi_{\theta_t})^{\top} \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \\ &\geq - \underbrace{\left| \left(\nabla_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] \right)^{\top} \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right|}_{(\text{I})} - \underbrace{\left| \left(\widehat{\nabla}_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] \right)^{\top} \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right|}_{(\text{II})} + \underbrace{\widehat{\nabla}_{\theta} J(\pi_{\theta_t})^{\top} \widehat{\nabla}_{\theta} J(\pi_{\theta_t})}_{(\text{III})}. \end{aligned}$$

The resulting three terms can be bounded as follows,

$$\begin{aligned}
 \text{Term (I):} \quad (\text{I}) &\leq \left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2 \cdot \left\| \nabla_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] \right\|_2 = \left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2 \cdot b_t, \\
 \text{Term (II):} \quad (\text{II}) &\leq \frac{\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] \right\|_2^2}{2} + \frac{\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2}{2}, \\
 \text{Term (III):} \quad (\text{III}) &\geq \left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2.
 \end{aligned}$$

Thus, by plugging the above three inequalities into (C.1), we have

$$\begin{aligned}
 J(\pi_{\theta_{t+1}}) - J(\pi_{\theta_t}) &\geq \frac{\eta}{2} \cdot \left(-\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2 \cdot 2b_t - \left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] \right\|_2^2 + \left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2 \right) \\
 &\quad - \frac{L\eta^2}{2} \cdot \left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2.
 \end{aligned} \tag{C.2}$$

By taking expectation on both sides of (C.2), we obtain

$$\mathbb{E}[J(\pi_{\theta_{t+1}}) - J(\pi_{\theta_t})] \geq -\eta \cdot \mathbb{E}[\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2] \cdot b_t - \frac{\eta}{2} \cdot v_t + \frac{\eta - L\eta^2}{2} \cdot \mathbb{E}[\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2].$$

Rearranging terms gives

$$\frac{\eta - L\eta^2}{2} \cdot \mathbb{E}[\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2] \leq \mathbb{E}[J(\pi_{\theta_{t+1}}) - J(\pi_{\theta_t})] + \eta \mathbb{E}[\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2] b_t + \frac{\eta}{2} v_t. \tag{C.3}$$

By establishing the connection between the minimum expected gradient norm and the average norm over T iterations, we are able to obtain the following bound,

$$\begin{aligned}
 \min_{t \in [T]} \mathbb{E}[\left\| \nabla_{\theta} J(\pi_{\theta_t}) \right\|_2^2] &\leq \frac{1}{T} \cdot \sum_{t=0}^{T-1} \mathbb{E}[\left\| \nabla_{\theta} J(\pi_{\theta_t}) \right\|_2^2] \\
 &\leq \frac{2}{T} \cdot \sum_{t=0}^{T-1} \left(\mathbb{E}[\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2] + \mathbb{E}[\left\| \nabla_{\theta} J(\pi_{\theta_t}) - \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2] \right),
 \end{aligned} \tag{C.4}$$

where the second inequality holds since for any vector $y, z \in \mathbb{R}^d$,

$$\|y + z\|_2^2 \leq \|y\|_2^2 + \|z\|_2^2 + 2\|y\|_2 \cdot \|z\|_2 \leq 2\|y\|_2^2 + 2\|z\|_2^2. \tag{C.5}$$

The last term on the right-hand side of (C.4) can be characterized by

$$\begin{aligned}
 \mathbb{E}[\left\| \nabla_{\theta} J(\pi_{\theta_t}) - \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2] &= \mathbb{E}[\left\| \nabla_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] + \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] - \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2^2] \\
 &\leq 2\left\| \nabla_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] \right\|_2^2 + 2\mathbb{E}[\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_{\theta} J(\pi_{\theta_t})] \right\|_2^2] \\
 &= 2b_t^2 + 2v_t,
 \end{aligned} \tag{C.6}$$

where the inequality follows from (C.5).

For $T \geq 4L^2$, by setting $\eta = 1/\sqrt{T}$, we have $\eta < 1/L$ and $(\eta - L\eta^2)/2 > 0$. Therefore, following the results in (C.3) and (C.6), we further have

$$\begin{aligned}
 &\min_{t \in [T]} \mathbb{E}[\left\| \nabla_{\theta} J(\pi_{\theta_t}) \right\|_2^2] \\
 &\leq \frac{4c}{T} \cdot \left(\mathbb{E}[J(\pi_{\theta_T}) - J(\pi_{\theta_1})] + \sum_{t=0}^{T-1} \left(\eta \cdot \mathbb{E}[\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2] \cdot b_t + \frac{\eta}{2} \cdot v_t \right) \right) + \frac{4}{T} \cdot \sum_{t=0}^{T-1} (b_t^2 + v_t) \\
 &= \frac{4}{T} \cdot \left(\sum_{t=0}^{T-1} c \cdot \left(\eta \cdot \mathbb{E}[\left\| \widehat{\nabla}_{\theta} J(\pi_{\theta_t}) \right\|_2] \cdot b_t + \frac{\eta}{2} \cdot v_t \right) + b_t^2 + v_t \right) + \frac{4c}{T} \cdot \mathbb{E}[J(\pi_{\theta_T}) - J(\pi_{\theta_1})],
 \end{aligned}$$

where the last step holds due to the definition $c = (\eta - L\eta^2)^{-1}$.

By noting that $\eta \widehat{\nabla}_\theta J(\pi_{\theta_t}) = \theta_{t+1} - \theta_t$, we conclude the proof by

$$\begin{aligned} & \min_{t \in [T]} \mathbb{E} \left[\left\| \nabla_\theta J(\pi_{\theta_t}) \right\|_2^2 \right] \\ & \leq \frac{4}{T} \cdot \left(\sum_{t=0}^{T-1} c \cdot \left(\mathbb{E} \left[\left\| \theta_{t+1} - \theta_t \right\|_2 \right] \cdot b_t + \frac{\eta}{2} \cdot v_t \right) + b_t^2 + v_t \right) + \frac{4c}{T} \cdot \mathbb{E} [J(\pi_{\theta_T}) - J(\pi_{\theta_1})] \\ & \leq \frac{4}{T} \cdot \left(\sum_{t=0}^{T-1} c \cdot (2\delta \cdot b_t + \frac{\eta}{2} \cdot v_t) + b_t^2 + v_t \right) + \frac{4c}{T} \cdot \mathbb{E} [J(\pi_{\theta_T}) - J(\pi_{\theta_1})]. \end{aligned}$$

where the second inequality holds since $\|\theta\|_2 \leq \delta$ for any θ . Therefore, we conclude the proof of Proposition 5.2. \square

C.2 Proof of Proposition 5.4

Proof. Since the RP-DP gradient $\widehat{\nabla}_\theta J(\pi_\theta) = \widehat{\nabla}_\theta^{\text{DP}} J(\pi_\theta)$ in (4.1) and the RP-DR gradient $\widehat{\nabla}_\theta J(\pi_\theta) = \widehat{\nabla}_\theta^{\text{DR}} J(\pi_\theta)$ in (4.2) share the same state transition $\widehat{s}_{i+1,n} = \widehat{f}(\widehat{s}_{i,n}, \xi_n)$, where recall that the only difference lies in the source of noise ξ_n , our subsequent analysis holds for both RP-DP and RP-DR.

To upper-bound the gradient variance $v_t = \mathbb{E}[\|\widehat{\nabla}_\theta J(\pi_{\theta_t}) - \mathbb{E}[\widehat{\nabla}_\theta J(\pi_{\theta_t})]\|_2^2]$, we characterize the norm inside the outer expectation.

We start with the case where the sample size $N = 1$, which naturally generalizes to $N > 1$. Specifically, we consider an *arbitrary* trajectory obtained by unrolling the model under policy π_{θ_t} . We denote the pathwise gradient $\widehat{\nabla}_\theta J(\pi_{\theta_t})$ of this trajectory as g' . Then we have

$$v_t \leq \max_{g'} \left\| g' - \mathbb{E}[\widehat{\nabla}_\theta J(\pi_{\theta_t})] \right\|_2^2 = \left\| g - \mathbb{E}[\widehat{\nabla}_\theta J(\pi_{\theta_t})] \right\|_2^2 = \left\| \mathbb{E}[g - \widehat{\nabla}_\theta J(\pi_{\theta_t})] \right\|_2^2,$$

where g is the pathwise gradient $\widehat{\nabla}_\theta J(\pi_{\theta_t})$ of a *fixed* (but unknown) trajectory $(\widehat{s}_{0,n}, \widehat{a}_{0,n}, \widehat{s}_{1,n}, \widehat{a}_{1,n}, \dots)$ such that the maximum is achieved.

Using the fact that $\|\mathbb{E}[\cdot]\|_2 \leq \mathbb{E}[\|\cdot\|_2]$, we further obtain

$$v_t \leq \mathbb{E} \left[\left\| g - \widehat{\nabla}_\theta J(\pi_{\theta_t}) \right\|_2^2 \right]. \quad (\text{C.7})$$

Let $\widehat{x}_{i,n} = (\widehat{s}_{i,n}, \widehat{a}_{i,n})$. By the triangular inequality, we have

$$\begin{aligned} \mathbb{E} \left[\left\| g - \widehat{\nabla}_\theta J(\pi_{\theta_t}) \right\|_2 \right] & \leq \sum_{i=0}^{h-1} \gamma^i \cdot \mathbb{E}_{\widehat{x}_i} \left[\left\| \nabla_\theta r(\widehat{x}_{i,n}) - \nabla_\theta r(\overline{x}_i) \right\|_2 \right] \\ & \quad + \gamma^h \cdot \mathbb{E}_{\widehat{x}_h} \left[\left\| \nabla \widehat{Q}(\widehat{x}_{h,n}) \nabla_\theta \widehat{x}_{h,n} - \nabla \widehat{Q}(\overline{x}_h) \nabla_\theta \overline{x}_h \right\|_2 \right]. \end{aligned} \quad (\text{C.8})$$

By the chain rule, we have for any $i \geq 1$ that

$$\frac{d\widehat{a}_{i,n}}{d\theta} = \frac{\partial \widehat{a}_{i,n}}{\partial \widehat{s}_{i,n}} \cdot \frac{d\widehat{s}_{i,n}}{d\theta} + \frac{\partial \widehat{a}_{i,n}}{\partial \theta}, \quad (\text{C.9})$$

$$\frac{d\widehat{s}_{i,n}}{d\theta} = \frac{\partial \widehat{s}_{i,n}}{\partial \widehat{s}_{i-1,n}} \cdot \frac{d\widehat{s}_{i-1,n}}{d\theta} + \frac{\partial \widehat{s}_{i,n}}{\partial \widehat{a}_{i-1,n}} \cdot \frac{d\widehat{a}_{i-1,n}}{d\theta}. \quad (\text{C.10})$$

Plugging $d\widehat{a}_{i-1,n}/d\theta$ in (C.9) into (C.10), we get

$$\begin{aligned} \left\| \frac{d\widehat{s}_{i,n}}{d\theta} \right\|_2 & = \left\| \left(\frac{\partial \widehat{s}_{i,n}}{\partial \widehat{s}_{i-1,n}} + \frac{\partial \widehat{s}_{i,n}}{\partial \widehat{a}_{i-1,n}} \cdot \frac{\partial \widehat{a}_{i-1,n}}{\partial \widehat{s}_{i-1,n}} \right) \cdot \frac{d\widehat{s}_{i-1,n}}{d\theta} + \frac{\partial \widehat{s}_{i,n}}{\partial \widehat{a}_{i-1,n}} \cdot \frac{\partial \widehat{a}_{i-1,n}}{\partial \theta} \right\|_2 \\ & \leq L_{\widehat{f}} \widetilde{L}_\pi \cdot \left\| \frac{d\widehat{s}_{i-1,n}}{d\theta} \right\|_2 + L_{\widehat{f}} L_\theta, \end{aligned} \quad (\text{C.11})$$

where the inequality follows from Assumption 5.3 and the Cauchy-Schwarz inequality.

Recursively applying (C.11), we obtain for any $i \geq 1$ that

$$\left\| \frac{d\widehat{s}_{i,n}}{d\theta} \right\|_2 \leq L_{\widehat{f}} L_\theta \cdot \sum_{j=0}^{i-1} L_{\widehat{f}}^j \widetilde{L}_\pi^j \leq i \cdot L_\theta L_{\widehat{f}}^{i+1} \widetilde{L}_\pi^i, \quad (\text{C.12})$$

where the first inequality follows from the induction

$$z_i = az_{i-1} + b = a \cdot (az_{i-2} + b) + b = a^i \cdot z_0 + b \cdot \sum_{j=0}^{i-1} a^j. \quad (\text{C.13})$$

In (C.13), $\{z_j\}_{0 \leq j \leq i}$ is the real sequence satisfying $z_j = az_{j-1} + b$. For $d\hat{a}_{i,n}/d\theta$ defined in (C.9), we further have

$$\left\| \frac{d\hat{a}_{i,n}}{d\theta} \right\|_2 \leq L_\pi \cdot \left\| \frac{d\hat{s}_{i,n}}{d\theta} \right\|_2 + L_\theta \leq i \cdot L_\theta L_{\hat{f}}^{i+1} \tilde{L}_\pi^{i+1} + L_\theta. \quad (\text{C.14})$$

Combining (C.12) and (C.14), we obtain

$$\left\| \frac{d\hat{x}_{i,n}}{d\theta} \right\|_2 = \left\| \frac{d\hat{s}_{i,n}}{d\theta} \right\|_2 + \left\| \frac{d\hat{a}_{i,n}}{d\theta} \right\|_2 \leq \underbrace{2i \cdot L_\theta L_{\hat{f}}^{i+1} \tilde{L}_\pi^{i+1} + L_\theta}_{\hat{K}(i)}. \quad (\text{C.15})$$

Therefore, we bound the second term on the right-hand side of (C.8) as follows,

$$\begin{aligned} & \mathbb{E}_{\bar{x}_h} \left[\left\| \nabla \hat{Q}(\hat{x}_{h,n}) \nabla_\theta \hat{x}_{h,n} - \nabla \hat{Q}(\bar{x}_h) \nabla_\theta \bar{x}_h \right\|_2 \right] \\ & \leq \mathbb{E}_{\bar{x}_h} \left[\left\| \nabla \hat{Q}(\hat{x}_{h,n}) \nabla_\theta \hat{x}_{h,n} - \nabla \hat{Q}(\bar{x}_h) \nabla_\theta \hat{x}_{h,n} \right\|_2 \right] + \mathbb{E}_{\bar{x}_h} \left[\left\| \nabla \hat{Q}(\bar{x}_h) \nabla_\theta \hat{x}_{h,n} - \nabla \hat{Q}(\bar{x}_h) \nabla_\theta \bar{x}_h \right\|_2 \right] \\ & \leq 2L_{\hat{Q}} \cdot \hat{K}(i) + L_{\hat{Q}} \cdot \left(\mathbb{E}_{\bar{s}_i} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} - \frac{d\bar{s}_i}{d\theta} \right\|_2 \right] + \mathbb{E}_{\bar{a}_i} \left[\left\| \frac{d\hat{a}_{i,n}}{d\theta} - \frac{d\bar{a}_i}{d\theta} \right\|_2 \right] \right), \end{aligned} \quad (\text{C.16})$$

where the last inequality follows from the Cauchy-Schwartz inequality and Assumption 5.3.

By the chain rule, we bound the first term on the right-hand side of (C.8) as follows,

$$\begin{aligned} & \mathbb{E}_{\bar{x}_i} \left[\left\| \nabla_\theta r(\hat{x}_{i,n}) - \nabla_\theta r(\bar{x}_i) \right\|_2 \right] \\ & = \mathbb{E}_{\bar{x}_i} \left[\left\| \nabla r(\hat{x}_{i,n}) \nabla_\theta \hat{x}_{i,n} - \nabla r(\bar{x}_i) \nabla_\theta \bar{x}_i \right\|_2 \right] \\ & \leq \mathbb{E}_{\bar{x}_i} \left[\left\| \nabla r(\hat{x}_{i,n}) \nabla_\theta \hat{x}_{i,n} - \nabla r(\hat{x}_{i,n}) \nabla_\theta \bar{x}_i \right\|_2 \right] + \mathbb{E}_{\bar{x}_i} \left[\left\| \nabla r(\hat{x}_{i,n}) \nabla_\theta \bar{x}_i - \nabla r(\bar{x}_i) \nabla_\theta \bar{x}_i \right\|_2 \right] \\ & \leq L_r \cdot \left(\mathbb{E}_{\bar{s}_i} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} - \frac{d\bar{s}_i}{d\theta} \right\|_2 \right] + \mathbb{E}_{\bar{a}_i} \left[\left\| \frac{d\hat{a}_{i,n}}{d\theta} - \frac{d\bar{a}_i}{d\theta} \right\|_2 \right] \right) + 2L_r \cdot \hat{K}(i). \end{aligned} \quad (\text{C.17})$$

Plugging (C.16) and (C.17) into (C.8) and (C.7), we obtain

$$\begin{aligned} v_t & \leq \left[\left(L_r \cdot \sum_{i=0}^{h-1} \gamma^i + \gamma^h \cdot L_{\hat{Q}} \right) \cdot \left(\mathbb{E}_{\bar{s}_h} \left[\left\| \frac{d\hat{s}_{h,n}}{d\theta} - \frac{d\bar{s}_h}{d\theta} \right\|_2 \right] + \mathbb{E}_{\bar{a}_h} \left[\left\| \frac{d\hat{a}_{h,n}}{d\theta} - \frac{d\bar{a}_h}{d\theta} \right\|_2 \right] + 2\hat{K}(h) \right) \right]^2 \\ & = O \left(h^4 \left(\frac{1-\gamma^h}{1-\gamma} \right)^2 \tilde{L}_{\hat{f}}^{4h} \tilde{L}_\pi^{4h} + \gamma^{2h} h^4 \tilde{L}_{\hat{f}}^{4h} \tilde{L}_\pi^{4h} \right), \end{aligned} \quad (\text{C.18})$$

where the inequality follows from Lemma C.2 and by plugging the definition of $\hat{K}(i)$ in (C.15).

Note that the variance v_t scales with the batch size N at the rate of $1/N$. Since the analysis above is established for $N = 1$, the bound of the gradient variance v_t is established by dividing the right-hand side of (C.18) by N , which concludes the proof of Proposition 5.4. \square

Lemma C.2. Denote $e = \sup \mathbb{E}_{\bar{s}_0} [\|d\hat{s}_{0,n}/d\theta - d\bar{s}_0/d\theta\|_2]$, which is a constant that only depends on the initial state distribution¹. For any timestep $i \geq 1$ and the corresponding state, action, we have the following results,

$$\begin{aligned} \mathbb{E}_{\bar{s}_i} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} - \frac{d\bar{s}_i}{d\theta} \right\|_2 \right] & \leq \tilde{L}_{\hat{f}}^i \tilde{L}_\pi^i \left(e + 4i \cdot \tilde{L}_{\hat{f}} \tilde{L}_\pi \cdot \hat{K}(i-1) + 2i \cdot \tilde{L}_{\hat{f}} L_\theta \right), \\ \mathbb{E}_{\bar{a}_i} \left[\left\| \frac{d\hat{a}_{i,n}}{d\theta} - \frac{d\bar{a}_i}{d\theta} \right\|_2 \right] & \leq \tilde{L}_{\hat{f}}^i \tilde{L}_\pi^{i+1} \left(e + 4i \cdot \tilde{L}_{\hat{f}} \tilde{L}_\pi \cdot \hat{K}(i-1) + 2i \cdot \tilde{L}_{\hat{f}} L_\theta \right) + 2L_\pi \hat{K}(i) + 2L_\theta. \end{aligned}$$

¹We define e to account for the stochasticity of the initial state distribution. $e = 0$ when the initial state is deterministic.

Proof. Firstly, from (C.10), we obtain for any $i \geq 1$ that

$$\begin{aligned}
 & \mathbb{E}_{\bar{s}_i} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} - \frac{d\bar{s}_i}{d\theta} \right\|_2 \right] \\
 &= \mathbb{E} \left[\left\| \frac{\partial \hat{s}_{i,n}}{\partial \hat{s}_{i-1,n}} \cdot \frac{d\hat{s}_{i-1,n}}{d\theta} + \frac{\partial \hat{s}_{i,n}}{\partial \hat{a}_{i-1,n}} \cdot \frac{d\hat{a}_{i-1,n}}{d\theta} - \frac{\partial \bar{s}_i}{\partial \bar{s}_{i-1}} \cdot \frac{d\bar{s}_{i-1}}{d\theta} - \frac{\partial \bar{s}_i}{\partial \bar{a}_{i-1}} \cdot \frac{d\bar{a}_{i-1}}{d\theta} \right\|_2 \right] \\
 &\text{According to the triangle inequality, we further have} \\
 &\leq \mathbb{E} \left[\left\| \frac{\partial \hat{s}_{i,n}}{\partial \hat{s}_{i-1,n}} \cdot \frac{d\hat{s}_{i-1,n}}{d\theta} - \frac{\partial \bar{s}_i}{\partial \bar{s}_{i-1}} \cdot \frac{d\bar{s}_{i-1,n}}{d\theta} \right\|_2 \right] + \mathbb{E} \left[\left\| \frac{\partial \bar{s}_i}{\partial \bar{s}_{i-1}} \cdot \frac{d\hat{s}_{i-1,n}}{d\theta} - \frac{\partial \bar{s}_i}{\partial \bar{s}_{i-1}} \cdot \frac{d\bar{s}_{i-1}}{d\theta} \right\|_2 \right] \\
 &\quad + \mathbb{E} \left[\left\| \frac{\partial \hat{s}_{i,n}}{\partial \hat{a}_{i-1,n}} \cdot \frac{d\hat{a}_{i-1,n}}{d\theta} - \frac{\partial \bar{s}_i}{\partial \bar{a}_{i-1}} \cdot \frac{d\hat{a}_{i-1,n}}{d\theta} \right\|_2 \right] + \mathbb{E} \left[\left\| \frac{\partial \bar{s}_i}{\partial \bar{a}_{i-1}} \cdot \frac{d\hat{a}_{i-1,n}}{d\theta} - \frac{\partial \bar{s}_i}{\partial \bar{a}_{i-1}} \cdot \frac{d\bar{a}_{i-1}}{d\theta} \right\|_2 \right] \\
 &\leq 2L_{\hat{f}} \cdot \left(\left\| \frac{d\hat{s}_{i-1,n}}{d\theta} \right\|_2 + \left\| \frac{d\hat{a}_{i-1,n}}{d\theta} \right\|_2 \right) + L_{\hat{f}} \cdot \mathbb{E}_{\bar{s}_{i-1}} \left[\left\| \frac{d\hat{s}_{i-1,n}}{d\theta} - \frac{d\bar{s}_{i-1}}{d\theta} \right\|_2 \right] \\
 &\quad + L_{\hat{f}} \cdot \mathbb{E}_{\bar{a}_{i-1}} \left[\left\| \frac{d\hat{a}_{i-1,n}}{d\theta} - \frac{d\bar{a}_{i-1}}{d\theta} \right\|_2 \right]. \tag{C.19}
 \end{aligned}$$

Similarly, we have from (C.9) that

$$\begin{aligned}
 & \mathbb{E}_{\bar{a}_i} \left[\left\| \frac{d\hat{a}_{i,n}}{d\theta} - \frac{d\bar{a}_i}{d\theta} \right\|_2 \right] \\
 &= \mathbb{E} \left[\left\| \frac{\partial \hat{a}_{i,n}}{\partial \hat{s}_{i,n}} \cdot \frac{d\hat{s}_{i,n}}{d\theta} + \frac{\partial \hat{a}_{i,n}}{\partial \theta} - \frac{\partial \bar{a}_i}{\partial \bar{s}_i} \cdot \frac{d\bar{s}_i}{d\theta} - \frac{\partial \bar{a}_i}{\partial \theta} \right\|_2 \right] \\
 &\leq \mathbb{E} \left[\left\| \frac{\partial \hat{a}_{i,n}}{\partial \hat{s}_{i,n}} \cdot \frac{d\hat{s}_{i,n}}{d\theta} - \frac{\partial \bar{a}_i}{\partial \bar{s}_i} \cdot \frac{d\bar{s}_{i,n}}{d\theta} \right\|_2 \right] + \mathbb{E} \left[\left\| \frac{\partial \bar{a}_i}{\partial \bar{s}_i} \cdot \frac{d\hat{s}_{i,n}}{d\theta} - \frac{\partial \bar{a}_i}{\partial \bar{s}_i} \cdot \frac{d\bar{s}_i}{d\theta} \right\|_2 \right] + \mathbb{E} \left[\left\| \frac{\partial \hat{a}_{i,n}}{\partial \theta} - \frac{\partial \bar{a}_i}{\partial \theta} \right\|_2 \right] \\
 &\leq 2L_{\pi} \cdot \mathbb{E} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} \right\|_2 \right] + L_{\pi} \cdot \mathbb{E} \left[\left\| \frac{d\bar{s}_{i,n}}{d\theta} - \frac{d\bar{s}_i}{d\theta} \right\|_2 \right] + 2L_{\theta}. \tag{C.20}
 \end{aligned}$$

Plugging (C.20) back to (C.19), we obtain

$$\begin{aligned}
 & \mathbb{E}_{\bar{s}_i} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} - \frac{d\bar{s}_i}{d\theta} \right\|_2 \right] \\
 &\lesssim 4L_{\hat{f}}\tilde{L}_{\pi} \cdot \left(\left\| \frac{d\hat{s}_{i-1,n}}{d\theta} \right\|_2 + \left\| \frac{d\hat{a}_{i-1,n}}{d\theta} \right\|_2 \right) + L_{\hat{f}}\tilde{L}_{\pi} \cdot \mathbb{E}_{\bar{s}_{i-1}} \left[\left\| \frac{d\hat{s}_{i-1,n}}{d\theta} - \frac{d\bar{s}_{i-1}}{d\theta} \right\|_2 \right] + 2L_{\hat{f}}L_{\theta} \\
 &\leq 4L_{\hat{f}}\tilde{L}_{\pi} \cdot \hat{K}(i-1) + L_{\hat{f}}\tilde{L}_{\pi} \cdot \mathbb{E}_{\bar{s}_{i-1}} \left[\left\| \frac{d\hat{s}_{i-1,n}}{d\theta} - \frac{d\bar{s}_{i-1}}{d\theta} \right\|_2 \right] + 2L_{\hat{f}}L_{\theta},
 \end{aligned}$$

where the last inequality follows from the definition of \hat{K} in (C.15).

Applying this recursion gives

$$\begin{aligned}
 \mathbb{E}_{\bar{s}_i} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} - \frac{d\bar{s}_i}{d\theta} \right\|_2 \right] &\leq e(L_{\hat{f}}\tilde{L}_{\pi})^i + (4L_{\hat{f}}\tilde{L}_{\pi} \cdot \hat{K}(i-1) + 2L_{\hat{f}}L_{\theta}) \cdot \sum_{j=0}^{i-1} (L_{\hat{f}}\tilde{L}_{\pi})^j \\
 &\leq \tilde{L}_{\hat{f}}\tilde{L}_{\pi}^i \left(e + 4i \cdot \tilde{L}_{\hat{f}}\tilde{L}_{\pi} \cdot \hat{K}(i-1) + 2i \cdot \tilde{L}_{\hat{f}}L_{\theta} \right),
 \end{aligned}$$

where the first equality follows from (C.13).

As a consequence, we have from (C.20) that

$$\mathbb{E}_{\bar{a}_i} \left[\left\| \frac{d\hat{a}_{i,n}}{d\theta} - \frac{d\bar{a}_i}{d\theta} \right\|_2 \right] \leq \tilde{L}_{\hat{f}}\tilde{L}_{\pi}^{i+1} \left(e + 4i \cdot \tilde{L}_{\hat{f}}\tilde{L}_{\pi} \cdot \hat{K}(i-1) + 2i \cdot \tilde{L}_{\hat{f}}L_{\theta} \right) + 2L_{\pi}\hat{K}(i) + 2L_{\theta}.$$

This concludes the proof. \square

C.3 Proof of Proposition 5.6

Proof. The analysis of gradient bias differs from that of gradient variance as it involves not only the distribution of approximate states but also the recurrent dependencies of the true value on future timesteps, which must be given extra attention.

In the following analysis, we will first apply similar techniques as those outlined in the previous section to establish an upper bound on the decomposed reward terms in the gradient bias. Afterward, we will address the distribution mismatch issue caused by the recursive structure of V^{π_θ} and the non-recursive structure of the value approximation \hat{V}_{ω_t} .

Step 1: Bound the cumulative reward terms in the gradient bias.

To begin with, we decompose the bias of the reward gradient at timestep $i \geq 0$ as follows,

$$\begin{aligned}
 & \mathbb{E}_{(s_i, a_i) \sim \mathbb{P}(s_i, a_i), (\hat{s}_{i,n}, \hat{a}_{i,n}) \sim \mathbb{P}(\hat{s}_i, \hat{a}_i)} \left[\left\| \frac{dr(\hat{x}_{i,n})}{d\theta} - \frac{dr(x_i)}{d\theta} \right\|_2 \right] \\
 &= \mathbb{E} \left[\left\| \frac{dr}{d\hat{x}_{i,n}} \cdot \frac{d\hat{x}_{i,n}}{d\theta} - \frac{dr}{dx_i} \cdot \frac{dx_i}{d\theta} \right\|_2 \right] \\
 &\leq \mathbb{E} \left[\left\| \frac{dr}{d\hat{x}_{i,n}} \cdot \frac{d\hat{x}_{i,n}}{d\theta} - \frac{dr}{dx_i} \cdot \frac{d\hat{x}_{i,n}}{d\theta} \right\|_2 + \left\| \frac{dr}{dx_i} \cdot \frac{d\hat{x}_{i,n}}{d\theta} - \frac{dr}{dx_i} \cdot \frac{dx_i}{d\theta} \right\|_2 \right] \\
 &\leq 2L_r \cdot \hat{K}(i) + L_r \cdot \left(\mathbb{E} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} - \frac{ds_i}{d\theta} \right\|_2 \right] + \mathbb{E} \left[\left\| \frac{d\hat{a}_{i,n}}{d\theta} - \frac{da_i}{d\theta} \right\|_2 \right] \right), \tag{C.21}
 \end{aligned}$$

where $\mathbb{P}(s_i, a_i)$ and $\mathbb{P}(\hat{s}_i, \hat{a}_i)$ are defined in (4.4) with respect to $s_0 \sim \nu_\pi, \hat{s}_0 \sim \nu_\pi$.

We have from (C.9) that for any $i \geq 1$,

$$\begin{aligned}
 & \mathbb{E} \left[\left\| \frac{d\hat{a}_{i,n}}{d\theta} - \frac{da_i}{d\theta} \right\|_2 \right] \\
 &= \mathbb{E} \left[\left\| \frac{\partial \hat{a}_{i,n}}{\partial \hat{s}_{i,n}} \cdot \frac{d\hat{s}_{i,n}}{d\theta} + \frac{\partial \hat{a}_{i,n}}{\partial \theta} - \frac{\partial a_i}{\partial s_i} \cdot \frac{ds_i}{d\theta} - \frac{\partial a_i}{\partial \theta} \right\|_2 \right]
 \end{aligned}$$

By the triangle inequality and the Lipschitz assumption, it then follows that

$$\begin{aligned}
 & \leq \mathbb{E} \left[\left\| \frac{\partial \hat{a}_{i,n}}{\partial \hat{s}_{i,n}} \cdot \frac{d\hat{s}_{i,n}}{d\theta} - \frac{\partial a_i}{\partial s_i} \cdot \frac{d\hat{s}_{i,n}}{d\theta} \right\|_2 \right] + \mathbb{E} \left[\left\| \frac{\partial a_i}{\partial s_i} \cdot \frac{d\hat{s}_{i,n}}{d\theta} - \frac{\partial a_i}{\partial s_i} \cdot \frac{ds_i}{d\theta} \right\|_2 \right] + \mathbb{E} \left[\left\| \frac{\partial \hat{a}_{i,n}}{\partial \theta} - \frac{\partial a_i}{\partial \theta} \right\|_2 \right] \\
 & \leq 2L_\pi \cdot \mathbb{E} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} \right\|_2 \right] + L_\pi \cdot \mathbb{E} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} - \frac{ds_i}{d\theta} \right\|_2 \right] + 2L_\theta. \tag{C.22}
 \end{aligned}$$

Similarly, we have from (C.10) that for any $i \geq 1$,

$$\begin{aligned}
 & \mathbb{E} \left[\left\| \frac{d\hat{s}_{i,n}}{d\theta} - \frac{ds_i}{d\theta} \right\|_2 \right] \\
 &= \mathbb{E} \left[\left\| \frac{\partial \hat{s}_{i,n}}{\partial \hat{s}_{i-1,n}} \cdot \frac{d\hat{s}_{i-1,n}}{d\theta} + \frac{\partial \hat{s}_{i,n}}{\partial \hat{a}_{i-1,n}} \cdot \frac{d\hat{a}_{i-1,n}}{d\theta} - \frac{\partial s_i}{\partial s_{i-1}} \cdot \frac{ds_{i-1}}{d\theta} - \frac{\partial s_i}{\partial a_{i-1}} \cdot \frac{da_{i-1}}{d\theta} \right\|_2 \right]
 \end{aligned}$$

Applying the triangle inequality to extract the $\epsilon_{f,t}$ term defined in (4.4), we proceed by

$$\begin{aligned}
 & \leq \mathbb{E} \left[\left\| \frac{\partial \hat{s}_{i,n}}{\partial \hat{s}_{i-1,n}} \cdot \frac{d\hat{s}_{i-1,n}}{d\theta} - \frac{\partial s_i}{\partial s_{i-1}} \cdot \frac{d\hat{s}_{i-1,n}}{d\theta} \right\|_2 \right] + \mathbb{E} \left[\left\| \frac{\partial s_i}{\partial s_{i-1}} \cdot \frac{d\hat{s}_{i-1,n}}{d\theta} - \frac{\partial s_i}{\partial s_{i-1}} \cdot \frac{ds_{i-1}}{d\theta} \right\|_2 \right] \\
 & \quad + \mathbb{E} \left[\left\| \frac{\partial \hat{s}_{i,n}}{\partial \hat{a}_{i-1,n}} \cdot \frac{d\hat{a}_{i-1,n}}{d\theta} - \frac{\partial s_i}{\partial a_{i-1}} \cdot \frac{d\hat{a}_{i-1,n}}{d\theta} \right\|_2 \right] + \mathbb{E} \left[\left\| \frac{\partial s_i}{\partial a_{i-1}} \cdot \frac{d\hat{a}_{i-1,n}}{d\theta} - \frac{\partial s_i}{\partial a_{i-1}} \cdot \frac{da_{i-1}}{d\theta} \right\|_2 \right] \\
 & \leq \epsilon_{f,t} \cdot \mathbb{E} \left[\left\| \frac{d\hat{s}_{i-1,n}}{d\theta} \right\|_2 + \left\| \frac{d\hat{a}_{i-1,n}}{d\theta} \right\|_2 \right] + L_f \cdot \mathbb{E} \left[\left\| \frac{d\hat{s}_{i-1,n}}{d\theta} - \frac{ds_{i-1}}{d\theta} \right\|_2 \right] \\
 & \quad + L_f \cdot \mathbb{E} \left[\left\| \frac{d\hat{a}_{i-1,n}}{d\theta} - \frac{da_{i-1}}{d\theta} \right\|_2 \right] \\
 & \leq \epsilon_{f,t} \cdot \hat{K}(i-1) + L_f \cdot \mathbb{E} \left[\left\| \frac{d\hat{s}_{i-1,n}}{d\theta} - \frac{ds_{i-1}}{d\theta} \right\|_2 \right] + L_f \cdot \mathbb{E} \left[\left\| \frac{d\hat{a}_{i-1,n}}{d\theta} - \frac{da_{i-1}}{d\theta} \right\|_2 \right], \tag{C.23}
 \end{aligned}$$

where the last inequality follows from the definition of $\widehat{K}(i-1)$ in (C.15).

Combining (C.22) and (C.23), we have

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{d\widehat{s}_{i,n}}{d\theta} - \frac{ds_i}{d\theta} \right\|_2 \right] &\lesssim (\epsilon_{f,t} + 2L_f L_\pi) \cdot \widehat{K}(i-1) + L_f \widetilde{L}_\pi \cdot \mathbb{E} \left[\left\| \frac{d\widehat{s}_{i-1,n}}{d\theta} - \frac{ds_{i-1}}{d\theta} \right\|_2 \right] + 2L_f L_\theta \\ &= ((\epsilon_{f,t} + 2L_f L_\pi) \cdot \widehat{K}(i-1) + 2L_f L_\theta) \cdot \sum_{j=0}^{i-1} L_f^j \widetilde{L}_\pi^j \\ &\leq ((\epsilon_{f,t} + 2L_f L_\pi) \cdot \widehat{K}(i-1) + 2L_f L_\theta) \cdot i \cdot \widetilde{L}_f^i \widetilde{L}_\pi^i, \end{aligned} \quad (\text{C.24})$$

where the last inequality follows from (C.13) and the fact that $s_0, \widehat{s}_{0,n}$ are sampled from the same initial distribution, and the equality holds by applying the recursion.

Plugging (C.24) into (C.22), we obtain

$$\mathbb{E} \left[\left\| \frac{d\widehat{a}_{i,n}}{d\theta} - \frac{da_i}{d\theta} \right\|_2 \right] \leq [(\epsilon_{f,t} + 2L_f L_\pi) \cdot \widehat{K}(i-1) + 2L_f L_\theta] \cdot i \cdot \widetilde{L}_f^i \widetilde{L}_\pi^{i+1} + 2L_\pi \widehat{K}(i) + 2L_\theta. \quad (\text{C.25})$$

Step 2: Address the state distribution mismatch issue.

The next step is to address the distribution mismatch issue caused by the recursive structure of the value function and the non-recursive structure of the value approximation, i.e., the critic.

We define $\bar{\sigma}_1(s, a) = \mathbb{P}(s_h = s, a_h = a)$ where $s_0 \sim \nu_\pi, a_i \sim \pi(\cdot | s_i)$, and $s_{i+1} \sim f(\cdot | s_i, a_i)$. In a similar way, we define $\widehat{\sigma}_1(s, a) = \mathbb{P}(\widehat{s}_h = s, \widehat{a}_h = a)$ where $\widehat{s}_0 \sim \nu_\pi, \widehat{a}_i \sim \pi(\cdot | \widehat{s}_i)$, and $\widehat{s}_{i+1} \sim \widehat{f}(\cdot | \widehat{s}_i, \widehat{a}_i)$.

Now we are ready to bound the gradient bias. From Lemma C.4, we know that

$$\begin{aligned} b_t &\leq \kappa \kappa' \cdot \mathbb{E}_{s_0 \sim \nu_\pi, \widehat{s}_{0,n} \sim \nu_\pi} \left[\left\| \nabla_\theta \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) - \nabla_\theta \sum_{i=0}^{h-1} \gamma^i \cdot r(\widehat{s}_{i,n}, \widehat{a}_{i,n}) \right\|_2 \right] \\ &\quad + \kappa' \gamma^h \cdot \mathbb{E}_{(s_h, a_h) \sim \bar{\sigma}_1, (\widehat{s}_{h,n}, \widehat{a}_{h,n}) \sim \widehat{\sigma}_1} \left[\left\| \frac{\partial Q^{\pi_\theta}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_\theta}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \nabla_\theta \widehat{Q}_t(\widehat{s}_{h,n}, \widehat{a}_{h,n}) \right\|_2 \right], \end{aligned} \quad (\text{C.26})$$

where recall that $\kappa' = \beta + \kappa \cdot (1 - \beta)$.

From Lemma C.5, we have for any policy π_θ that the state-action value function is L_Q -Lipschitz continuous, which gives for any $s \in \mathcal{S}$ and $a \in \mathcal{A}$ that

$$\left\| \frac{\partial Q^{\pi_\theta}}{\partial a} \right\|_2 \leq L_Q, \quad \left\| \frac{\partial Q^{\pi_\theta}}{\partial s} \right\|_2 \leq L_Q. \quad (\text{C.27})$$

The bias brought by the critic, i.e., the last term on the right-hand side of (C.26), can be further bounded by

$$\begin{aligned} &\mathbb{E}_{(s_h, a_h) \sim \bar{\sigma}_1, (\widehat{s}_{h,n}, \widehat{a}_{h,n}) \sim \widehat{\sigma}_1} \left[\left\| \frac{\partial Q^{\pi_\theta}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_\theta}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \nabla_\theta \widehat{Q}_t(\widehat{s}_{h,n}, \widehat{a}_{h,n}) \right\|_2 \right] \\ &= \mathbb{E}_{\bar{\sigma}_1, \widehat{\sigma}_1} \left[\left\| \frac{\partial Q^{\pi_\theta}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_\theta}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \frac{\partial \widehat{Q}_t}{\partial \widehat{a}_{h,n}} \cdot \frac{d\widehat{a}_{h,n}}{d\theta} - \frac{\partial \widehat{Q}_t}{\partial \widehat{s}_{h,n}} \cdot \frac{d\widehat{s}_{h,n}}{d\theta} \right\|_2 \right] \\ &\leq \mathbb{E}_{\bar{\sigma}_1, \widehat{\sigma}_1} \left[\left\| \frac{\partial Q^{\pi_\theta}}{\partial a_h} \cdot \frac{da_h}{d\theta} - \frac{\partial Q^{\pi_\theta}}{\partial a_h} \cdot \frac{d\widehat{a}_{h,n}}{d\theta} \right\|_2 + \left\| \frac{\partial Q^{\pi_\theta}}{\partial a_h} \cdot \frac{d\widehat{a}_{h,n}}{d\theta} - \frac{\partial \widehat{Q}_t}{\partial \widehat{a}_{h,n}} \cdot \frac{d\widehat{a}_{h,n}}{d\theta} \right\|_2 \right. \\ &\quad \left. + \left\| \frac{\partial Q^{\pi_\theta}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \frac{\partial Q^{\pi_\theta}}{\partial s_h} \cdot \frac{d\widehat{s}_{h,n}}{d\theta} \right\|_2 + \left\| \frac{\partial Q^{\pi_\theta}}{\partial s_h} \cdot \frac{d\widehat{s}_{h,n}}{d\theta} - \frac{\partial \widehat{Q}_t}{\partial \widehat{s}_{h,n}} \cdot \frac{d\widehat{s}_{h,n}}{d\theta} \right\|_2 \right] \\ &\leq L_Q \cdot \left(\mathbb{E}_{\bar{\sigma}_1, \widehat{\sigma}_1} \left[\left\| \frac{da_h}{d\theta} - \frac{d\widehat{a}_{h,n}}{d\theta} \right\|_2 + \left\| \frac{ds_h}{d\theta} - \frac{d\widehat{s}_{h,n}}{d\theta} \right\|_2 \right] \right) + \left(\frac{\gamma^h}{1-\gamma} \right)^2 \widehat{K}(h) \cdot \epsilon_{v,t}, \end{aligned} \quad (\text{C.28})$$

where the equality follows from the chain rule and the fact that the critic \widehat{Q}_t has a non-recursive structure, the last inequality follows from (C.27), (C.15) and the definition of $\epsilon_{v,t}$ in (4.5).

Plugging (C.21) and (C.28) into (C.26), we obtain

$$b_t \leq \kappa \kappa' \cdot h \cdot \left(L_r \cdot \left(\mathbb{E}_{\bar{\sigma}_1, \hat{\sigma}_1} \left[\left\| \frac{d\hat{s}_{h,n}}{d\theta} - \frac{ds_h}{d\theta} \right\|_2 \right] + \mathbb{E}_{\bar{\sigma}_1, \hat{\sigma}_1} \left[\left\| \frac{d\hat{a}_{h,n}}{d\theta} - \frac{da_h}{d\theta} \right\|_2 \right] \right) + 2L_r \cdot \hat{K}(h) \right) \\ + \kappa' \gamma^h \cdot \left(L_Q \cdot \left(\mathbb{E}_{\bar{\sigma}_1, \hat{\sigma}_1} \left[\left\| \frac{d\hat{s}_{h,n}}{d\theta} - \frac{ds_h}{d\theta} \right\|_2 \right] + \mathbb{E}_{\bar{\sigma}_1, \hat{\sigma}_1} \left[\left\| \frac{d\hat{a}_{h,n}}{d\theta} - \frac{da_h}{d\theta} \right\|_2 \right] \right) + \hat{K}(h) \cdot \left(\frac{\gamma^h}{1-\gamma} \right)^2 \epsilon_{v,t} \right), \quad (\text{C.29})$$

Plugging (C.24), (C.25), and (C.15) into the (C.29), we conclude the proof by obtaining

$$b_t = O \left(\kappa \kappa' h^2 \frac{1-\gamma^h}{1-\gamma} \tilde{L}_f^h \tilde{L}_f^h \tilde{L}_\pi^{2h} \epsilon_{f,t} + \kappa' h \gamma^h \left(\frac{\gamma^h}{1-\gamma} \right)^2 \tilde{L}_f^h \tilde{L}_\pi^h \epsilon_{v,t} \right). \quad (\text{C.30})$$

□

Lemma C.3. The expected value gradient over the state distribution $\mathbb{P}(s_h)$ can be represented by

$$\mathbb{E}_{s_h \sim \mathbb{P}(s_h)} [\nabla_\theta V^{\pi_\theta}(s_h)] = \mathbb{E}_{(s,a) \sim \bar{\sigma}_1} \left[\frac{\partial Q^{\pi_\theta}}{\partial a} \cdot \frac{da}{d\theta} + \frac{\partial Q^{\pi_\theta}}{\partial s} \cdot \frac{ds}{d\theta} \right],$$

where $\mathbb{P}(s_h)$ is the state distribution at timestep h when $s_0 \sim \zeta$, $a_i \sim \pi(\cdot | s_i)$, and $s_{i+1} \sim f(\cdot | s_i, a_i)$.

Proof. At state s_h , the value gradient can be rewritten as

$$\begin{aligned} \nabla_\theta V^{\pi_\theta}(s_h) &= \nabla_\theta \mathbb{E} \left[r(s_h, a_h) + \gamma \cdot \int_{\mathcal{S}} f(s_{h+1} | s_h, a_h) \cdot V^\pi(s_{h+1}) ds_{h+1} \right] \\ &= \nabla_\theta \mathbb{E} [r(s_h, a_h)] + \gamma \cdot \mathbb{E} \left[\nabla_\theta \int_{\mathcal{S}} f(s_{h+1} | s_h, a_h) \cdot V^\pi(s_{h+1}) ds_{h+1} \right] \\ &= \mathbb{E} \left[\frac{\partial r_h}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial r_h}{\partial s_h} \cdot \frac{ds_h}{d\theta} \right. \\ &\quad \left. + \gamma \int_{\mathcal{S}} \left(\nabla_\theta f(s_{h+1} | s_h, a_h) \cdot V^\pi(s_{h+1}) + f(s_{h+1} | s_h, a_h) \cdot \nabla_\theta V^\pi(s_{h+1}) \right) ds_{h+1} \right] \\ &= \mathbb{E} \left[\frac{\partial r_h}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial r_h}{\partial s_h} \cdot \frac{ds_h}{d\theta} + \gamma \int_{\mathcal{S}} \left(\nabla_a f(s_{h+1} | s_h, a) \cdot \frac{da_h}{d\theta} \cdot V^\pi(s_{h+1}) \right. \right. \\ &\quad \left. \left. + \nabla_s f(s_{h+1} | s_h, a_h) \cdot \frac{ds_h}{d\theta} \cdot V^\pi(s_{h+1}) + f(s_{h+1} | s_h, a_h) \cdot \nabla_\theta V^\pi(s_{h+1}) \right) ds_{h+1} \right], \quad (\text{C.31}) \end{aligned}$$

where the first equation follows from the Bellman equation and the last two equations hold due to the chain rule. Here, it is worth noting that when $h \geq 1$, both a_h and s_h have dependencies on all previous timesteps. For any $h \geq 1$, we have from the chain rule that $\nabla_\theta r(s_h, a_h) = \partial r_h / \partial a_h \cdot da_h / d\theta + \partial r_h / \partial s_h \cdot ds_h / d\theta$. This differs from the case when $h = 0$, e.g., in the deterministic policy gradient theorem (Silver et al., 2014), where we can simply write $\nabla_\theta r(s_h, a_h) = \partial r_h / \partial a_h \cdot \partial a_h / \partial \theta$.

Rearranging terms in (C.31) gives

$$\begin{aligned} \nabla_\theta V^{\pi_\theta}(s_h) &= \mathbb{E} \left[\nabla_a \left(r(s_h, a_h) + \gamma \int_{\mathcal{S}} f(s_{h+1} | s_h, a_h) \cdot V^\pi(s_{h+1}) ds_{h+1} \right) \cdot \frac{da_h}{d\theta} \right. \\ &\quad \left. + \nabla_s \left(r(s_h, a_h) + \gamma \int_{\mathcal{S}} f(s_{h+1} | s_h, a_h) \cdot V^\pi(s_{h+1}) ds_{h+1} \right) \cdot \frac{ds_h}{d\theta} \right. \\ &\quad \left. + \gamma \int_{\mathcal{S}} f(s_{h+1} | s_h, a_h) \cdot \nabla_\theta V^\pi(s_{h+1}) ds_{h+1} \right] \\ &= \mathbb{E} \left[\frac{\partial Q^{\pi_\theta}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_\theta}}{\partial s_h} \cdot \frac{ds_h}{d\theta} + \gamma \int_{\mathcal{S}} f(s_{h+1} | s_h, a_h) \cdot \nabla_\theta V^\pi(s_{h+1}) ds_{h+1} \right], \quad (\text{C.32}) \end{aligned}$$

where the last equation holds since $Q^{\pi_\theta}(s_h, a_h) = r(s_h, a_h) + \gamma \int_{\mathcal{S}} f(s_{h+1} | s_h, a_h) \cdot V^\pi(s_{h+1}) ds_{h+1}$.

By recursively applying (C.32), we obtain

$$\nabla_\theta V^{\pi_\theta}(s_h) = \mathbb{E} \left[\int_{\mathcal{S}} \sum_{i=h}^{\infty} \gamma^{i-h} \cdot f(s_{i+1} | s_i, a_i) \cdot \left(\frac{\partial Q^{\pi_\theta}}{\partial a_i} \cdot \frac{da_i}{d\theta} + \frac{\partial Q^{\pi_\theta}}{\partial s_i} \cdot \frac{ds_i}{d\theta} \right) ds_{i+1} \right]. \quad (\text{C.33})$$

Let $\bar{\sigma}_2(s, a) = (1 - \gamma) \cdot \sum_{i=h}^{\infty} \gamma^{i-h} \cdot \mathbb{P}(s_i = s, a_i = a)$, where $s_0 \sim \zeta$, $a_i \sim \pi(\cdot | s_i)$, and $s_{i+1} \sim f(\cdot | s_i, a_i)$. By definition we have

$$(1 - \gamma) \cdot \sum_{i=0}^{h-1} \gamma^i \cdot \mathbb{P}(s_i = s, a_i = a) + \gamma^h \cdot \bar{\sigma}_1(s, a) = \sigma(s, a) = (1 - \gamma) \cdot \sum_{i=0}^{h-1} \gamma^i \cdot \mathbb{P}(s_i = s, a_i = a) + \gamma^h \cdot \bar{\sigma}_2(s, a).$$

Therefore we have the equivalence $\bar{\sigma}_1(s, a) = \bar{\sigma}_2(s, a)$.

By taking the expectation over s_h in (C.33), we have the stated result, i.e.,

$$\mathbb{E}_{s_h \sim \mathbb{P}(s_h)} [\nabla_{\theta} V^{\pi_{\theta}}(s_h)] = \mathbb{E}_{(s,a) \sim \bar{\sigma}_2} \left[\frac{\partial Q^{\pi_{\theta}}}{\partial a} \cdot \frac{da}{d\theta} + \frac{\partial Q^{\pi_{\theta}}}{\partial s} \cdot \frac{ds}{d\theta} \right] = \mathbb{E}_{(s,a) \sim \bar{\sigma}_1} \left[\frac{\partial Q^{\pi_{\theta}}}{\partial a} \cdot \frac{da}{d\theta} + \frac{\partial Q^{\pi_{\theta}}}{\partial s} \cdot \frac{ds}{d\theta} \right].$$

□

Lemma C.4. Recall that the state distribution μ_{π} where $\hat{s}_{0,n}$ is sampled from is of the form $\mu_{\pi}(s) = \beta \cdot \nu_{\pi}(s) + (1 - \beta) \cdot \zeta(s)$. The gradient bias b_t at any iteration t satisfies

$$\begin{aligned} b_t &\leq \kappa(\beta + \kappa \cdot (1 - \beta)) \cdot \mathbb{E}_{s_0 \sim \nu_{\pi}, \hat{s}_{0,n} \sim \nu_{\pi}} \left[\left\| \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) \right\|_2 \right] \\ &\quad + (\beta + \kappa \cdot (1 - \beta)) \gamma^h \cdot \mathbb{E}_{(s_h, a_h) \sim \bar{\sigma}_1, (\hat{s}_{h,n}, \hat{a}_{h,n}) \sim \bar{\sigma}_1} \left[\left\| \frac{\partial Q^{\pi_{\theta}}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_{\theta}}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \nabla_{\theta} \hat{Q}_t(\hat{s}_{h,n}, \hat{a}_{h,n}) \right\|_2 \right]. \end{aligned}$$

Proof. To begin, we decompose the gradient bias by

$$\begin{aligned} b_t &= \left\| \nabla_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\hat{\nabla}_{\theta} J(\pi_{\theta_t})] \right\|_2 \\ &= \left\| \mathbb{E}[\nabla_{\theta} J(\pi_{\theta_t}) - \hat{\nabla}_{\theta} J(\pi_{\theta_t})] \right\|_2 \\ &= \left\| \mathbb{E}_{s_0 \sim \zeta, \hat{s}_{0,n} \sim \mu_{\pi}} \left[\nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) + \gamma^h \cdot \nabla_{\theta} V^{\pi_{\theta}}(s_h) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) - \gamma^h \cdot \nabla_{\theta} \hat{V}_t(\hat{s}_{h,n}) \right] \right\|_2, \end{aligned} \tag{C.34}$$

where we note that s_0 and $\hat{s}_{0,n}$ are sampled from ζ and μ_{π} following the definition of the RL objective and the form of gradient estimator, respectively.

For $\mu_{\pi}(s) = \beta \cdot \nu_{\pi}(s) + (1 - \beta) \cdot \zeta(s)$, let Z be the random variable satisfying $\mathbb{P}(Z = 0) = \beta$ and $\mathbb{P}(Z = 1) = 1 - \beta$, i.e., the event $Z = 0$ and $Z = 1$ corresponds to that the state s is sampled from ν_{π} and ζ , respectively. For any random variable Y , following the law of total expectation, we know that

$$\begin{aligned} \mathbb{E}_{\mu_{\pi}}[Y] &= \mathbb{E}[\mathbb{E}[Y|Z]] = \mathbb{E}[Y|Z = 0]\mathbb{P}(Z = 0) + \mathbb{E}[Y|Z = 1]\mathbb{P}(Z = 1) \\ &= \beta \mathbb{E}[Y|Z = 0] + (1 - \beta) \mathbb{E}[Y|Z = 1] \\ &= \beta \mathbb{E}_{\nu_{\pi}}[Y] + (1 - \beta) \mathbb{E}_{\zeta}[Y]. \end{aligned} \tag{C.35}$$

Therefore, we have from (C.34) that

$$\begin{aligned} b_t &\leq \mathbb{E}_{\hat{s}_{0,n} \sim \mu_{\pi}} \left[\left\| \mathbb{E}_{s_0 \sim \zeta} \left[\nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) + \gamma^h \cdot \nabla_{\theta} V^{\pi_{\theta}}(s_h) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) - \gamma^h \cdot \nabla_{\theta} \hat{V}_t(\hat{s}_{h,n}) \right] \right\|_2 \right] \\ &\leq \beta \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} \left[\left\| \mathbb{E}_{s_0 \sim \zeta} \left[\nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) + \gamma^h \cdot \nabla_{\theta} V^{\pi_{\theta}}(s_h) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) - \gamma^h \cdot \nabla_{\theta} \hat{V}_t(\hat{s}_{h,n}) \right] \right\|_2 \right] \\ &\quad + (1 - \beta) \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \zeta} \left[\left\| \mathbb{E}_{s_0 \sim \zeta} \left[\nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) + \gamma^h \cdot \nabla_{\theta} V^{\pi_{\theta}}(s_h) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) - \gamma^h \cdot \nabla_{\theta} \hat{V}_t(\hat{s}_{h,n}) \right] \right\|_2 \right], \end{aligned} \tag{C.36}$$

where the first inequality holds since $\|\mathbb{E}[\cdot]\|_2 \leq \mathbb{E}[\|\cdot\|_2]$ and the second inequality holds due to (C.35).

Using the result from Lemma C.3, we know that

$$\begin{aligned} & \mathbb{E}_{s_0 \sim \zeta} \left[\nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) + \gamma^h \cdot \nabla_{\theta} V^{\pi_{\theta}}(s_h) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) - \gamma^h \cdot \nabla_{\theta} \hat{V}_t(\hat{s}_{h,n}) \right] \\ &= \underbrace{\mathbb{E}_{s_0 \sim \zeta} \left[\nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) \right]}_{B_r} + \underbrace{\gamma^h \mathbb{E}_{(s_h, a_h) \sim \bar{\sigma}_1} \left[\frac{\partial Q^{\pi_{\theta}}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_{\theta}}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \nabla_{\theta} \hat{V}_t(\hat{s}_{h,n}) \right]}_{B_v}. \end{aligned}$$

Here, the shorthand notation B_r denotes the bias introduced by the h -step model expansion and B_v denotes the bias introduced by using a critic for tail estimation. Then we may rewrite (C.36) as

$$\begin{aligned} b_t &\leq \beta \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} [\|B_r + B_v\|_2] + (1 - \beta) \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \zeta} [\|B_r + B_v\|_2] \\ &\leq \left(\beta \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} [\|B_r\|_2] + (1 - \beta) \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \zeta} [\|B_r\|_2] \right) + \left(\beta \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} [\|B_v\|_2] + (1 - \beta) \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \zeta} [\|B_v\|_2] \right). \end{aligned} \quad (\text{C.37})$$

For the first term on the right-hand side of (C.37), we have

$$\begin{aligned} & \beta \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} [\|B_r\|_2] + (1 - \beta) \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \zeta} [\|B_r\|_2] \\ &= \beta \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} \left[\left\| \mathbb{E}_{s_0 \sim \zeta} \left[\nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) \right] \right\|_2 \right] \\ &\quad + (1 - \beta) \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} \left[\left\| \mathbb{E}_{s_0 \sim \zeta} \left[\nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) \right] \right\|_2 \right] \cdot \left\{ \mathbb{E}_{\nu_{\pi}} \left[\left(\frac{d\zeta}{d\nu_{\pi}}(s) \right)^2 \right] \right\}^{1/2} \\ &\leq (\beta + \kappa \cdot (1 - \beta)) \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} \left[\left\| \mathbb{E}_{s_0 \sim \zeta} \left[\nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) \right] \right\|_2 \right] \\ &\leq \kappa (\beta + \kappa \cdot (1 - \beta)) \cdot \mathbb{E}_{s_0 \sim \nu_{\pi}, \hat{s}_{0,n} \sim \nu_{\pi}} \left[\left\| \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(s_i, a_i) - \nabla_{\theta} \sum_{i=0}^{h-1} \gamma^i \cdot r(\hat{s}_{i,n}, \hat{a}_{i,n}) \right\|_2 \right], \end{aligned} \quad (\text{C.38})$$

where the first and second inequalities follow from the definition of κ in Proposition 5.6.

Similarly, for the second term on the right-hand side of (C.37), we have

$$\begin{aligned} & \beta \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} [\|B_v\|_2] + (1 - \beta) \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \zeta} [\|B_v\|_2] \\ &= \beta \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \nu_{\pi}} \left[\gamma^h \cdot \left\| \mathbb{E}_{(s_h, a_h) \sim \bar{\sigma}_1} \left[\frac{\partial Q^{\pi_{\theta}}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_{\theta}}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \nabla_{\theta} \hat{V}_t(\hat{s}_{h,n}) \right] \right\|_2 \right] \\ &\quad + (1 - \beta) \cdot \mathbb{E}_{\hat{s}_{0,n} \sim \zeta} \left[\gamma^h \cdot \left\| \mathbb{E}_{(s_h, a_h) \sim \bar{\sigma}_1} \left[\frac{\partial Q^{\pi_{\theta}}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_{\theta}}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \nabla_{\theta} \hat{V}_t(\hat{s}_{h,n}) \right] \right\|_2 \right] \\ &\leq (\beta + \kappa \cdot (1 - \beta)) \gamma^h \cdot \mathbb{E}_{(s_h, a_h) \sim \bar{\sigma}_1, \hat{s}_{0,n} \sim \nu_{\pi}} \left[\left\| \frac{\partial Q^{\pi_{\theta}}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_{\theta}}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \nabla_{\theta} \hat{V}_t(\hat{s}_{h,n}) \right\|_2 \right] \\ &= (\beta + \kappa \cdot (1 - \beta)) \gamma^h \cdot \mathbb{E}_{(s_h, a_h) \sim \bar{\sigma}_1, (\hat{s}_{h,n}, \hat{a}_{h,n}) \sim \bar{\sigma}_1} \left[\left\| \frac{\partial Q^{\pi_{\theta}}}{\partial a_h} \cdot \frac{da_h}{d\theta} + \frac{\partial Q^{\pi_{\theta}}}{\partial s_h} \cdot \frac{ds_h}{d\theta} - \nabla_{\theta} \hat{Q}_t(\hat{s}_{h,n}, \hat{a}_{h,n}) \right\|_2 \right]. \end{aligned} \quad (\text{C.39})$$

Plugging (C.38) and (C.39) into (C.37) completes the proof. \square

Lemma C.5 (Lipschitz Value Function (Rachelson & Lagoudakis, 2010) Theorem 1). Under Assumption 5.3, for $\gamma L_f(1 + L_{\pi}) < 1$, then the state-action value function is L_Q -Lipschitz continuous, such that for any policy π_{θ} , state $s_1, s_2 \in \mathcal{S}$ and action $a_1, a_2 \in \mathcal{A}$, $|Q^{\pi_{\theta}}(s_1, a_1) - Q^{\pi_{\theta}}(s_2, a_2)| \leq L_Q \cdot \|(s_1 - s_2, a_1 - a_2)\|_2$, and

$$L_Q = L_r / (1 - \gamma L_f(1 + L_{\pi})).$$

C.4 Proof of Proposition 5.7

Proof. When $\gamma \approx 1$, we have

$$\frac{1 - \gamma^h}{1 - \gamma} = \sum_{i=0}^{h-1} \gamma^i \approx h, \quad \frac{\gamma^h}{1 - \gamma} = \frac{1}{1 - \gamma} - \frac{1 - \gamma^h}{1 - \gamma} \approx \frac{1}{1 - \gamma} - h.$$

We denote by $H = 1/(1 - \gamma) = \sum_{i=0}^{\infty} \gamma^i$ the effective task horizon.

To find the optimal unroll length h^* that minimizes the upper bound of the convergence, we define $g(h)$ as follows,

$$g(h) = c \cdot (2\delta \cdot b'_t + \frac{\eta}{2} \cdot v_t'^2) + b_t'^2 + v_t'^2.$$

Here, $v_t'^2$ and b_t' are the leading terms in the variance, bias bound (i.e., (C.18) and (C.30)) when L_f , $L_{\hat{f}}$, and L_{π} are less than or equal to 1. Formally, $v_t' = h^3$ and $b_t' = h^3 \epsilon_{f,t} + h(H-h)^2 \epsilon_{v,t}$. We consider the terms that are only dependent on h , H , $\epsilon_{f,t}$, and $\epsilon_{v,t}$ to simplify the analysis and determine the order of h^* .

Our first problem is to find the optimal model unroll h'^* that minimizes $g(h)$. We notice that $g(h)$ increases monotonically with respect to b_t' and v_t' when they are non-negative. This further simplifies the problem to find

$$h'^* = \underset{h}{\operatorname{argmin}} b_t' + c' v_t' = \underset{h}{\operatorname{argmin}} \underbrace{h^3(\epsilon_{f,t} + c') + h(H-h)^2 \epsilon_{v,t}}_{g_1(h)}, \quad (\text{C.40})$$

where c' is some constant that does not affect the order of h'^* .

By taking the derivative of the right-hand side of (C.40) with respect to h and setting it to zero, we obtain

$$\frac{\partial}{\partial h} g_1(h) = 3h^2 \cdot (\epsilon_{f,t} + c') + (3h^2 - 4Hh + H^2) \cdot \epsilon_{v,t} = 0. \quad (\text{C.41})$$

Solve the above quadratic equation with respect to h , we have the two non-negative roots $h_1'^*$ and $h_2'^*$ as follows,

$$h_1'^* = \frac{4H\epsilon_{v,t} + \sqrt{(4H\epsilon_{v,t})^2 - 12c_1\epsilon_{v,t}H^2}}{6c_1}, \quad h_2'^* = \frac{4H\epsilon_{v,t} - \sqrt{(4H\epsilon_{v,t})^2 - 12c_1\epsilon_{v,t}H^2}}{6c_1},$$

where we define $c_1 = \epsilon_{f,t} + \epsilon_{v,t} + c'$.

Now we study the resulting two cases. If $(4H\epsilon_{v,t})^2 - 12c_1\epsilon_{v,t}H^2 \geq 0$, we have

$$h'^* = h_1'^* = O(\epsilon_{v,t}/(\epsilon_{f,t} + \epsilon_{v,t}) \cdot H).$$

We can verify that $h_1'^*$ is indeed the minimum by calculating the second-order derivative at $h_1'^*$ as follows,

$$\begin{aligned} \frac{\partial^2 g_1(h_1'^*)}{\partial h^2} &= \frac{4H\epsilon_{v,t} + \sqrt{(4H\epsilon_{v,t})^2 - 4c_1 \cdot H^2}}{6c_1} * 6(\epsilon_{f,t} + \epsilon_{v,t} + c') - 4H\epsilon_{v,t} \\ &= \sqrt{(4H\epsilon_{v,t})^2 - 4c_1 \cdot H^2} > 0. \end{aligned}$$

The other case is $(4H\epsilon_{v,t})^2 - 12c_1\epsilon_{v,t}H^2 < 0$. When this happens, (C.41) does not have a real solution h'^* and we set h^* to 0. This concludes the proof of Proposition 5.7. \square

C.5 Proof of Corollary 5.9

Proof. We let the learning rate $\eta = 1/\sqrt{T}$. Then for $T \geq 4L^2$, we have $c = (\eta - L\eta^2)^{-1} \leq 2\sqrt{T}$ and $L\eta \leq 1/2$. By setting $N = O(\sqrt{T})$, we obtain

$$\begin{aligned} \min_{t \in [T]} \mathbb{E} \left[\|\nabla_{\theta} J(\pi_{\theta_t})\|_2^2 \right] &\leq \frac{4}{T} \cdot \left(\sum_{t=0}^{T-1} c \cdot (2\delta \cdot b_t + \frac{\eta}{2} \cdot v_t) + b_t^2 + v_t \right) + \frac{4c}{T} \cdot \mathbb{E}[J(\pi_{\theta_T}) - J(\pi_{\theta_1})] \\ &\leq \frac{4}{T} \left(\sum_{t=0}^{T-1} 4\sqrt{T}\delta \cdot b_t + b_t^2 + 2v_t \right) + \frac{8}{\sqrt{T}} \cdot \mathbb{E}[J(\pi_{\theta_T}) - J(\pi_{\theta_1})] \\ &\leq \frac{4}{T} \left(\sum_{t=0}^{T-1} 4\sqrt{T}\delta \cdot b_t + b_t^2 \right) + O(1/\sqrt{T}) \\ &\leq \frac{16\delta}{\sqrt{T}} \varepsilon(T) + \frac{4}{T} \varepsilon^2(T) + O(1/\sqrt{T}). \end{aligned}$$

This concludes the proof. \square

D Experimental Details

D.1 Implementations and Comparisons with More RL Baselines

For the model-based baseline Model-Based Policy Optimization (MBPO) (Janner et al., 2019), we use the implementation in the Mbrl-lib (Pineda et al., 2021). For all other model-free baselines, we use the implementations in Tianshou (Weng et al., 2021) that have state-of-the-art results.

We observe that the RP-DP has competitive performance in all the evaluation tasks compared to the popular baselines, suggesting the importance of studying model-based RP PGMs. In experiments, we implement RP-DR as the on-policy SVG(1) (Heess et al., 2015). We observe that the training can be unstable when using the off-policy SVG implementation, which requires a carefully chosen policy update rate as well as a proper size of the experience replay buffer. This is because when the learning rate is large, the magnitude of the inferred policy noise (from the previous data samples in the experience replay) can be huge. Implementing an on-policy version of RP-DR can avoid such an issue, following (Heess et al., 2015). This, however, can degrade the performance of RP-DR compared to the off-policy RP-DP algorithm in several tasks. We conjecture that implementing the off-policy version of RP-DR can boost its performance, which requires techniques to stabilize training and we leave it as future work. For RP-DP, we implement it as Model-Augmented Actor-Critic (MAAC) (Clavera et al., 2020) with entropy regularization (Haarnoja et al., 2018), as suggested by (Amos et al., 2021). RP(0) represents setting $h = 0$ in the RP PGM formulas (Amos et al., 2021), which is a model-free algorithm that is a stochastic counterpart of deterministic policy gradients.

For model-free baselines, we compare with Likelihood Ratio (LR) policy gradient methods (c.f. (2.2)), including REINFORCE (Sutton et al., 1999), Natural Policy Gradient (NPG) (Kakade, 2001), Advantage Actor Critic (A2C), Actor Critic using Kronecker-Factored Trust Region (ACKTR) (Wu et al., 2017), and Proximal Policy Optimization (PPO) (Schulman et al., 2017). We also evaluate algorithms that are built upon DDPG (Lillicrap et al., 2015), including Soft Actor-Critic (SAC) (Haarnoja et al., 2018) and Twin Delayed Deep Deterministic policy gradient (TD3) (Fujimoto et al., 2018).

D.2 Implementation and Ablation of Spectral Normalization

In experiments, we use Multilayer Perceptrons (MLPs) for the critic, policy, and model. Besides, we adopt Gaussian dynamical models and policies as the source of stochasticity. To test the benefit of smooth function approximations in model-based RP policy gradient algorithms, spectral normalization is applied to all layers of the policy MLP and all except the final layers of the model MLP. The number of layers for the policy and the dynamics model is 4 and 5, respectively.

Our code is based on PyTorch (Paszke et al., 2019), which has an out-of-the-shelf implementation of spectral normalization. Thus, applying SN to the MLP is pretty simple and no additional lines of code are needed. Specifically, we only need to import and apply SN to each layer:

```
from torch.nn.utils.parametrizations import spectral_norm
layer = [spectral_norm(nn.Linear(in_dim, hidden_dim)), nn.ReLU()]
```

Moreover, we conduct ablation study on the functions that spectral normalization is applied to: Both the model and the policy (default setting); Only the model; Only the policy; No SN is applied (vanilla setting). The results are shown in Figure 8.

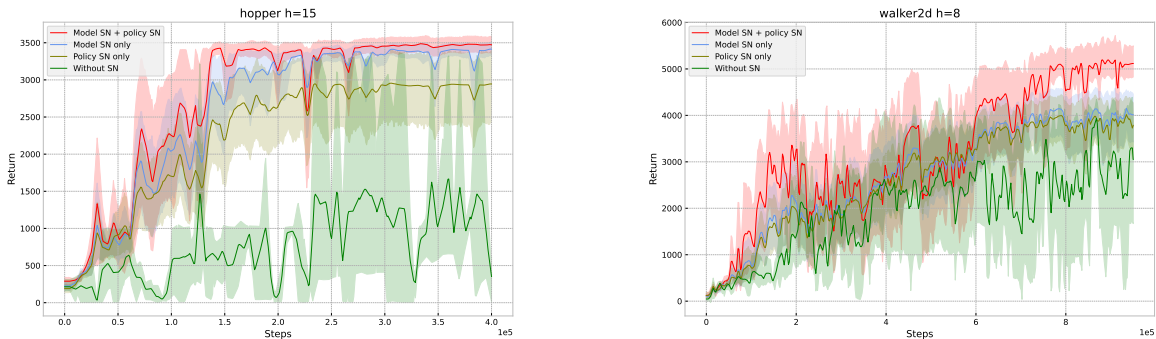


Figure 8. Ablation on different model learners: single-step and multi-step state prediction models, and multi-step state prediction models trained with an additional directional derivative error.

We observe in Fig. 8 that for both the hopper and the walker2d tasks, applying SN to the model and policy simultaneously achieves the best performance, which supports our theoretical results. Besides, learning a smooth transition kernel by applying SN to the neural network model only is slightly better than only applying SN to the policy. At the same time, the vanilla implementation of model-based RP PGM fails to give acceptable result.

D.3 Ablation on Different Model Learners

Our main theoretical results in Section 5 depend on the model error defined in (4.4), which, however, cannot directly serve as the model training objective. For this reason, we evaluate different model learners: single- and multi-step (h -step) state prediction models, as well as multi-step predictive models integrated with the directional derivative error (Li et al., 2021). The results are reported in Figure 9. We observe that enlarging the prediction steps benefits training. The algorithm also converges faster in walker2d when considering derivative error, which approximately minimizes 4.4 and supports our analysis. However, calculating the directional derivative error by searching k nearest points in the buffer significantly increases the computational cost, for which reason we use h -step state predictive models as default in experiments.

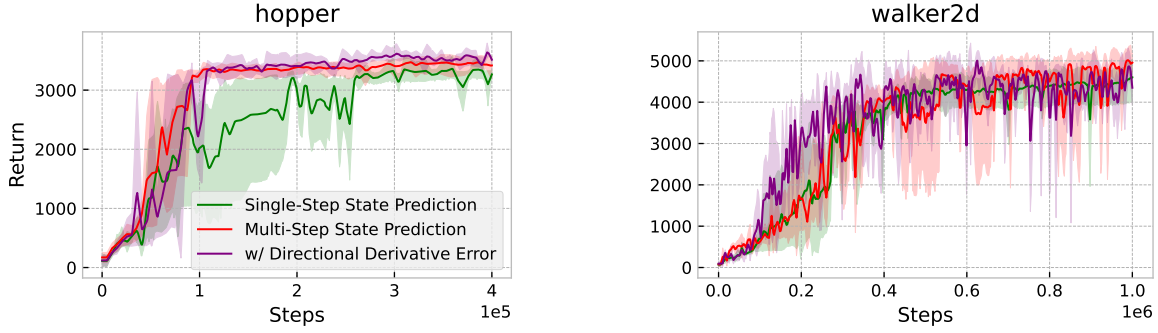


Figure 9. Ablation on different model learners: single-step and multi-step state prediction models, and multi-step state prediction models trained with an additional directional derivative error.

D.4 Figures in the Main Text in Larger Sizes

Here, we provide identical figures that are larger in size. Figure 10, 11, 12, 13 correspond to Figure 1, 4, 6, 7 in the main text, respectively.

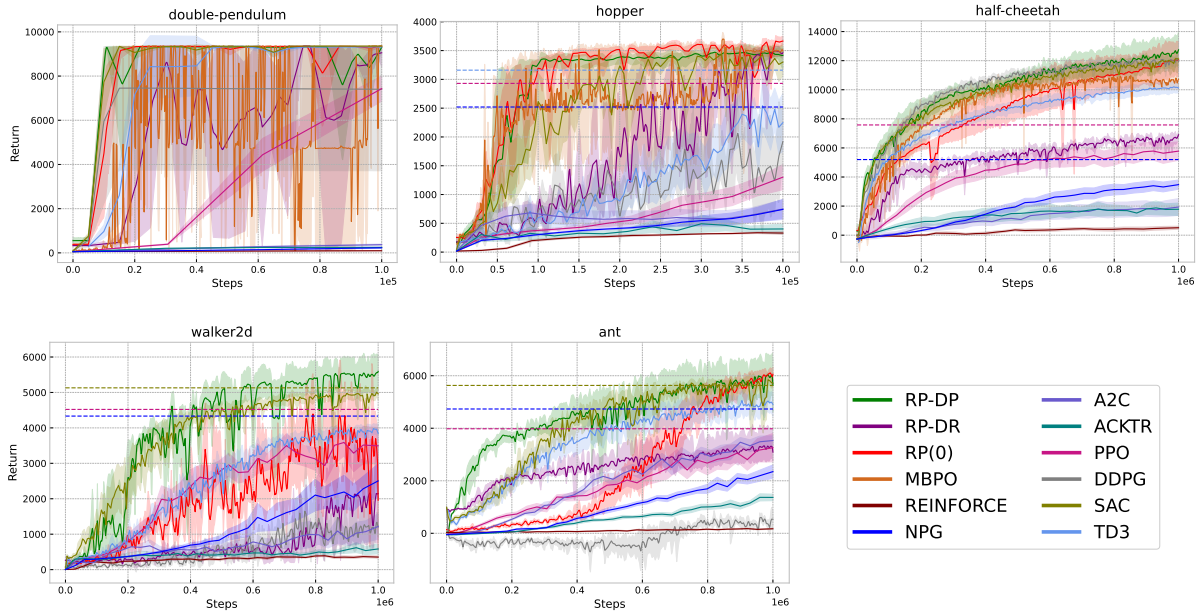


Figure 10. Evaluation of model-based RP PGMs in MuJoCo tasks. The dashed lines represent the value at the convergence of the corresponding model-free algorithms.

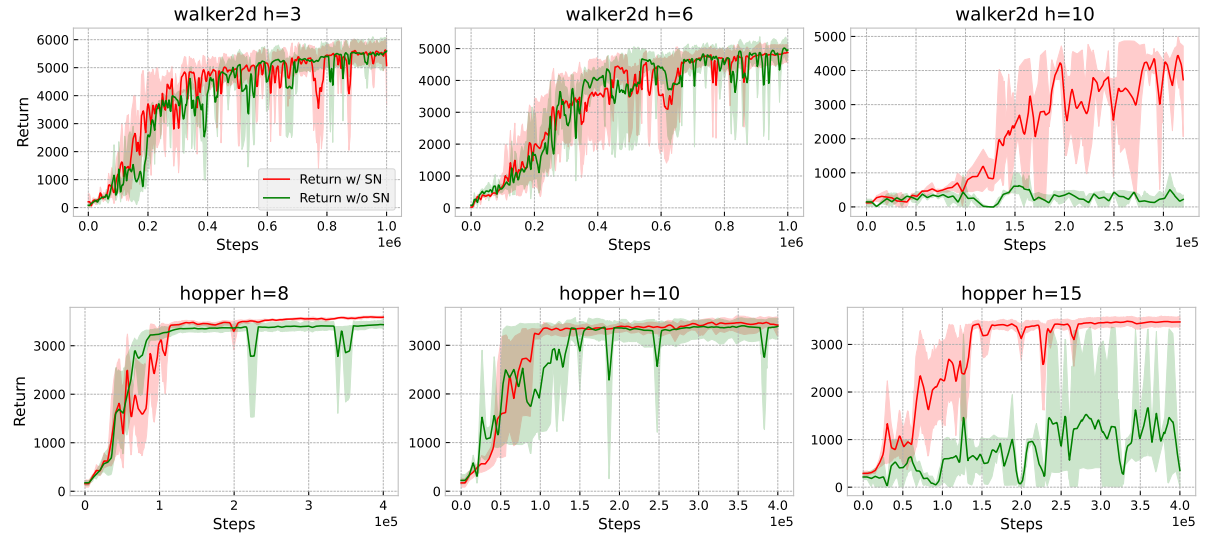


Figure 11. Performance of model-based RP PG methods with and without spectral normalization.

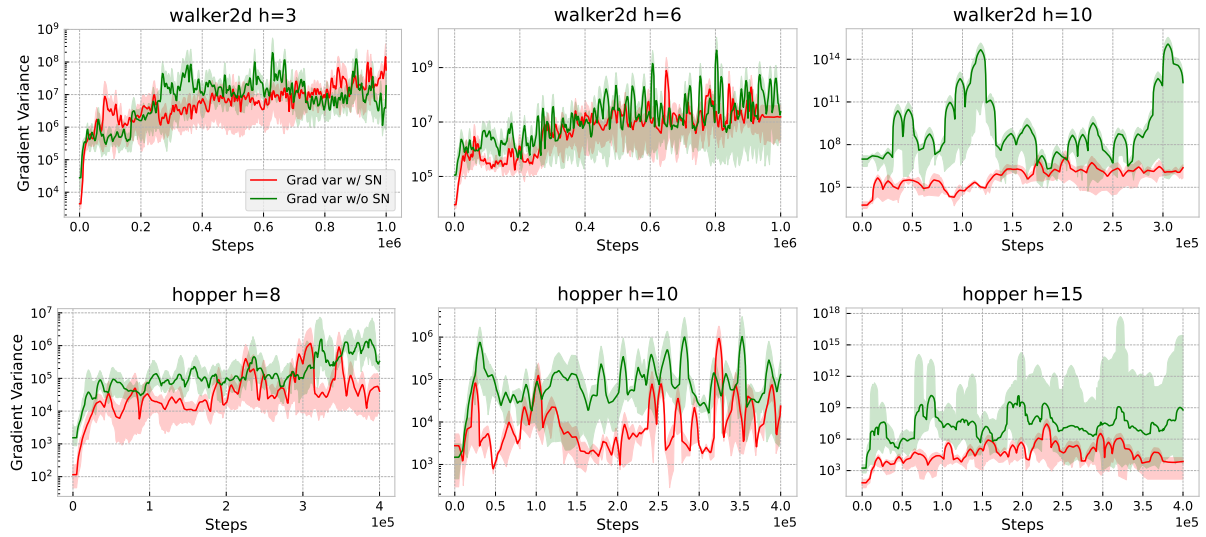


Figure 12. Ablation on the gradient variance.

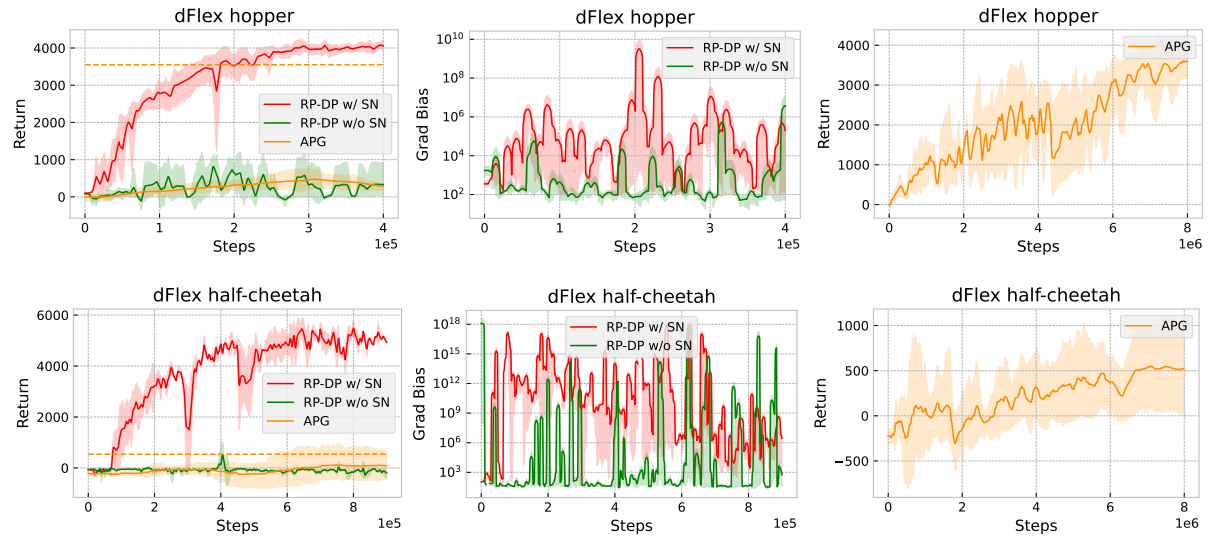


Figure 13. Performance and gradient bias in differentiable simulation. The last column is the full training curves of APG, which needs 20 times more steps in hopper to reach a comparable return with RP-DP-SN in the first column.