
NEORL: Efficient Exploration for Nonepisodic RL

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the problem of nonepisodic reinforcement learning (RL) for nonlinear
2 dynamical systems, where the system dynamics are unknown and the RL agent
3 has to learn from a single trajectory, i.e., without resets. We propose *Nonepisodic*
4 *Optimistic RL* (NEORL), an approach based on the principle of optimism in
5 the face of uncertainty. NEORL uses well-calibrated probabilistic models and
6 plans optimistically w.r.t. the epistemic uncertainty about the unknown dynamics.
7 Under continuity and bounded energy assumptions on the system, we provide a
8 first-of-its-kind regret bound of $\mathcal{O}(\beta_T \sqrt{TT_T})$ for general nonlinear systems with
9 Gaussian process dynamics. We compare NEORL to other baselines on several
10 deep RL environments and empirically demonstrate that NEORL achieves the
11 optimal average cost while incurring the least regret.

12 1 Introduction

13 In recent years, data-driven control approaches, such as reinforcement learning (RL), have demon-
14 strated remarkable achievements. However, most RL algorithms are devised for an episodic setting,
15 where during each episode, the agent interacts in the environment for a predetermined episode
16 length or until a termination condition is met. After the episode, the agent is reset back to an initial
17 state from where the next episode commences. Episodes prevent the system from blowing up, i.e.,
18 maintain stability, while also restricting exploration to states that are relevant to the task at hand.
19 Moreover, resets ensure that the agent explores close to the initial states and does not end up at
20 undesirable parts of the state space that exhibit low reward. In simulation, resetting is typically
21 straightforward. However, if we wish to enable agents to learn by interacting with the real world,
22 resets are often prohibitive since they typically involve manual intervention. Instead, agents should
23 be able to learn autonomously (Sharma et al., 2021) i.e., from a single trajectory. While several works
24 in the Deep RL community have addressed this challenge, (c.f., Section 5), the theoretical results for
25 this setting are fairly limited. In particular, the setting has been extensively studied for finite state and
26 action spaces (Kearns & Singh, 2002; Brafman & Tennenholtz, 2002; Jaksch et al., 2010) and linear
27 systems (Abbasi-Yadkori & Szepesvári, 2011; Simchowitz & Foster, 2020; Dean et al., 2020; Lale
28 et al., 2020). However, the extension to nonlinear systems is much less understood. In our work, we
29 address this gap and propose a practical RL algorithm that is grounded in theory. In particular, we
30 make the following contributions.

31 Contributions

- 32 1. We propose, NEORL, a novel model-based RL algorithm based on the principle of optimism in
33 the face of uncertainty. NEORL operates in a nonepisodic setting and picks average cost optimal
34 policies optimistically w.r.t. to the model’s epistemic uncertainty.
- 35 2. We show that when the dynamics lie in a reproducing kernel Hilbert space (RKHS) of kernel
36 k , NEORL exhibits a regret of $\mathcal{O}(\beta_T \sqrt{TT_T})$, where the regret, akin to prior work, is measured
37 w.r.t to the optimal average cost under known dynamics, T is the number of environment steps

38 and Γ_T the maximum information gain of kernel k (Srinivas et al., 2012). Our regret bound is
 39 similar to the ones obtained in the episodic setting (Kakade et al., 2020; Curi et al., 2020; Sukhija
 40 et al., 2024; Treven et al., 2024) and Gaussian process (GP) bandit optimization (Srinivas et al.,
 41 2012; Chowdhury & Gopalan, 2017; Scarlett et al., 2017). To the best of our knowledge, we are
 42 the first to obtain regret bounds for the setting.

43 3. We evaluate NEORL on several RL benchmarks against common model-based RL baselines.
 44 Our experimental results demonstrate that NEORL consistently achieves sublinear regret, also
 45 when neural networks are employed instead of GPs for modeling dynamics. Moreover, in all
 46 our experiments, NEORL converges to the optimal average cost.

47 2 Problem Setting

48 We consider a discrete-time dynamical system with running costs c .

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathbf{f}^*(\mathbf{x}_t, \mathbf{u}_t) + \mathbf{w}_t, (\mathbf{x}_t, \mathbf{u}_t) \in \mathcal{X} \times \mathcal{U}, \mathbf{x}(0) = \mathbf{x}_0 & (1) \\ c(\mathbf{x}, \mathbf{u}) &\in \mathbb{R}_{\geq 0} & \text{(Running cost)} \end{aligned}$$

49 Here $\mathbf{x}_t \in \mathcal{X} \subseteq \mathbb{R}^{d_x}$ is the state, $\mathbf{u}_t \in \mathcal{U} \subseteq \mathbb{R}^{d_u}$ the control input, and $\mathbf{w}_t \in \mathcal{W} \subseteq \mathbb{R}^w$ the process
 50 noise. The dynamics \mathbf{f}^* are unknown and the cost c is assumed to be known.

51 **Task** In this work, we study the average cost RL problem (Puterman, 2014), i.e., we want to learn
 52 the solution to the following minimization problem

$$A(\boldsymbol{\pi}^*, \mathbf{x}_0) = \min_{\boldsymbol{\pi} \in \Pi} A(\boldsymbol{\pi}, \mathbf{x}_0) = \min_{\boldsymbol{\pi} \in \Pi} \limsup_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\pi}} \left[\sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t) \right]. \quad (2)$$

53 Moreover, we consider the nonepisodic RL setting where the system starts at an initial state $\mathbf{x}_0 \in \mathcal{X}$
 54 but never resets back during learning, that is, we seek to learn from a single trajectory. After each step
 55 t in the environment, the RL system receives a transition tuple $(\mathbf{x}_t, \mathbf{u}_t, \mathbf{x}_{t+1})$ and updates its policy
 56 based on the data \mathcal{D}_t collected thus far during learning. The average cost formulation is common
 57 for the nonepisodic setting (Jaksch et al., 2010; Abbasi-Yadkori & Szepesvári, 2011; Simchowitz &
 58 Foster, 2020), and the cumulative regret for the learning algorithm in this case is defined as

$$R_T = \sum_{t=0}^{T-1} \mathbb{E}_{\mathbf{x}_t, \mathbf{u}_t | \mathbf{x}_0} [c(\mathbf{x}_t, \mathbf{u}_t) - A(\boldsymbol{\pi}^*, \mathbf{x}_0)]. \quad (3)$$

59 Studying the average cost criterion for general continuous state-action spaces is challenging even
 60 when the dynamics are known, since the average cost exists only for special classes of nonlinear
 61 systems (Arapostathis et al., 1993). In the following, we impose assumptions on the dynamics and
 62 policy class Π that enable our theoretical analysis.

63 2.1 Assumptions

64 Imposing continuity on \mathbf{f}^* is quite common in the control theory (Khalil, 2015) and reinforcement
 65 learning literature (Curi et al., 2020; Sussex et al., 2023; Sukhija et al., 2024). To this end, for our
 66 analysis, we make the following assumption.

67 **Assumption 2.1** (Continuity of \mathbf{f}^* and $\boldsymbol{\pi}$). The dynamics model \mathbf{f}^* and all $\boldsymbol{\pi} \in \Pi$ are continuous.

68 Next, we make an assumption on the system’s stochasticity.

69 **Assumption 2.2** (Process noise distribution). The process noise is i.i.d. Gaussian with variance
 70 σ^2 , i.e., $\mathbf{w}_t \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$.

71 For simplicity, we focus on the homoscedastic setting. However, the analysis can be extended for the
 72 more general heteroscedastic case. In the following, we make assumptions on our policy class. To
 73 this end, we first introduce the class of \mathcal{K}_∞ functions.

74 **Definition 2.3** (\mathcal{K}_∞ -functions). The function $\xi : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ is of class \mathcal{K}_∞ , if it is continuous,
 75 strictly increasing, $\xi(0) = 0$ and $\xi(s) \rightarrow \infty$ for $s \rightarrow \infty$.

76 **Assumption 2.4** (Policies with bounded energy). We assume there exists $\kappa, \xi \in \mathcal{K}_\infty$, positive
 77 constants K, C_u, C_l with $C_u > C_l$, and $\gamma \in (0, 1)$ such that for each $\boldsymbol{\pi} \in \Pi$ we have,

78 *Bounded energy:* There exists a Lyapunov function $V^\pi : \mathcal{X} \rightarrow [0, \infty)$ for which

$$\begin{aligned} |V^\pi(\mathbf{x}) - V^\pi(\mathbf{x}')| &\leq \kappa(\|\mathbf{x} - \mathbf{x}'\|) && \text{(uniform continuity)} \\ C_l \xi(\|\mathbf{x}\|) &\leq V^\pi(\mathbf{x}) \leq C_u \xi(\|\mathbf{x}\|) && \text{(positive definiteness)} \\ \mathbb{E}_{\mathbf{x}'|\mathbf{x}, \pi}[V^\pi(\mathbf{x}')] &\leq \gamma V^\pi(\mathbf{x}) + K && \text{(drift condition)} \end{aligned}$$

79 *Bounded norm of cost:*

$$\sup_{\mathbf{x} \in \mathcal{X}} \frac{c(\mathbf{x}, \pi(\mathbf{x}))}{1 + V^\pi(\mathbf{x})} < \infty$$

80 *Boundedness of the noise with respect to κ :*

$$\mathbb{E}_{\mathbf{w}}[\kappa(\|\mathbf{w}\|)] < \infty, \quad \mathbb{E}_{\mathbf{w}}[\kappa^2(\|\mathbf{w}\|)] < \infty$$

81 The bounded energy assumption is introduced to ensure that the system does not end up in states from
 82 which it cannot recover. In particular, the Lyapunov function V^π can be viewed as an energy function
 83 for the dynamical system, and the drift condition above ensures that in expectation the energy at
 84 the next state \mathbf{x}' is not increasing to ∞ , that is, the system is not “blowing up”. Other works that
 85 study learning nonlinear dynamics (Foster et al., 2020; Sattar & Oymak, 2022; Lale et al., 2021)
 86 in the nonepisodic setting also make stability assumptions such as global exponential stability for
 87 their analysis. In similar spirit, we make the bounded energy assumption for our policy class. The
 88 drift condition on the Lyapunov function is also used to study the ergodicity of Markov chains for
 89 continuous state spaces (Meyn & Tweedie, 2012; Hairer & Mattingly, 2011), which is crucial for our
 90 analysis of the infinite horizon behavior of the system. Moreover, for a very rich class of problems,
 91 the drift condition is satisfied. We highlight this in the corollary below.

92 **Corollary 2.5.** *Assume f^* is uniformly continuous and for all $\pi \in \Pi$, $\mathbf{x} \in \mathcal{X}$, $\|\pi(\mathbf{x})\| \leq u_{\max}$.
 93 Further assume, there exists $\pi_s \in \Pi$ such that we have constants K, C_u, C_l with $C_u > C_l$, $\gamma \in (0, 1)$,
 94 $\kappa, \alpha \in \mathcal{K}_\infty$ and a Lyapunov function $V : \mathcal{X} \rightarrow [0, \infty)$ for which*

$$\begin{aligned} |V(\mathbf{x}) - V(\mathbf{x}')| &\leq \kappa(\|\mathbf{x} - \mathbf{x}'\|) \\ C_l \xi(\|\mathbf{x}\|) &\leq V(\mathbf{x}) \leq C_u \xi(\|\mathbf{x}\|) \\ \mathbb{E}_{\mathbf{x}'|\mathbf{x}, \pi_s}[V(\mathbf{x}')] &\leq \gamma V(\mathbf{x}) + K. \end{aligned}$$

95 *Then, V also satisfies the drift condition for all $\pi \in \Pi$, i.e., is a Lyapunov function for all policies.*

96 We prove this corollary in Appendix A. Intuitively, if the inputs are bounded, the energy inserted into
 97 the system by another policy is also bounded. Nearly all real-world systems have bounded inputs due
 98 to the physical limitations of actuators. For these systems, it suffices if only one policy in Π satisfies
 99 the drift condition. In Appendix A.4, we discuss an alternative set of assumptions on the costs, that
 100 relaxes the bounded energy requirement on the policy class Π .

101 The boundedness assumptions for the cost and the noise in Assumption 2.4 are satisfied for a rich
 102 class of cost and \mathcal{K}_∞ functions.

103 Under these assumptions, we can show the existence of the average cost solution.

104 **Theorem 2.6** (Existence of Average Cost Solution). *Let Assumption 2.1 – 2.4 hold. Consider any
 105 $\pi \in \Pi$ and let P^π denote its transition kernel, i.e., $P^\pi(\mathbf{x}, \mathcal{A}) = \mathbb{P}(\mathbf{x}' \in \mathcal{A} | \mathbf{x}, \pi(\mathbf{x}))$. Then P^π
 106 admits a unique invariant measure \bar{P}^π and there exists $C_2, C_3 \in (0, \infty)$, $\lambda \in (0, 1)$ such that*

107 *Average Cost;*

$$A(\pi) = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t) \right] = \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [c(\mathbf{x}, \pi(\mathbf{x}))]$$

108 *Bias Cost; Letting $B(\pi, \mathbf{x}_0) = \lim_{T \rightarrow \infty} \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t) - A(\pi) \right]$ denote the bias, we have*

$$|B(\pi, \mathbf{x}_0)| = \left| \lim_{T \rightarrow \infty} \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t) - A(\pi) \right] \right| \leq C_2(1 + V^\pi(\mathbf{x}_0)) \frac{1}{1 - \lambda}$$

109 *for all $\mathbf{x}_0 \in \mathcal{X}$.*

110 Theorem 2.6 is a crucial result for our analysis since it implies that the average cost is bounded and
 111 *independent of the initial state* \mathbf{x}_0 . Furthermore, it also shows that the bias is bounded. Similar to the
 112 discounted case, the average cost criterion satisfies the following Bellman equation (Puterman, 2014)

$$B(\boldsymbol{\pi}, \mathbf{x}) + A(\boldsymbol{\pi}) = c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \mathbb{E}_{\mathbf{x}'}[B(\boldsymbol{\pi}, \mathbf{x}') | \mathbf{x}, \boldsymbol{\pi}] \quad (4)$$

113 Accordingly, the bias term plays an important role in the regret analysis (also notice its similarity to
 114 our regret term in Equation (3)).

115 Thus far, we have only made assumptions that make the average cost problem tractable. In the
 116 following, we make an assumption on the dynamics that allow us to learn it from data. We start with
 117 the definition of a well-calibrated statistical model of \mathbf{f}^* .

118 **Definition 2.7** (Well-calibrated statistical model of \mathbf{f}^* , Rothfuss et al. (2023)). Let $\mathcal{Z} \stackrel{\text{def}}{=} \mathcal{X} \times \mathcal{U}$. An
 119 all-time well-calibrated statistical model of the function \mathbf{f}^* is a sequence $\{\mathcal{M}_n(\delta)\}_{n \geq 0}$, where

$$\mathcal{M}_n(\delta) \stackrel{\text{def}}{=} \{ \mathbf{f} : \mathcal{Z} \rightarrow \mathbb{R}^{d_x} \mid \forall \mathbf{z} \in \mathcal{Z}, \forall j \in 1, \dots, d_x : |\mu_{n,j}(\mathbf{z}) - f_j(\mathbf{z})| \leq \beta_n(\delta) \sigma_{n,j}(\mathbf{z}) \},$$

120 if, with probability at least $1 - \delta$, we have $\mathbf{f}^* \in \bigcap_{n \geq 0} \mathcal{M}_n(\delta)$. Here, $\mu_{n,j}$ and $\sigma_{n,j}$ denote the
 121 j -th element in the vector-valued mean and standard deviation functions $\boldsymbol{\mu}_n$ and $\boldsymbol{\sigma}_n$ respectively,
 122 and $\beta_n(\delta) \in \mathbb{R}_{\geq 0}$ is a scalar function that depends on the confidence level $\delta \in (0, 1]$ and which is
 123 monotonically increasing in n .

124 Next, we assume that \mathbf{f}^* resides in a Reproducing Kernel Hilbert Space (RKHS) of vector-valued
 125 functions and show that this is sufficient for us to obtain a well-calibrated model.

126 **Assumption 2.8.** We assume that the functions f_j^* , $j \in 1, \dots, d_x$ lie in a RKHS with kernel k
 127 and have a bounded norm B , that is $\mathbf{f}^* \in \mathcal{H}_{k,B}^{d_x}$, with $\mathcal{H}_{k,B}^{d_x} = \{ \mathbf{f} \mid \|f_j\|_k \leq B, j = 1, \dots, d_x \}$.
 128 Moreover, we assume that $k(\mathbf{x}, \mathbf{x}) \leq \sigma_{\max}$ for all $\mathbf{x} \in \mathcal{X}$.

129 The mean and epistemic uncertainty of the vector-valued function \mathbf{f}^* are denoted with $\boldsymbol{\mu}_n(\mathbf{z}) =$
 130 $[\mu_{n,j}(\mathbf{z})]_{j \leq d_x}$, and $\boldsymbol{\sigma}_n(\mathbf{z}) = [\sigma_{n,j}(\mathbf{z})]_{j \leq d_x}$ and have an analytical solution

$$\begin{aligned} \mu_{n,j}(\mathbf{z}) &= \bar{\mu}^j(\mathbf{z}) + \mathbf{k}_n^\top(\mathbf{z})(\mathbf{K}_n + \sigma^2 \mathbf{I})^{-1}(\mathbf{y}_{1:n}^j - \bar{\mu}_{1:n}^j), \\ \sigma_{n,j}^2(\mathbf{z}) &= k(\mathbf{x}, \mathbf{x}) - \mathbf{k}_n^\top(\mathbf{z})(\mathbf{K}_n + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_n(\mathbf{x}), \end{aligned} \quad (5)$$

131 Here, $\mathbf{y}_{1:n}^j$ corresponds to the noisy measurements of f_j^* , i.e., the observed next state from the
 132 transitions dataset $\mathcal{D}_{1:n}$, $\bar{\mu}^j(\mathbf{z})$ corresponds to the fixed mean function, e.g., $\bar{\mu}^j(\mathbf{z}) = \mathbf{x}$, $\bar{\mu}_{1:n}^j$ its
 133 values on the dataset, $\mathbf{k}_n = [k(\mathbf{z}, \mathbf{z}_i)]_{i \leq nT}$, $\mathbf{z}_i \in \mathcal{D}_{1:n}$, and $\mathbf{K}_n = [k(\mathbf{z}_i, \mathbf{z}_l)]_{i,l \leq nT}$, $\mathbf{z}_i, \mathbf{z}_l \in \mathcal{D}_{1:n}$
 134 is the data kernel matrix. The restriction on the kernel $k(\mathbf{x}, \mathbf{x}) \leq \sigma_{\max}$ has also appeared in works
 135 studying the episodic setting for nonlinear systems (Mania et al., 2020; Kakade et al., 2020; Curi
 136 et al., 2020; Sukhija et al., 2024; Wagenmaker et al., 2023).

137 **Lemma 2.9** (Well calibrated confidence intervals for RKHS, Rothfuss et al. (2023)). Let $\mathbf{f}^* \in \mathcal{H}_{k,B}^{d_x}$.
 138 Suppose $\boldsymbol{\mu}_n$ and $\boldsymbol{\sigma}_n$ are the posterior mean and variance of a GP with kernel k , c.f., Equation (5).
 139 There exists $\beta_n(\delta)$, for which the tuple $(\boldsymbol{\mu}_n, \boldsymbol{\sigma}_n, \beta_n(\delta))$ is a well-calibrated statistical model of \mathbf{f}^* .

140 In summary, in the RKHS setting, a GP is a well-calibrated model. For more general models like
 141 BNNs, methods such as Kuleshov et al. (2018) can be used for calibration. Our results can also be
 142 extended beyond the RKHS setting to other classes of well-calibrated models similar to Curi et al.
 143 (2020).

144 3 NEORL

145 In the following, we present our algorithm: **Nonepisodic Optimistic RL** (NEORL) for efficient
 146 nonepisodic exploration in continuous state-action spaces. NEORL builds on recent advances in
 147 episodic RL (Kakade et al., 2020; Curi et al., 2020; Sukhija et al., 2024; Treven et al., 2024) and
 148 leverages the optimism in the face of uncertainty paradigm to pick policies that are optimistic w.r.t. the
 149 dynamics within our calibrated statistical model. Moreover, NEORL suggests policies according to
 150 the following decision rule

$$(\boldsymbol{\pi}_n, \mathbf{f}_n) \stackrel{\text{def}}{=} \arg \min_{\boldsymbol{\pi} \in \Pi, \mathbf{f} \in \mathcal{M}_{n-1} \cap \mathcal{M}_0} A(\boldsymbol{\pi}, \mathbf{f}). \quad (6)$$

Algorithm 1 NEORL: NONEPISODIC OPTIMISTIC RL

Init: Aleatoric uncertainty σ , Probability δ , Statistical model $(\mu_0, \sigma_0, \beta_0(\delta))$, H_0
for $n = 1, \dots, N$ **do**

$\pi_n = \arg \min_{\pi \in \Pi} \min_{f \in \mathcal{M}_{n-1} \cap \mathcal{M}_0} A(\pi, f)$	► Prepare policy
$H_n = 2H_{n-1}$	► Set horizon
$\mathcal{D}_n \leftarrow \text{ROLLOUT}(\pi_n)$	► Collect measurements for horizon H_n
Update $(\mu_n, \sigma_n, \beta_n) \leftarrow \mathcal{D}_n$	► Update model

end for

151 Here, f_n is a dynamical system such that the cost by controlling f_n with its optimal policy π_n is
152 the lowest among all the plausible systems from $\mathcal{M}_{n-1} \cap \mathcal{M}_0$. Note, from Lemma 2.9 we have
153 that $f^* \in \mathcal{M}_{n-1} \cap \mathcal{M}_0$ (with high probability) and therefore the solution to Equation (6) gives an
154 optimistic estimate for the average cost.

155 NEORL proceeds in the following manner. Similar to Jaksch et al. (2010), we bin the total time T
156 the agent spends interacting in the environment into N “artificial” episodes. At each episode, we
157 pick a policy according to Equation (6) and roll it out for H_n steps on the system. Next, we use
158 the data collected during the rollout to update our statistical model. Finally, we double the horizon
159 $H_{n+1} = 2H_n$, akin to Simchowitz & Foster (2020), and continue to the next episode *without resetting*
160 the system back to the initial state x_0 . The algorithm is summarized in Algorithm 1.

161 3.1 Theoretical Results

162 In the following, we study the theoretical properties for NEORL and provide a first-of-its-kind bound
163 on the cumulative regret for the average cost criterion for general nonlinear dynamical systems. Our
164 bound depends on the *maximum information gain* of kernel k (Srinivas et al., 2012), defined as

$$\Gamma_T(k) = \max_{\mathcal{A} \subset \mathcal{X} \times \mathcal{U}; |\mathcal{A}| \leq T} \frac{1}{2} \log |\mathbf{I} + \sigma^{-2} \mathbf{K}_T|.$$

165 Γ_T represents the complexity of learning f^* and is sublinear for a very rich class of kernels (Vakili
166 et al., 2021). In Appendix A, we report the dependence of Γ_T on T in Table 1.

167 **Theorem 3.1** (Cumulative Regret of NEORL). *Let Assumption 2.1 – 2.8 hold, and define H_0 as the*
168 *smallest integer such that*

$$H_0 > \frac{\log(C_u/C_l)}{\log(1/\gamma)}.$$

169 *Then with probability at least $1 - \delta$, we have the following regret for NEORL*

$$R_T \leq D_4(x_0, K, \gamma) \beta_T \sqrt{T \Gamma_T} + D_5(x_0, K, \gamma) \log_2 \left(\frac{T}{H_0} + 1 \right). \quad (7)$$

170 *with $D_4(x_0, K, \gamma), D_5(x_0, K, \gamma) \in (0, \infty)$ when $\|x_0\| < \infty, K < \infty$, and $\gamma < 1$.*

171 Theorem 3.1 gives sublinear regret for a rich class of RKHS functions. Moreover, it also gives a
172 minimal horizon H_0 that we need to maintain before switching to the next policy. Even for the
173 linear case, fast switching between stable controllers can destabilize the closed-loop system. We
174 ensure this does not happen in our case by having a minimal horizon of H_0 . Lastly, the regret
175 bound depends on constants D_4 and D_5 . The constants are finite when $\gamma < 1, K < \infty$ (bounded
176 energy from Assumption 2.4 is satisfied), and $\|x_0\| < \infty$. Theorem 3.1 can also be derived beyond
177 the RKHS setting for a more general class of well-calibrated models. In this case, the maximum
178 information gain is replaced by the model complexity from Curi et al. (2020) (c.f., Curi et al. (2020);
179 Sukhija et al. (2024); Treven et al. (2024) for further detail).

180 3.2 Practical Modifications

181 For testing NEORL, we make three modifications that simplify its deployment in practice in terms
182 of implementation and computation time. First, instead of doubling the horizon H_n we pick a fixed
183 horizon H during the experiment. This makes the planning and training of the agent easier. Next,

Algorithm 2 Practical NEORL:

Init: Aleatoric uncertainty σ , Probability δ , Statistical model $(\mu_0, \sigma_0, \beta_0(\delta))$
for $n = 1, \dots, N$ **do**
 for $h = 1, \dots, H$ **do**

$$\min_{\mathbf{u}_0: H_{\text{MPC}}-1, \boldsymbol{\eta}_0: H_{\text{MPC}}-1} \mathbb{E} \left[\sum_{h=0}^{H_{\text{MPC}}-1} c(\hat{\mathbf{x}}_h, \mathbf{u}_h) \right]; \mathbf{x}_0 = \mathbf{x}_h^n \quad \blacktriangleright \text{Solve MPC problem}$$

 $(\mathbf{x}_n^h, \mathbf{u}_0^*, \mathbf{x}_n^{h+1}) \leftarrow \text{ROLLOUT}(\mathbf{u}_0^*) \quad \blacktriangleright \text{Collect transition}$
 end for
 Update $(\mu_n, \sigma_n, \beta_n) \leftarrow \mathcal{D}_n$
end for

184 we use a receding horizon controller, i.e., model predictive control (MPC) (García et al., 1989),
185 instead of directly optimizing for the average cost in Equation (6). MPC is widely used to obtain a
186 feedback controller for the infinite horizon setting. Moreover, while for linear systems, the Riccati
187 equations (Anderson & Moore, 2007) provide an analytical solution to Equation (2), no such solution
188 exists for the nonlinear case and MPC is commonly used as an approximation. Further, under
189 additional assumptions on the cost and dynamics, MPC also obtains a policy with bounded average
190 cost, which is crucial for the nonepisodic case (c.f., Assumption 2.4). We use the iCEM optimizer for
191 planning (Pinneri et al., 2021). Finally, instead of optimizing over $\mathcal{M}_n \cap \mathcal{M}_0$, we optimize directly
192 over \mathcal{M}_n . This allows us to use the reparameterization trick from Curi et al. (2020) and obtain a
193 simple and tractable optimization problem. In summary, for each step t in the environment, we solve
194 the following optimization problem

$$\min_{\mathbf{u}_0: H_{\text{MPC}}-1, \boldsymbol{\eta}_0: H_{\text{MPC}}-1} \mathbb{E} \left[\sum_{h=0}^{H_{\text{MPC}}-1} c(\hat{\mathbf{x}}_h, \mathbf{u}_h) \right], \quad (8)$$

s.t. $\hat{\mathbf{x}}_{h+1} = \mu_{n-1}(\hat{\mathbf{x}}_h, \mathbf{u}_h) + \beta_{n-1}(\delta)\sigma_{n-1}(\hat{\mathbf{x}}_h, \mathbf{u}_h)\boldsymbol{\eta}_h + \mathbf{w}_h$ and $\hat{\mathbf{x}}_0 = \mathbf{x}_t$.

195 Here H_{MPC} is the MPC horizon. We take the first input from the solution of the problem above,
196 i.e., \mathbf{u}_0^* , and execute this in the system. We then repeat this procedure for H steps and then update
197 our statistical model \mathcal{M}_n . The resulting optimization above considers a larger action space as it
198 includes the hallucinated controls $\boldsymbol{\eta}$ (Curi et al., 2020) as additional input variables. Moreover, the
199 final algorithm can be seen as a natural extension to H-UCRL (Curi et al., 2020) for the nonepisodic
200 setting. We summarize the algorithm in Algorithm 2. Note while these modifications deviate from
201 our theoretical analysis, empirically they work well for GP and BNN models, c.f., Section 4.

202 4 Experiments

203 We evaluate NEORL on the Pendulum-v1 and MountainCar environment from the OpenAI gym
204 benchmark suite (Brockman et al., 2016), Cartpole, Reacher, and Swimmer from the DeepMind
205 control suite (Tassa et al., 2018), the racecar simulator from Kabzan et al. (2020), and a soft robotic
206 arm from Tekinalp et al. (2024). The swimmer and the soft robotic arm are fairly high-dimensional
207 systems – the swimmer has a 28-dimensional state and 5-dimensional action space, and the soft arm
208 is represented by a 58-dimensional state and has a 12-dimensional action space. All environments
209 are never reset during learning. Moreover, the Pendulum-v1, MountainCar, CartPole, and Reacher
210 environments operate within a bounded domain and thus inherently satisfy Assumption 2.4. The
211 swimmer, racecar, and soft arm can operate in an unbounded domain but have a cost function that
212 penalizes the distance between the system’s state \mathbf{x}_t and a target state \mathbf{x}^* . Therefore, the cost
213 encourages the system to move towards the target and remain within a bounded domain, as elaborated
214 on further in Appendix A.4.

215 **Baselines** In this work, we focus on model-based RL (MBRL) algorithms due to their sample
216 efficiency. To this end, we consider common techniques for planning with unknown dynamics, such
217 as planning with the mean, trajectory sampling (Chua et al., 2018), and Thompson sampling (Osband
218 & Van Roy, 2017). We adapt these three for our setting similar to as discussed in Section 3.2. For
219 all experiments with probabilistic ensembles, we consider TS1 from Chua et al. (2018) for trajectory
220 sampling, and for the GP experiment, we use distribution sampling from Chua et al. (2018). We

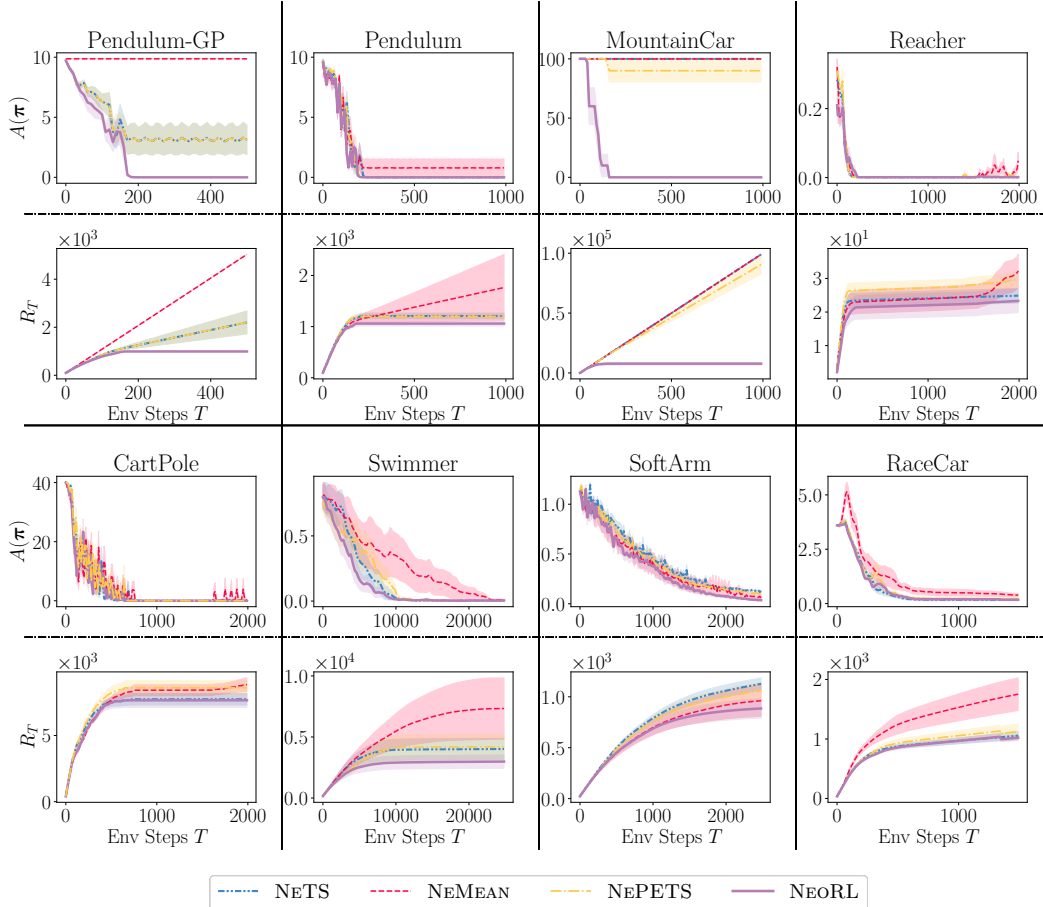


Figure 1: Average reward $A(\pi)$ and cumulative regret R_T over ten different seeds for all environments. We report the mean performance with one standard error as shaded regions. During all experiments, the environment is never reset. For all baselines, we model the dynamics with probabilistic ensembles, except in the Pendulum-GP experiment, where GPs are used instead. NEORL significantly outperforms all baselines and converges to the optimal average reward, $A(\pi^*) = 0$, showing sublinear cumulative regret R_T for all environments.

221 call the three baselines NEMEAN (nonepisodic mean), NEPETS (nonepisodic PETS), and NETS
 222 (nonepisodic Thompson sampling). NEMEAN and NEPETS are greedy w.r.t. the current estimate
 223 of the dynamics, i.e., do not explicitly encourage exploration. In our experiments, we show that
 224 being greedy does not suffice to converge to the optimal average cost, that is, obtain sublinear regret.

225 **Convergence to the optimal average cost** In Figure 1 we report the normalized average cost
 226 and cumulative regret of NEORL, NEMEAN, NEPETS, and NETS. The normalized average cost
 227 is defined such that $A(\pi^*) = 0$ for all environments. We observe that NEMEAN fails to converge
 228 to the optimal average cost for the Pendulum-v1 environment for both probabilistic ensembles and
 229 a GP model. It also fails to solve the MountainCar environment and is unstable for the Reacher
 230 and CartPole. In general, NEMEAN performs the worst among all methods. This is similar to the
 231 episodic case, where using the mean model often leads to the policy “overfitting” to the model
 232 inaccuracies (Chua et al., 2018). NEPETS performs better than the mean, however still significantly
 233 worse than NEORL. Even in the episodic setting, PETS tends to underexplore (Curi et al., 2020). We
 234 observe the same for the nonepisodic case, especially for the MountainCar task, which is a challenging
 235 RL environment with a sparse cost. Here NEPETS is also not able to achieve the optimal average
 236 cost and thus does not have sublinear cumulative regret. NETS performs similarly to NEPETS and is
 237 also not able to solve the MountainCar task.

238 NEORL performs the best among the baselines and converges to the optimal average cost achieving
 239 sublinear cumulative regret using only $\sim 10^3$ environment interactions. Moreover, this observation is
 240 consistent between different dynamics models (GPs and probabilistic ensembles) and environments.

241 Even in environments that are unbounded, i.e., Swimmer, SoftArm, and RaceCar, we observe that
 242 NEORL converges to the optimal average cost the fastest. We believe this is due to the MPC, which
 243 encourages the system to move closer to the target.

244 5 Related Work

245 **Average cost RL for finite state-action spaces** A significant amount of work studies the average
 246 cost/reward RL setting for finite-state action spaces. Moreover, seminal algorithms such as E^3 (Kearns
 247 & Singh, 2002) and R-max (Brafman & Tenenbholz, 2002) have established PAC bounds for
 248 the nonepisodic setting. These bounds are further improved for communicating MDPs by the
 249 UCRL2 (Jaksch et al., 2010) algorithm, which, similar to NEORL, is based on the optimism in the
 250 face of uncertainty paradigm and picks policies that are optimistic w.r.t. to the estimated dynamics.
 251 Their result is extended for weakly-communicating MDPs by REGAL (Bartlett & Tewari, 2012),
 252 similar results are derived for Thompson sampling based exploration (Ouyang et al., 2017), and
 253 for factored-MDP (Xu & Tewari, 2020). Albeit the significant amount of work for the finite case,
 254 progress for continuous state-action spaces has mostly been limited to linear dynamical systems.

255 **Nonepisodic RL for linear systems** There is a large body of work for nonepisodic learning with
 256 linear systems (Abbasi-Yadkori & Szepesvári, 2011; Cohen et al., 2019; Simchowitz & Foster,
 257 2020; Dean et al., 2020; Lale et al., 2020; Faradonbeh et al., 2020; Abeille & Lazaric, 2020; Treven
 258 et al., 2021). For linear systems with quadratic costs, the average reward problem, also known as
 259 the linear quadratic-Gaussian (LQG), has a closed-form solution which is obtained via the Riccati
 260 equations (Anderson & Moore, 2007). Moreover, for LQG, stability and optimality are intertwined,
 261 making studying linear systems much easier than their nonlinear counterpart. For studying nonlinear
 262 systems, additional assumptions on their stability are usually made.

263 **Nonepisodic RL beyond linear systems** In the case of nonlinear systems, guarantees have mostly
 264 been established for the episodic setting (Mania et al., 2020; Kakade et al., 2020; Curi et al., 2020;
 265 Wagenmaker et al., 2023; Sukhija et al., 2024; Treven et al., 2024). Only a few works consider the
 266 nonepisodic/single-trajectory case. For instance, Foster et al. (2020); Sattar & Oymak (2022) study the
 267 problem of system identification of a closed-loop globally exponentially stable dynamical system from
 268 a single trajectory. Lale et al. (2021) study the nonepisodic setting for nonlinear systems with MPC.
 269 Moreover, they consider finite-order or exponentially fading NARX systems that lie in the RKHS
 270 of infinitely smooth functions, which they further approximate with random Fourier features (Rahimi
 271 & Recht, 2007) ϕ with feature size D . Further, they assume access to bounded persistently exciting
 272 inputs w.r.t. the feature matrix $\Phi_t \Phi_t^\top$. This assumption is generally tough to verify and common exci-
 273 tation strategies such as random exploration often don't perform well for nonlinear systems (Sukhija
 274 et al., 2024). Further, the algorithm acts greedily w.r.t. the estimated dynamics, akin to NEMEAN,
 275 and requires the feature size D to increase with the horizon T . They give a regret bound of $\mathcal{O}(T^{2/3})$
 276 where the regret is measured w.r.t. to the oracle MPC with access to the true dynamics. Lale et al.
 277 (2021) also assume exponential input-to-output stability of the system to avoid blow-up during explo-
 278 ration. Our work considers more general RKHS, does not require apriori knowledge of persistently
 279 exciting inputs, and gives a regret bound of $\mathcal{O}(\beta_T \sqrt{TT_T})$ w.r.t. the optimal average cost criterion.
 280 Moreover, our regret bound is similar to the ones obtained for nonlinear systems in the episodic case
 281 and Gaussian process bandits (Srinivas et al., 2012; Chowdhury & Gopalan, 2017; Scarlett et al.,
 282 2017). To the best of our knowledge, we are the first to give such a regret bound for nonlinear systems.

283 6 Conclusion

284 We propose, NEORL, a novel model-based RL algorithm for the nonepisodic setting with nonlinear
 285 dynamics and continuous state and action spaces. NEORL seeks for average-cost optimal policies
 286 and leverages the model's epistemic uncertainty to perform optimistic exploration. Similar to the
 287 episodic case (Kakade et al., 2020; Curi et al., 2020), we provide a regret bound for NEORL of
 288 $\mathcal{O}(\beta_T \sqrt{TT_T})$ for Gaussian process dynamics. To our knowledge, we are the first to obtain this result
 289 in the nonepisodic setting. We compare NEORL to other model-based RL methods on standard
 290 deep RL benchmarks. Our experiments demonstrate that NEORL, converges to the optimal average
 291 cost of $A(\pi^*) = 0$ across all environments, suffering sublinear regret even when Bayesian neural
 292 networks are used to model the dynamics. Moreover, NEORL outperforms all our baselines across
 293 all environments requiring only $\sim 10^3$ samples for learning.

294 **References**

- 295 Abbasi-Yadkori, Y. and Szepesvári, C. Regret bounds for the adaptive control of linear quadratic systems. In
296 *Conference on Learning Theory*, 2011.
- 297 Abeille, M. and Lazaric, A. Efficient optimistic exploration in linear-quadratic regulators via lagrangian
298 relaxation. In *International Conference on Machine Learning*, 2020.
- 299 Anderson, B. D. and Moore, J. B. *Optimal control: linear quadratic methods*. Courier Corporation, 2007.
- 300 Arapostathis, A., Borkar, V. S., Fernández-Gaucherand, E., Ghosh, M. K., and Marcus, S. I. Discrete-time
301 controlled markov processes with average cost criterion: A survey. *SIAM Journal on Control and Optimization*,
302 1993.
- 303 Bartlett, P. L. and Tewari, A. Regal: A regularization based algorithm for reinforcement learning in weakly
304 communicating mdps. *arXiv preprint arXiv:1205.2661*, 2012.
- 305 Brafman, R. I. and Tennenholtz, M. R-max-a general polynomial time algorithm for near-optimal reinforcement
306 learning. *Journal of Machine Learning Research*, 2002.
- 307 Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym.
308 *arXiv preprint arXiv:1606.01540*, 2016.
- 309 Chowdhury, S. R. and Gopalan, A. On kernelized multi-armed bandits. In *ICML*, 2017.
- 310 Chua, K., Calandra, R., McAllister, R., and Levine, S. Deep reinforcement learning in a handful of trials using
311 probabilistic dynamics models. In *NeurIPS*, 2018.
- 312 Cohen, A., Koren, T., and Mansour, Y. Learning linear-quadratic regulators efficiently with only \sqrt{T} regret. In
313 *International Conference on Machine Learning*, 2019.
- 314 Curi, S., Berkenkamp, F., and Krause, A. Efficient model-based reinforcement learning through optimistic policy
315 search and planning. *NeurIPS*, 33:14156–14170, 2020.
- 316 Dean, S., Mania, H., Matni, N., Recht, B., and Tu, S. On the sample complexity of the linear quadratic regulator.
317 *Foundations of Computational Mathematics*, 20(4):633–679, 2020.
- 318 Faradonbeh, M. K. S., Tewari, A., and Michailidis, G. Optimism-based adaptive regulation of linear-quadratic
319 systems. *IEEE Transactions on Automatic Control*, 2020.
- 320 Foster, D., Sarkar, T., and Rakhlin, A. Learning nonlinear dynamical systems from a single trajectory. In
321 *Learning for Dynamics and Control*, 2020.
- 322 García, C. E., Prett, D. M., and Morari, M. Model predictive control: Theory and practice - a survey. *Automatica*,
323 pp. 335–348, 1989.
- 324 Hairer, M. and Mattingly, J. C. Yet another look at harris’ ergodic theorem for markov chains. In *Seminar on*
325 *Stochastic Analysis, Random Fields and Applications VI: Centro Stefano Franscini, Ascona, May 2008*, pp.
326 109–117. Springer, 2011.
- 327 Jaksch, T., Ortner, R., and Auer, P. Near-optimal regret bounds for reinforcement learning. *Journal of Machine*
328 *Learning Research*, 2010.
- 329 Kabzan, J., Valls, M. I., Reijgwart, V. J., Hendriks, H. F., Ehmke, C., Prajapat, M., Bühler, A., Gosala, N., Gupta,
330 M., Sivanesan, R., et al. Amz driverless: The full autonomous racing system. *Journal of Field Robotics*, 2020.
- 331 Kakade, S., Krishnamurthy, A., Lowrey, K., Ohnishi, M., and Sun, W. Information theoretic regret bounds for
332 online nonlinear control. *NeurIPS*, 33:15312–15325, 2020.
- 333 Kearns, M. and Singh, S. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 2002.
- 334 Khalil, H. K. *Nonlinear control*, volume 406. Pearson New York, 2015.
- 335 Kuleshov, V., Fenner, N., and Ermon, S. Accurate uncertainties for deep learning using calibrated regression. In
336 *ICML*, pp. 2796–2804. PMLR, 2018.
- 337 Lale, S., Azizzadenesheli, K., Hassibi, B., and Anandkumar, A. Logarithmic regret bound in partially observable
338 linear dynamical systems. *Advances in Neural Information Processing Systems*, 2020.
- 339 Lale, S., Azizzadenesheli, K., Hassibi, B., and Anandkumar, A. Model learning predictive control in nonlinear
340 dynamical systems. In *Conference on Decision and Control (CDC)*. IEEE, 2021.

341 Mania, H., Jordan, M. I., and Recht, B. Active learning for nonlinear system identification with guarantees.
342 *arXiv preprint arXiv:2006.10277*, 2020.

343 Meyn, S. P. and Tweedie, R. L. *Markov chains and stochastic stability*. Springer Science & Business Media,
344 2012.

345 Osband, I. and Van Roy, B. Why is posterior sampling better than optimism for reinforcement learning? In
346 *International conference on machine learning*, 2017.

347 Ouyang, Y., Gagrani, M., Nayyar, A., and Jain, R. Learning unknown markov decision processes: A thompson
348 sampling approach. *Advances in neural information processing systems*, 30, 2017.

349 Pinneri, C., Sawant, S., Blaes, S., Achterhold, J., Stueckler, J., Rolinek, M., and Martius, G. Sample-efficient
350 cross-entropy method for real-time planning. In *CORL*, Proceedings of Machine Learning Research, pp.
351 1049–1065, 2021.

352 Puterman, M. L. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons,
353 2014.

354 Rahimi, A. and Recht, B. Random features for large-scale kernel machines. *Advances in neural information
355 processing systems*, 20, 2007.

356 Rothfuss, J., Sukhija, B., Birchler, T., Kassraie, P., and Krause, A. Hallucinated adversarial control for
357 conservative offline policy evaluation. *UAI*, 2023.

358 Sattar, Y. and Oymak, S. Non-asymptotic and accurate learning of nonlinear dynamical systems. *Journal of
359 Machine Learning Research*, 2022.

360 Scarlett, J., Bogunovic, I., and Cevher, V. Lower bounds on regret for noisy Gaussian process bandit optimization.
361 In *Conference on Learning Theory*, 2017.

362 Sharma, A., Xu, K., Sardana, N., Gupta, A., Hausman, K., Levine, S., and Finn, C. Autonomous reinforcement
363 learning: Formalism and benchmarking. *arXiv preprint arXiv:2112.09605*, 2021.

364 Simchowitz, M. and Foster, D. Naive exploration is optimal for online lqr. In *International Conference on
365 Machine Learning*. PMLR, 2020.

366 Srinivas, N., Krause, A., Kakade, S. M., and Seeger, M. W. Information-theoretic regret bounds for gaussian
367 process optimization in the bandit setting. *IEEE Transactions on Information Theory*, 2012.

368 Sukhija, B., Treven, L., Sancaktar, C., Blaes, S., Coros, S., and Krause, A. Optimistic active exploration of
369 dynamical systems. *NeurIPS*, 2024.

370 Sussex, S., Makarova, A., and Krause, A. Model-based causal bayesian optimization. In *ICLR*, May 2023.

371 Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D. d. L., Budden, D., Abdolmaleki, A., Merel, J.,
372 Lefrancq, A., et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.

373 Tekinalp, A., Kim, S. H., Bhosale, Y., Parthasarathy, T., Naughton, N., Albazroun, A., Joon, R., Cui, S.,
374 Nasiriziba, I., Stölzle, M., Shih, C.-H. C., and Gazzola, M. Gazzolalab/pyelastica: v0.3.2, 2024.

375 Treven, L., Curi, S., Mutnỳ, M., and Krause, A. Learning stabilizing controllers for unstable linear quadratic
376 regulators from a single trajectory. In *Learning for Dynamics and Control*, 2021.

377 Treven, L., Hübotter, J., Sukhija, B., Dörfler, F., and Krause, A. Efficient exploration in continuous-time
378 model-based reinforcement learning. *NeurIPS*, 2024.

379 Vakili, S., Khezeli, K., and Picheny, V. On information gain and regret bounds in gaussian process bandits. In
380 *AISTATS*, 2021.

381 Wagenmaker, A., Shi, G., and Jamieson, K. Optimal exploration for model-based rl in nonlinear systems. *arXiv
382 preprint arXiv:2306.09210*, 2023.

383 Xu, Z. and Tewari, A. Reinforcement learning in factored mdps: Oracle-efficient algorithms and tighter regret
384 bounds for the non-episodic setting. *Advances in Neural Information Processing Systems*, 2020.

385 Appendices

386 A Proofs

387 In this section, we prove Theorem 2.6 and Theorem 3.1. First, we start with the proof of Corollary 2.5.

388 *Proof of Corollary 2.5.* We first analyze the following term $\mathbb{E}_{\mathbf{w}}[V(\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \mathbf{w}) -$
 389 $V(\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}_s(\mathbf{x})) + \mathbf{w})]$ for any $\boldsymbol{\pi} \in \Pi$.

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{w}}[V(\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \mathbf{w}) - V(\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}_s(\mathbf{x})) + \mathbf{w})] \\
 & \leq \mathbb{E}_{\mathbf{w}}[\kappa(\|\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \mathbf{w} - (\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}_s(\mathbf{x})) + \mathbf{w})\|)] \quad (\text{Uniform continuity of } V) \\
 & = \kappa(\|\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) - \mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}_s(\mathbf{x}))\|) \\
 & \leq \kappa(\kappa_{\mathbf{f}^*}(\|\boldsymbol{\pi}(\mathbf{x}) - \boldsymbol{\pi}_s(\mathbf{x})\|)) \quad (\text{Uniform continuity of } \mathbf{f}^*) \\
 & \leq \kappa(\kappa_{\mathbf{f}^*}(2u_{\max})). \quad (\text{Bounded inputs})
 \end{aligned}$$

390 Therefore,

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x}'|\boldsymbol{\pi}, \mathbf{x}}[V(\mathbf{x}')] &= \mathbb{E}_{\mathbf{w}}[V(\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \mathbf{w})] \\
 & \leq \mathbb{E}_{\mathbf{w}}[V(\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}_s(\mathbf{x})) + \mathbf{w})] + \kappa(\kappa_{\mathbf{f}^*}(2u_{\max})) \\
 & = \mathbb{E}_{\mathbf{x}'|\boldsymbol{\pi}_s, \mathbf{x}}[V(\mathbf{x}')] + \kappa(\kappa_{\mathbf{f}^*}(2u_{\max})) \\
 & \leq \gamma V(\mathbf{x}) + K + \kappa(\kappa_{\mathbf{f}^*}(2u_{\max})) \\
 & = \gamma V(\mathbf{x}) + \tilde{K} \quad (\tilde{K} = K + \kappa(\kappa_{\mathbf{f}^*}(2u_{\max})))
 \end{aligned}$$

391 Hence, V satisfies the drift condition for $\boldsymbol{\pi}$. Furthermore, since V also satisfies positive definiteness
 392 by assumption, the bounded energy condition holds for all $\boldsymbol{\pi} \in \Pi$. \square

393 A.1 Proof of Theorem 2.6

394 For proving Theorem 2.6, we invoke the results from (Hairer & Mattingly, 2011, Theorem 1.2 – 1.3).
 395 For this we require that the Markov chain induced by a policy $\boldsymbol{\pi}$ satisfies the drift condition. In our
 396 setting, this corresponds to Assumption 2.4. Next, we show that the chain satisfies the following
 397 minorisation condition.

398 **Lemma A.1** (Minorisation condition). *Consider the system in Equation (1) and let Assump-*
 399 *tion 2.1 – 2.4 hold. Let P^π denote the transition kernel for the policy $\boldsymbol{\pi} \in \Pi$, i.e., $P^\pi(\mathbf{x}, \mathcal{A}) =$
 400 $\mathbb{P}(\mathbf{x}' \in \mathcal{A} | \mathbf{x}, \boldsymbol{\pi}(\mathbf{x}))$. Then, for all $\boldsymbol{\pi} \in \Pi$, exists a constant $\alpha \in (0, 1)$ and a probability measure
 401 $\zeta(\cdot)$ s.t.,*

$$\inf_{\mathbf{x} \in \mathcal{C}} P^\pi(\mathbf{x}, \cdot) \geq \alpha \zeta(\cdot) \quad (9)$$

402 with $\mathcal{C} \stackrel{\text{def}}{=} \{\mathbf{x} \in \mathcal{X}; V^\pi(\mathbf{x}) \leq R\}$ for some $R > 2K/1-\gamma$

403 *Proof.* We prove it in 3 steps. First, we show that \mathcal{C} is contained in a compact domain. From the
 404 Assumption 2.4 we pick the function $\xi \in \mathcal{K}_\infty$. Since $C_l \xi(0) = 0$, $\lim_{s \rightarrow \infty} \xi(s) = +\infty$ and $C_l \xi$ is
 405 continuous, there exists M such that $C_l \xi(M) = R$. Then for $\|\mathbf{x}\| > M$ we have:

$$V^\pi(\mathbf{x}) \geq C_l \xi(\|\mathbf{x}\|) > \xi(M) = R.$$

406 Therefore we have: $\mathcal{C} \subseteq \mathcal{B}(\mathbf{0}, M) \stackrel{\text{def}}{=} \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{0}\| \leq M\}$. In the second step we show that
 407 $\mathbf{f}(\mathcal{C}, \boldsymbol{\pi}(\mathcal{C}))$ is bounded, in particular we show that there exists $B > 0$ such that: $\mathbf{f}(\mathcal{C}, \boldsymbol{\pi}(\mathcal{C})) \subseteq$
 408 $\mathcal{B}(\mathbf{0}, B)$. This is true since continuous image of compact set is compact and the observation:

$$\mathcal{C} \subseteq \mathcal{B}(\mathbf{0}, M) \implies \mathbf{f}(\mathcal{C}, \boldsymbol{\pi}(\mathcal{C})) \subseteq \mathbf{f}(\mathcal{B}(\mathbf{0}, M), \boldsymbol{\pi}(\mathcal{B}(\mathbf{0}, M))).$$

409 Since $\mathbf{f}(\mathcal{B}(\mathbf{0}, M), \boldsymbol{\pi}(\mathcal{B}(\mathbf{0}, M)))$ is compact there exists B such that $\mathbf{f}(\mathcal{C}, \boldsymbol{\pi}(\mathcal{C})) \subseteq \mathcal{B}(\mathbf{0}, B)$. In
 410 the last step we prove that $\alpha \stackrel{\text{def}}{=} 2^{-d_{\mathbf{x}}} e^{-B^2/\sigma^2}$ and ζ with law of $\mathcal{N}\left(0, \frac{\sigma^2}{2}\right)$ satisfy condition of

411 Lemma A.1. It is enough to show that $\forall \boldsymbol{\mu} \in \mathcal{B}(\mathbf{0}, B), \forall \mathbf{x} \in \mathbb{R}^{d_{\mathbf{x}}}$ we have:

$$\alpha \frac{1}{(2\pi)^{\frac{d_{\mathbf{x}}}{2}} \left(\frac{\sigma^2}{2}\right)^{\frac{d_{\mathbf{x}}}{2}}} e^{-\frac{\|\mathbf{x}\|^2}{\sigma^2}} \leq \frac{1}{(2\pi)^{\frac{d_{\mathbf{x}}}{2}} (\sigma^2)^{\frac{d_{\mathbf{x}}}{2}}} e^{-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}}$$

412 which can be proven with simple algebraic manipulations. \square

413 Through the minorisation condition and Assumption 2.4, we can prove the ergodicity of the closed-
414 loop system for a given policy $\pi \in \Pi$.

415 **Theorem A.2** (Ergodicity of closed-loop system). *Let Assumption 2.1 – 2.4, consider any probability*
416 *measures ζ_1, ζ_2 , and $\theta > 0$, define $P^\pi \zeta, \|\varphi\|_{1+\theta V^\pi}, \rho_\theta^\pi$ as*

$$\begin{aligned} (P^\pi \zeta)(\mathcal{A}) &= \int_{\mathcal{X}} P^\pi(\mathbf{x}, \mathcal{A}) \zeta(d\mathbf{x}) \\ \|\varphi\|_{1+\theta V^\pi} &= \sup_{\mathbf{x} \in \mathcal{X}} \frac{|\varphi(\mathbf{x})|}{1 + \theta V^\pi(\mathbf{x})} \\ \rho_\theta^\pi(\zeta_1, \zeta_2) &= \sup_{\varphi: \|\varphi\|_{1+\theta V^\pi} \leq 1} \int_{\mathcal{X}} \varphi(\mathbf{x})(\zeta_1 - \zeta_2)(d\mathbf{x}) = \int_{\mathcal{X}} (1 + \theta V^\pi(\mathbf{x})) |\zeta_1 - \zeta_2|(d\mathbf{x}). \end{aligned}$$

417 *We have for all $\pi \in \Pi$, that P^π admits a unique invariant measure \bar{P}^π . Furthermore, there exist*
418 *constants $C_1 > 0, \theta > 0, \lambda \in (0, 1)$ such that*

$$\rho_\theta^\pi(P^\pi \zeta_1, P^\pi \zeta_2) \leq \lambda \rho_\theta^\pi(\zeta_1, \zeta_2) \quad (1)$$

$$\|\mathbb{E}_{\mathbf{x} \sim (P^\pi)^t} [\varphi(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [\varphi(\mathbf{x})]\|_{1+V^\pi} \leq C_1 \lambda^t \|\varphi - \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [\varphi(\mathbf{x})]\|_{1+V^\pi}. \quad (2)$$

419 *holds for every measurable function $\varphi : \mathcal{X} \rightarrow \mathcal{R}$ with $\|\varphi\|_{1+V^\pi} < \infty$. Here $(P^\pi)^t$ denotes the*
420 *t -step transition kernel under the policy π .*

421 *Moreover, $\theta = \alpha_0/K$, and*

$$\lambda = \max \left\{ 1 - (\alpha - \alpha_0), \frac{2 + R/K\alpha_0\gamma_0}{2 + R/K\alpha_0} \right\} \quad (10)$$

422 *for any $\alpha_0 \in (0, \alpha)$ and $\gamma_0 \in (\gamma + 2K/R, 1)$.*

423 *Proof.* From Assumption 2.4, we have a value function for each policy that satisfies the drift condition.
424 Furthermore, in Lemma A.1 we show that our system also satisfies the minorisation condition for all
425 policies. Under these conditions, we can use the results from [Hairer & Mattingly \(2011, Theorem 1.2.](#)
426 [– 1.3.\)](#). \square

427 Note that $\|\cdot\|_{1+\theta V^\pi}$ represents a family of equivalent norms for any $\theta > 0$. Now we prove Theo-
428 rem 2.6.

429 *Proof of Theorem 2.6.* From Theorem A.2, we have

$$\rho_\theta^\pi((P^\pi)^{t+1}, (P^\pi)^t) = \rho_\theta^\pi(P^\pi (P^\pi)^t, P^\pi (P^\pi)^{t-1}) \leq \lambda^t \rho_\theta^\pi(P^\pi \delta_{\mathbf{x}_0}, \delta_{\mathbf{x}_0}),$$

430 where $\delta_{\mathbf{x}_0}$ is the dirac measure. Therefore, $(P^\pi)^t$ is a Cauchy sequence. Furthermore, ρ_θ^π is complete
431 for the set of probability measures integrating V , thus $\rho_\theta^\pi((P^\pi)^t, P^\pi) \rightarrow 0$ for $t \rightarrow \infty$ (c.f., [Hairer &](#)
432 [Mattingly \(2011\)](#) for more details). In particular, we have for φ such that $\|\varphi\|_{1+\theta V^\pi} \leq 1$,

$$\lim_{t \rightarrow \infty} \int_{\mathcal{X}} \varphi(\mathbf{x})(P^\pi)^t(d\mathbf{x}) = \int_{\mathcal{X}} \varphi(\mathbf{x}) \bar{P}^\pi(d\mathbf{x}).$$

433 Note that since all $\|\cdot\|_{1+\theta V^\pi}$ norms are equivalent for $\theta > 0$, if $\|c\|_{1+V^\pi} \leq C$ (Assumption 2.4),
434 then $\|c\|_{1+\theta V^\pi} \leq C'$ for some $C' \in (0, \infty)$. Furthermore, note that $c(\cdot) \geq 0$. Therefore,

$$\begin{aligned} \int_{\mathcal{X}} c(\mathbf{x}) \bar{P}^\pi(d\mathbf{x}) &= \lim_{t \rightarrow \infty} \int_{\mathcal{X}} c(\mathbf{x})(P^\pi)^t(d\mathbf{x}) \\ &\leq C \lim_{t \rightarrow \infty} \int_{\mathcal{X}} (1 + V^\pi(\mathbf{x}))(P^\pi)^t(d\mathbf{x}) \\ &= C + C \lim_{t \rightarrow \infty} \mathbb{E}_{\mathbf{x} \sim (P^\pi)^t} [V^\pi(\mathbf{x})] \\ &= C + C \lim_{t \rightarrow \infty} \mathbb{E}_{\mathbf{x} \sim (P^\pi)^{t-1}} [\mathbb{E}_{\mathbf{x}' \sim (P^\pi)} [V^\pi(\mathbf{x}') | \mathbf{x}]] \\ &\leq C + C \left(\lim_{t \rightarrow \infty} \gamma \mathbb{E}_{\mathbf{x} \sim (P^\pi)^{t-1}} [V^\pi(\mathbf{x})] + K \right) \quad (\text{Assumption 2.4}) \end{aligned}$$

$$\begin{aligned}
&\leq C + C \lim_{t \rightarrow \infty} \gamma^t V^\pi(\mathbf{x}_0) + K \frac{1 - \gamma^t}{1 - \gamma} \\
&= C \left(1 + K \frac{1}{1 - \gamma} \right)
\end{aligned}$$

435 In summary, we have $\mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [c(\mathbf{x})] \leq C \left(1 + K \frac{1}{1 - \gamma} \right)$

436 Consider any $t > 0$, and note that from Theorem A.2 we have

$$\begin{aligned}
\left\| \mathbb{E}_{\mathbf{x} \sim (P^\pi)^t} [c(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [c(\mathbf{x})] \right\|_{1+V^\pi} &= \sup_{\mathbf{x}_0 \in \mathcal{X}} \frac{|\mathbb{E}_{\mathbf{x} \sim (P^\pi)^t} [c(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [c(\mathbf{x})]|}{1 + V^\pi(\mathbf{x}_0)} \\
&\leq C_1 \lambda^t \|c - \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [c(\mathbf{x})]\|_{1+V^\pi} \quad (\text{Theorem A.2}) \\
&\leq C_1 \lambda^t \|c\|_{1+V^\pi} + C_1 \lambda^t \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [c(\mathbf{x})] \\
&= C_2 \lambda^t,
\end{aligned}$$

437 where $C_2 = C_1 (\|c\|_{1+V^\pi} + CK \frac{1}{1-\gamma})$.

438 Moreover, since the inequality holds for all \mathbf{x}_0 , we have

$$\frac{|\mathbb{E}_{\mathbf{x} \sim (P^\pi)^t} [c(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [c(\mathbf{x})]|}{1 + V^\pi(\mathbf{x}_0)} \leq C_2 \lambda^t.$$

439 In summary,

$$|\mathbb{E}_{\mathbf{x} \sim (P^\pi)^t} [c(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [c(\mathbf{x})]| \leq C_2 (1 + V^\pi(\mathbf{x}_0)) \lambda^t.$$

440 Consider any $T \geq 0$, and define with $\bar{c} = \mathbb{E}_{\mathbf{x} \sim \bar{P}^\pi} [c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}))]$.

$$\begin{aligned}
\mathbb{E}_\pi \left[\sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t) - \bar{c} \right] &= \sum_{t=0}^{T-1} \mathbb{E}_{(P^\pi)^t} [c(\mathbf{x}_t, \mathbf{u}_t)] - \bar{c} \\
&\leq \sum_{t=0}^{T-1} |\mathbb{E}_{(P^\pi)^t} [c(\mathbf{x}_t, \mathbf{u}_t)] - \bar{c}| \\
&\leq C_2 (1 + V^\pi(\mathbf{x}_0)) \sum_{t=0}^{T-1} \lambda^t \\
&= C_2 (1 + V^\pi(\mathbf{x}_0)) \frac{1 - \lambda^T}{1 - \lambda}
\end{aligned}$$

441 Hence, we have

$$\lim_{T \rightarrow \infty} \left| \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t) - \bar{c} \right] \right| \leq C_2 (1 + V^\pi(\mathbf{x}_0)) \frac{1}{1 - \lambda},$$

442 and for any \mathbf{x}_0 in a compact subset of \mathcal{X}

$$\lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t) - \bar{c} \right] = 0.$$

443 Moreover,

$$|B(\boldsymbol{\pi}, \mathbf{x}_0)| \leq C_2 (1 + V^\pi(\mathbf{x}_0)) \frac{1}{1 - \lambda}.$$

444

□

445 Another interesting inequality that follows from the proof above is the difference in bias inequality.

$$|\mathbb{E}_{\mathbf{x}_0 \sim \zeta_1} [B(\boldsymbol{\pi}, \mathbf{x}_0)] - \mathbb{E}_{\mathbf{x}_0 \sim \zeta_2} [B(\boldsymbol{\pi}, \mathbf{x}_0)]| \leq \frac{C_3}{1 - \lambda} \int_{\mathcal{X}} (1 + V^\pi(\mathbf{x})) |\zeta_1 - \zeta_2| (d\mathbf{x})$$

446 for all probability measures ζ_1, ζ_2 . To show this holds, define $C' = \max_{\pi \in \Pi} \|c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}))\|_{1+\theta V^\pi}$.
 447 Furthermore, note that $C' < \infty$ from Assumption 2.4 and $\|c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}))\|_{1+\theta V^\pi} \leq 1$.

$$\begin{aligned}
 & \left| \mathbb{E}_{\mathbf{x} \sim (P^\pi)^t \zeta_1} c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) - \mathbb{E}_{\mathbf{x} \sim (P^\pi)^t \zeta_2} c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) \right| = \left| \int_{\mathcal{X}} c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) ((P^\pi)^t \zeta_1 - (P^\pi)^t \zeta_2)(d\mathbf{x}) \right| \\
 & = C' \left| \int_{\mathcal{X}} \frac{1}{C'} c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) ((P^\pi)^t \zeta_1 - (P^\pi)^t \zeta_2)(d\mathbf{x}) \right| \\
 & \leq C' \sup_{\varphi: \|\varphi\|_{1+\theta V^\pi} \leq 1} \int_{\mathcal{X}} \varphi(\mathbf{x}) ((P^\pi)^t \zeta_1 - (P^\pi)^t \zeta_2)(d\mathbf{x}) = C' \rho_\theta^\pi((P^\pi)^t \zeta_1, (P^\pi)^t \zeta_2) \\
 & \leq C' \lambda \rho_\theta^\pi((P^\pi)^{t-1} \zeta_1, (P^\pi)^{t-1} \zeta_2) \quad (\text{Theorem A.2}) \\
 & \leq C' \lambda^t \rho_\theta^\pi(\zeta_1, \zeta_2).
 \end{aligned}$$

448 Also, note that there exists $C_\theta \in (0, \infty)$ such that $C_\theta \|\varphi\|_{1+\theta V^\pi} \geq \|\varphi\|_{1+V^\pi}$ due to the equivalence
 449 of the two norms.

$$\begin{aligned}
 \rho_\theta^\pi(\zeta_1, \zeta_2) & = \sup_{\varphi: \|\varphi\|_{1+\theta V^\pi} \leq 1} \int_{\mathcal{X}} \varphi(\mathbf{x}) (\zeta_1 - \zeta_2)(d\mathbf{x}) \\
 & \leq \sup_{\varphi: \|\varphi\|_{1+V^\pi} \leq C_\theta} \int_{\mathcal{X}} \varphi(\mathbf{x}) (\zeta_1 - \zeta_2)(d\mathbf{x}) \\
 & = C_\theta \sup_{\varphi: \|\varphi\|_{1+V^\pi} \leq 1} \int_{\mathcal{X}} \varphi(\mathbf{x}) (\zeta_1 - \zeta_2)(d\mathbf{x}) \\
 & = C_\theta \rho_1^\pi(\zeta_1, \zeta_2)
 \end{aligned}$$

450 Therefore, for the bias we have

$$\begin{aligned}
 & \left| \mathbb{E}_{\mathbf{x}_0 \sim \zeta_1} [B(\boldsymbol{\pi}, \mathbf{x}_0)] - \mathbb{E}_{\mathbf{x}_0 \sim \zeta_2} [B(\boldsymbol{\pi}, \mathbf{x}_0)] \right| \\
 & \leq \lim_{T \rightarrow \infty} \sum_{t=0}^{T-1} \left| \mathbb{E}_{\mathbf{x} \sim (P^\pi)^t \zeta_1} c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) - \mathbb{E}_{\mathbf{x} \sim (P^\pi)^t \zeta_2} c(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) \right| \\
 & \leq C' \rho_\theta^\pi(\zeta_1, \zeta_2) \lim_{T \rightarrow \infty} \sum_{t=0}^{T-1} \lambda^t = \frac{C'}{1-\lambda} \rho_\theta^\pi(\zeta_1, \zeta_2) \\
 & \leq \frac{C' C_\theta}{1-\lambda} \rho_1^\pi(\zeta_1, \zeta_2) = \frac{C' C_\theta}{1-\lambda} \int_{\mathcal{X}} (1 + V^\pi(\mathbf{x})) |\zeta_1 - \zeta_2|(d\mathbf{x})
 \end{aligned}$$

451 Set $C_3 = C' C_\theta$.

452 A.2 Proof of bounded average cost for the optimistic system

453 In this section, we show that the results from Theorem 2.6 also transfer over to the optimistic
 454 dynamics.

455 **Theorem A.3** (Existence of Average Cost Solution for the Optimistic System). *Let Assumption 2.1 –*
 456 *2.8 hold. Consider any $n > 0$ and let $\boldsymbol{\pi}_n, \mathbf{f}_n$ denote the solution to Equation (6), $P^{\boldsymbol{\pi}_n, \mathbf{f}_n}$ its transition*
 457 *kernel. Then $P^{\boldsymbol{\pi}_n, \mathbf{f}_n}$ admits a unique invariant measure $\bar{P}^{\boldsymbol{\pi}_n, \mathbf{f}_n}$ and there exists $C_2, C_3 \in (0, \infty)$,*
 458 *$\hat{\lambda} \in (0, 1)$ such that*

459 Average Cost;

$$A(\boldsymbol{\pi}_n, \mathbf{f}_n) = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\boldsymbol{\pi}_n, \mathbf{f}_n} \left[\sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t) \right] = \mathbb{E}_{\mathbf{x} \sim \bar{P}^{\boldsymbol{\pi}_n, \mathbf{f}_n}} [c(\mathbf{x}, \boldsymbol{\pi}_n(\mathbf{x}))]$$

460 Bias Cost;

$$|B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_0)| = \left| \lim_{T \rightarrow \infty} \mathbb{E}_{\boldsymbol{\pi}_n, \mathbf{f}_n} \left[\sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t) - A(\boldsymbol{\pi}_n, \mathbf{f}_n) \right] \right| \leq C_2 (1 + V^{\boldsymbol{\pi}_n}(\mathbf{x}_0)) \frac{1}{1-\hat{\lambda}}$$

461 for all $\mathbf{x}_0 \in \mathcal{X}$.

462 Difference in Bias;

$$|\mathbb{E}_{\mathbf{x}_0 \sim \zeta_1}[B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_0)] - \mathbb{E}_{\mathbf{x}_0 \sim \zeta_2}[B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_0)]| \leq \frac{C_3}{1 - \lambda} \int_{\mathcal{X}} (1 + V^\pi(\mathbf{x})) |\zeta_1 - \zeta_2| (d\mathbf{x})$$

463 for all probability measures ζ_1, ζ_2 .

464 Theorem A.3 shows that the optimistic dynamics \mathbf{f}_n retain the boundedness property from the
 465 true dynamics \mathbf{f}^* and give a well-defined solution w.r.t. average cost and the bias cost. To prove
 466 Theorem A.3 we show that the optimistic system also satisfies the drift and minorisation condition.
 467 Then we can invoke the result from Hairer & Mattingly (2011) similar to the proof of Theorem 2.6.

468 **Lemma A.4** (Stability of optimistic system). *Let Assumption 2.1 – 2.8 hold, then we have with*
 469 *probability at least $1 - \delta$ for all $n \geq 0$, $\boldsymbol{\pi} \in \Pi$, $\mathbf{f} \in \mathcal{M}_n \cap \mathcal{M}_0$, that there exists a constant $\widehat{K} > 0$;*

$$\mathbb{E}_{\mathbf{x}'|\mathbf{x}, \mathbf{f}, \boldsymbol{\pi}}[V^\pi(\mathbf{x}')] \leq \gamma V^\pi(\mathbf{x}) + \widehat{K}.$$

470

471 *Proof.* Note, that V^π is uniformly continuous w.r.t. κ

$$|V^\pi(\mathbf{x}) - V^\pi(\mathbf{x}')| \leq \kappa(\|\mathbf{x} - \mathbf{x}'\|).$$

472 Furthermore, since $\mathbf{f} \in \mathcal{M}_n \cap \mathcal{M}_0$ and therefore $\mathbf{f} \in \mathcal{M}_0$, we have that there exists some $\boldsymbol{\eta} \in$
 473 $[-1, 1]^{d_x}$ such that

$$\mathbf{f}(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) = \boldsymbol{\mu}_0(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \beta_0 \boldsymbol{\sigma}_0(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) \boldsymbol{\eta}(\mathbf{x}).$$

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}}[V^\pi(\boldsymbol{\mu}_0(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \beta_0 \boldsymbol{\sigma}_0(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) \boldsymbol{\eta}(\mathbf{x}) + \mathbf{w})] - \mathbb{E}_{\mathbf{w}}[V^\pi(\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \mathbf{w})] \\ & \leq \kappa(\|\boldsymbol{\mu}_0(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \beta_0 \boldsymbol{\sigma}_0(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) \boldsymbol{\eta}(\mathbf{x}) - \mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}))\|) \\ & \leq \kappa(\|\boldsymbol{\mu}_0(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) - \mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}))\| + \|\beta_0 \boldsymbol{\sigma}_0(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) \boldsymbol{\eta}(\mathbf{x})\|) \\ & \leq \kappa\left(\left(1 + \sqrt{d_x}\right) \beta_0 \sqrt{d_x} \sigma_{\max}\right). \end{aligned} \quad (\text{Assumption 2.8})$$

474 Therefore,

$$\begin{aligned} \mathbb{E}_{\mathbf{x}'|\mathbf{x}, \mathbf{f}, \boldsymbol{\pi}}[V^\pi(\mathbf{x}')] & \leq \mathbb{E}_{\mathbf{x}'|\mathbf{x}, \mathbf{f}^*, \boldsymbol{\pi}}[V^\pi(\mathbf{x}')] + \kappa\left(\left(1 + \sqrt{d_x}\right) \beta_0 \sqrt{d_x} \sigma_{\max}\right) \\ & = \mathbb{E}_{\mathbf{x}'|\mathbf{x}, \boldsymbol{\pi}}[V^\pi(\mathbf{x}')] + \kappa\left(\left(1 + \sqrt{d_x}\right) \beta_0 \sqrt{d_x} \sigma_{\max}\right) \\ & \leq \gamma V^\pi(\mathbf{x}) + K + \kappa\left(\left(1 + \sqrt{d_x}\right) \beta_0 \sqrt{d_x} \sigma_{\max}\right). \end{aligned} \quad (\text{Assumption 2.4})$$

475 Define $\widehat{K} = K + \kappa\left(\left(1 + \sqrt{d_x}\right) \beta_0 \sqrt{d_x} \sigma_{\max}\right)$. □

476 **Lemma A.5** (Minorisation condition optimistic system). *Consider the system*

$$\mathbf{x}' = \mathbf{f}(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x})) + \mathbf{w}$$

477 for any $n \geq 0$, $\boldsymbol{\pi} \in \Pi$ and $\mathbf{f} \in \mathcal{M}_n \cap \mathcal{M}_0$. Let Assumption 2.1 – 2.8 hold. Let $P^{\boldsymbol{\pi}, \mathbf{f}}$ denote the
 478 transition kernel for the policy $\boldsymbol{\pi} \in \Pi$ i.e., $P^{\boldsymbol{\pi}, \mathbf{f}}(\mathbf{x}, \mathcal{A}) = \mathbb{P}(\mathbf{x}' \in \mathcal{A} | \mathbf{x}, \boldsymbol{\pi}(\mathbf{x}), \mathbf{f})$. Then, there exists
 479 a constant $\hat{\alpha} \in (0, 1)$ and a probability measure $\hat{\zeta}(\cdot)$ independent of n s.t.,

$$\inf_{\mathbf{x} \in \mathcal{C}} P^{\boldsymbol{\pi}, \mathbf{f}}(\mathbf{x}, \cdot) \geq \hat{\alpha} \hat{\zeta}(\cdot) \quad (11)$$

480 with $\mathcal{C} \stackrel{\text{def}}{=} \{\mathbf{x} \in \mathcal{X}; V^\pi(\mathbf{x}) < \hat{R}\}$ for some $\hat{R} > 2\widehat{K}/(1 - \gamma)$

481 *Proof.* First, we show that \mathcal{C} is contained in a compact domain. From the Assumption 2.4 we pick
 482 the function $\xi \in \mathcal{K}_\infty$. Since $C_l \xi(0) = 0$, $\lim_{s \rightarrow \infty} \xi(s) = +\infty$ and $C_l \xi$ is continuous, there exists
 483 M such that $C_l \xi(M) = \hat{R}$. Then for $\|\mathbf{x}\| > M$ we have:

$$V^\pi(\mathbf{x}) \geq C_l \xi(\|\mathbf{x}\|) > \xi(M) = \hat{R}.$$

484 Therefore we have: $\mathcal{C} \subseteq \mathcal{B}(\mathbf{0}, M) \stackrel{\text{def}}{=} \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{0}\| \leq M\}$. Since for any $\mathbf{x} \in \mathcal{C}$ we have
485 $\|\mathbf{f}(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}))\| \leq \|\mathbf{f}^*(\mathbf{x}, \boldsymbol{\pi}(\mathbf{x}))\| + \beta_0 \sigma_{\max}$. Since \mathbf{f}^* is continuous, there exists a B such that
486 $\mathbf{f}^*(\mathcal{C}, \boldsymbol{\pi}(\mathcal{C})) \subset \mathcal{B}(\mathbf{0}, B)$. Therefore we have: $\mathbf{f}(\mathcal{C}, \boldsymbol{\pi}(\mathcal{C})) \subset \mathcal{B}(\mathbf{0}, B_1)$, where $B_1 = B + \beta_0 \sigma_{\max}$.
487 In the last step we prove that $\alpha \stackrel{\text{def}}{=} 2^{-d_x} e^{-B_1^2/\sigma^2}$ and ζ with law of $\mathcal{N}\left(0, \frac{\sigma^2}{2}\right)$ satisfy condition of
488 Lemma A.1. It is enough to show that $\forall \boldsymbol{\mu} \in \mathcal{B}(\mathbf{0}, B_1), \forall \mathbf{x} \in \mathbb{R}^{d_x}$ we have:

$$\alpha \frac{1}{(2\pi)^{\frac{d_x}{2}} \left(\frac{\sigma^2}{2}\right)^{\frac{d_x}{2}}} e^{-\frac{\|\mathbf{x}\|^2}{\sigma^2}} \leq \frac{1}{(2\pi)^{\frac{d_x}{2}} (\sigma^2)^{\frac{d_x}{2}}} e^{-\frac{\|\mathbf{x}-\boldsymbol{\mu}\|^2}{2\sigma^2}}$$

489 which can be proven with simple algebraic manipulations. \square

490 *Proof of Theorem A.3.* As for the true system, the drift condition from Lemma A.4 and the mi-
491 norisation condition from Lemma A.5 are sufficient to show ergodicity of the optimistic system
492 (c.f., Theorem A.2 or Hairer & Mattingly (2011)). The rest of the proof is similar to Theorem 2.6. \square

493 A.3 Proof of Theorem 3.1

494 Since NEORL works in artificial episodes $n \in \{0, N-1\}$ of varying horizons H_n . We denote with
495 \mathbf{x}_k^n the state visited during episode n at time step $k \leq H_n$. Crucial, to our regret analysis is bounding
496 the first and second moment of $V^{\boldsymbol{\pi}^n}(\mathbf{x}_k^n)$ for all n, k . Given the nature of Assumption 2.4, this
497 requires analyzing geometric series. Thus, we start with the following elementary result of geometric
498 series.

499 **Corollary A.6.** *Consider the sequence $\{S_n\}_{n \geq 0}$ with $S_n \geq 0$ for all n . Let the following hold*

$$S_n \leq \rho S_{n-1} + C$$

500 for $\rho \in (0, 1)$ and $C > 0$. Then we have

$$S_n \leq \rho^n S_0 + C \frac{1}{1-\rho}.$$

501

Proof.

$$S_n \leq \rho S_{n-1} + C \leq \rho^2 S_{n-2} + C(1+\rho) \leq \rho^n S_0 + C \sum_{i=0}^{n-1} \rho^i \leq \rho^n S_0 + C \frac{1}{1-\rho}.$$

502

\square

503 **Lemma A.7.** *Let Assumption 2.1 – 2.8 hold and let H_0 be the smallest integer such that*

$$H_0 > \frac{\log(C_u/C_l)}{\log(1/\gamma)}.$$

504 Moreover, define $\nu = \frac{C_u}{C_l} \gamma^{H_0}$. Note, by definition of H_0 , $\nu < 1$. Then we have for all $k \in$
505 $\{0, \dots, H_n\}$ and $n > 0$

506 Bounded expectation over horizon

$$\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^n | \mathbf{x}_0} [V^{\boldsymbol{\pi}^n}(\mathbf{x}_k^n)] \leq \gamma^k \mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^n | \mathbf{x}_0} [V^{\boldsymbol{\pi}^n}(\mathbf{x}_0^n)] + K/(1-\gamma). \quad (12)$$

507 Bounded expectation over episodes

$$\mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^n | \mathbf{x}_0} [V^{\boldsymbol{\pi}^n}(\mathbf{x}_0^n)] \leq \nu^n V^{\boldsymbol{\pi}^0}(\mathbf{x}_0) + \frac{C_u}{C_l} K/(1-\gamma) \frac{1}{1-\nu}. \quad (13)$$

508 Moreover, we have

$$\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^n | \mathbf{x}_0} [V^{\boldsymbol{\pi}^n}(\mathbf{x}_k^n)] \leq D(\mathbf{x}_0, K, \gamma, \nu), \quad (14)$$

509 with $D(\mathbf{x}_0, K, \gamma, \nu) = V^{\boldsymbol{\pi}^0}(\mathbf{x}_0) + K/(1-\gamma) \left(\frac{C_u}{C_l} \frac{1}{1-\nu} + 1 \right)$

510 *Proof.* We start with proving the first claim

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^n}(\mathbf{x}_k^n)] &= \mathbb{E}_{\mathbf{x}_{k-1}^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} [V^{\pi^n}(\mathbf{x}_k^n)]] \\
&\leq \mathbb{E}_{\mathbf{x}_{k-1}^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\gamma V^{\pi^n}(\mathbf{x}_{k-1}^n) + K] && \text{(Assumption 2.4)} \\
&= \gamma \mathbb{E}_{\mathbf{x}_{k-1}^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^n}(\mathbf{x}_{k-1}^n)] + K
\end{aligned}$$

511 We can apply Corollary A.6 to prove the claim. For the second claim, we note that for any π , π' and
512 $\mathbf{x} \in \mathcal{X}$ we have from Assumption 2.4

$$V^\pi(\mathbf{x}) \leq C_u \alpha(\|\mathbf{x}\|) \leq \frac{C_u}{C_l} V^{\pi'}(\mathbf{x}).$$

513 Therefore,

$$\begin{aligned}
&\mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^n}(\mathbf{x}_0^n)] \\
&\leq \frac{C_u}{C_l} \mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^{n-1}}(\mathbf{x}_0^n)] \\
&= \frac{C_u}{C_l} \mathbb{E}_{\mathbf{x}_{H_n}^{n-1}, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^{n-1}}(\mathbf{x}_{H_n}^{n-1})] && \text{(Since } \mathbf{x}_0^n = \mathbf{x}_{H_n}^{n-1}\text{)} \\
&\leq \left(\frac{C_u}{C_l} \gamma^{H_n} \right) \mathbb{E}_{\mathbf{x}_0^{n-1}, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^{n-1}}(\mathbf{x}_0^{n-1})] + \frac{C_u}{C_l} K / (1 - \gamma) && \text{(Equation (12))}
\end{aligned}$$

514 For our choice of H_0 , we have for all $n \geq 0$ that $\frac{C_u}{C_l} \gamma^{H_n} \leq \frac{C_u}{C_l} \gamma^{H_0} \leq \nu < 1$. From Corollary A.6,
515 we get

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^n}(\mathbf{x}_0^n)] &\leq \left(\frac{C_u}{C_l} \gamma^{H_n} \right) \mathbb{E}_{\mathbf{x}_0^{n-1}, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^{n-1}}(\mathbf{x}_0^{n-1})] + \frac{C_u}{C_l} K / (1 - \gamma) \\
&\leq \nu \mathbb{E}_{\mathbf{x}_0^{n-1}, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^{n-1}}(\mathbf{x}_0^{n-1})] + \frac{C_u}{C_l} K / (1 - \gamma) \\
&\leq \nu^n V^{\pi^0}(\mathbf{x}_0) + \frac{C_u}{C_l} K / (1 - \gamma) \frac{1}{1 - \nu}. && \text{(Corollary A.6)}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^n}(\mathbf{x}_k^n)] &\leq \gamma^k \mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^n}(\mathbf{x}_0^n)] + K / (1 - \gamma) && \text{(Equation (12))} \\
&\leq \mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi^n}(\mathbf{x}_0^n)] + K / (1 - \gamma) \\
&\leq \nu^n V^{\pi^0}(\mathbf{x}_0) + \frac{C_u}{C_l} K / (1 - \gamma) \frac{1}{1 - \nu} + K / (1 - \gamma) && \text{(Equation (13))} \\
&\leq V^{\pi^0}(\mathbf{x}_0) + \frac{C_u}{C_l} K / (1 - \gamma) \frac{1}{1 - \nu} + K / (1 - \gamma)
\end{aligned}$$

516 □

517 **Lemma A.8.** Let Assumption 2.1 – 2.8 hold and let H_0 be the smallest integer such that

$$H_0 > \frac{\log(C_u/C_l)}{\log(1/\gamma)}.$$

518 Moreover, define $\nu = \frac{C_u}{C_l} \gamma^{H_0}$. Note, by definition of H_0 , $\nu < 1$.

519 Then we have for all $k \in \{0, \dots, H_n\}$ and $n > 0$

520 Bounded second moment over horizon

$$\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi^n}(\mathbf{x}_k^n))^2 \right] \leq \gamma^{2k} \mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi^n}(\mathbf{x}_0^n))^2 \right] + \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2} \quad (15)$$

521 with $D_2(\mathbf{x}_0, K, \gamma, \nu) = 2K\gamma D(\mathbf{x}_0, K, \gamma, \nu) + K^2 + C_w$, and $C_w = \mathbb{E}_w [\kappa^2(\|w\|)] +$
522 $3(\mathbb{E}_w [\kappa(\|w\|)])^2$.

523 Bounded second moment over episodes

$$\mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_0^n))^2 \right] \leq \nu^{2n} (V^{\pi_0}(\mathbf{x}_0))^2 + \left(\frac{C_u}{C_l} \right)^2 \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2} \frac{1}{1 - \nu^2}. \quad (16)$$

524 Moreover, let $D_3(\mathbf{x}_0, K, \gamma, \nu) = (V^{\pi_0}(\mathbf{x}_0))^2 + D_2(\mathbf{x}_0, K, \gamma, \nu) \left(\left(\frac{C_u}{C_l} \right)^2 \frac{1}{1 - \gamma^2} \frac{1}{1 - \nu^2} + \frac{1}{1 - \gamma^2} \right)$.

$$\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_k^n))^2 \right] \leq D_3(\mathbf{x}_0, K, \gamma, \nu)$$

525

526 *Proof.* Note that,

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} \left[(V^{\pi_n}(\mathbf{x}_k^n))^2 \right] &= \left(\mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} [V^{\pi_n}(\mathbf{x}_k^n)] \right)^2 \\ &\quad + \mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} \left[\left(V^{\pi_n}(\mathbf{x}_k^n) - \mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} [V^{\pi_n}(\mathbf{x}_k^n)] \right)^2 \right]. \end{aligned}$$

527 We first bound the second term. Let $\bar{\mathbf{x}}_k^n = \mathbf{f}^*(\mathbf{x}_{k-1}^n, \pi_n(\mathbf{x}_{k-1}^n))$, i.e., the next state in the absence of
528 transition noise.

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} \left[\left(V^{\pi_n}(\mathbf{x}_k^n) - \mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} [V^{\pi_n}(\mathbf{x}_k^n)] \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} \left[\left(V^{\pi_n}(\mathbf{x}_k^n) - V^{\pi_n}(\bar{\mathbf{x}}_k^n) + V^{\pi_n}(\bar{\mathbf{x}}_k^n) - \mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} [V^{\pi_n}(\mathbf{x}_k^n)] \right)^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} \left[\left(V^{\pi_n}(\mathbf{x}_k^n) - V^{\pi_n}(\bar{\mathbf{x}}_k^n) + \mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} [V^{\pi_n}(\bar{\mathbf{x}}_k^n) - V^{\pi_n}(\mathbf{x}_k^n)] \right)^2 \right] \\ &\leq \mathbb{E}_{\mathbf{w}} \left[(\kappa(\|w\|) + \mathbb{E}_{\mathbf{w}}[\kappa(\|w\|)])^2 \right] \quad (\text{uniform continuity of } V^{\pi_n}) \\ &= \mathbb{E}_{\mathbf{w}} \left[\kappa^2(\|w\|) + 3(\mathbb{E}_{\mathbf{w}}[\kappa(\|w\|)])^2 \right] \\ &= C_{\mathbf{w}} \quad (\text{Assumption 2.4}) \end{aligned}$$

529 Therefore we have

$$\begin{aligned} \mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} \left[(V^{\pi_n}(\mathbf{x}_k^n))^2 \right] &= \left(\mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} [V^{\pi_n}(\mathbf{x}_k^n)] \right)^2 + C_{\mathbf{w}} \\ &\leq (\gamma V^{\pi_n}(\mathbf{x}_k^n) + K)^2 + C_{\mathbf{w}} \\ &= \gamma^2 (V^{\pi_n}(\mathbf{x}_{k-1}^n))^2 + 2K\gamma V^{\pi_n}(\mathbf{x}_{k-1}^n) + K^2 + C_{\mathbf{w}}. \end{aligned}$$

$$\begin{aligned} &\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_k^n))^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_{k-1}^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[\mathbb{E}_{\mathbf{x}_k^n | \mathbf{x}_{k-1}^n} \left[(V^{\pi_n}(\mathbf{x}_k^n))^2 \right] \right] \\ &\leq \gamma^2 \mathbb{E}_{\mathbf{x}_{k-1}^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_{k-1}^n))^2 \right] + 2K\gamma \mathbb{E}_{\mathbf{x}_{k-1}^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [V^{\pi_n}(\mathbf{x}_{k-1}^n)] + K^2 + C_{\mathbf{w}} \\ &\leq \gamma^2 \mathbb{E}_{\mathbf{x}_{k-1}^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_{k-1}^n))^2 \right] + 2K\gamma D(\mathbf{x}_0, K, \gamma, \nu) + K^2 + C_{\mathbf{w}}. \quad (\text{Lemma A.7}) \end{aligned}$$

530 Let $D_2(\mathbf{x}_0, K, \gamma, \nu) = 2K\gamma D(\mathbf{x}_0, K, \gamma, \nu) + K^2 + C_{\mathbf{w}}$. Applying Corollary A.6 we get

$$\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_k^n))^2 \right] \leq \gamma^{2k} \mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_0^n))^2 \right] + \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2}$$

531 Similar to the first moment, we leverage that $V^{\pi_n}(\mathbf{x}) \leq \frac{C_u}{C_l} V^{\pi_{n-1}}(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}$, $\frac{C_u}{C_l} \gamma^{H_{n-1}} \leq \nu$,
532 and get,

$$\mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_0^n))^2 \right]$$

$$\begin{aligned}
&\leq \left(\frac{C_u}{C_l}\right)^2 \mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_{n-1}}(\mathbf{x}_0^n))^2 \right] \\
&= \left(\frac{C_u}{C_l}\right)^2 \mathbb{E}_{\mathbf{x}_{H_n}^{n-1}, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_{n-1}}(\mathbf{x}_{H_n}^{n-1}))^2 \right] \quad (\text{Since } \mathbf{x}_0^n = \mathbf{x}_{H_n}^{n-1}) \\
&\leq \left(\frac{C_u}{C_l} \gamma^{H_n}\right)^2 \mathbb{E}_{\mathbf{x}_0^{n-1}, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_{n-1}}(\mathbf{x}_0^{n-1}))^2 \right] + \left(\frac{C_u}{C_l}\right)^2 \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2} \quad (\text{Equation (15)}) \\
&\leq \nu^2 \mathbb{E}_{\mathbf{x}_0^{n-1}, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_{n-1}}(\mathbf{x}_0^{n-1}))^2 \right] + \left(\frac{C_u}{C_l}\right)^2 \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2} \\
&\leq \nu^{2n} (V^{\pi_0}(\mathbf{x}_0))^2 + \left(\frac{C_u}{C_l}\right)^2 \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2} \frac{1}{1 - \nu^2} \quad (\text{Corollary A.6})
\end{aligned}$$

533 Moreover,

$$\begin{aligned}
&\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_k^n))^2 \right] \\
&\leq \gamma^{2k} \mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_0^n))^2 \right] + \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2} \quad (\text{Equation (15)}) \\
&\leq \mathbb{E}_{\mathbf{x}_0^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[(V^{\pi_n}(\mathbf{x}_0^n))^2 \right] + \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2} \\
&\leq \nu^{2n} (V^{\pi_0}(\mathbf{x}_0))^2 + \left(\frac{C_u}{C_l}\right)^2 \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2} \frac{1}{1 - \nu^2} + \frac{D_2(\mathbf{x}_0, K, \gamma, \nu)}{1 - \gamma^2} \quad (\text{Equation (16)}) \\
&\leq (V^{\pi_0}(\mathbf{x}_0))^2 + D_2(\mathbf{x}_0, K, \gamma, \nu) \left(\left(\frac{C_u}{C_l}\right)^2 \frac{1}{1 - \gamma^2} \frac{1}{1 - \nu^2} + \frac{1}{1 - \gamma^2} \right)
\end{aligned}$$

534

□

535 Finally, we prove the regret bound of NEORL.

536 *Proof of Theorem 3.1.* In the following, let $\hat{\mathbf{x}}_{k+1}^n = \mathbf{f}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)) + \mathbf{w}_k^n$ denote the state predicted
537 under the optimistic dynamics and $\mathbf{x}_{k+1}^n = \mathbf{f}_n^*(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)) + \mathbf{w}_k^n$ the true state.

$$\begin{aligned}
&\mathbb{E} \left[\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} c(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)) - A(\boldsymbol{\pi}^*) \right] \\
&\leq \mathbb{E} \left[\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} c(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)) - A(\boldsymbol{\pi}_n, \mathbf{f}_n) \right] \quad (\text{Optimism}) \\
&= \mathbb{E} \left[\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_k^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \hat{\mathbf{x}}_{k+1}^n) \right] \quad (\text{Bellman equation (Equation (4))}) \\
&= \mathbb{E} \left[\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_k^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n) + B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \hat{\mathbf{x}}_{k+1}^n) \right] \\
&= \sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \hat{\mathbf{x}}_{k+1}^n)] \quad (\text{A}) \\
&\quad + \sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_k^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n)] \quad (\text{B})
\end{aligned}$$

538 First, we study the term (A).

539 **Proof for (A):** Note that because $\mathbf{f}_n \in \mathcal{M}_n$, there exists a $\boldsymbol{\eta} \in [-1, 1]^{d_x}$ such that $\hat{\mathbf{x}}_{k+1}^n =$
540 $\boldsymbol{\mu}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)) + \beta_n \boldsymbol{\sigma}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)) \boldsymbol{\eta}(\mathbf{x}_k^n) + \mathbf{w}_k^n$. Furthermore, $\mathbf{x}_{k+1}^n = \mathbf{f}_n^*(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)) + \mathbf{w}_k^n$
541 and the transition noise is Gaussian. Let $\zeta_{2,k}^n$ and $\zeta_{1,k}^n$ denote the respective distributions of the

542 two random variables, i.e., $\zeta_{1,k}^n \sim \mathcal{N}(\mathbf{f}^*(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)), \sigma^2 \mathbb{I})$ and $\zeta_{2,k}^n \sim \mathcal{N}(\mathbf{f}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)), \sigma^2 \mathbb{I})$.
 543 Next, define $\bar{B} = \mathbb{E}_{\mathbf{x} \sim \zeta_{2,k}^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x})]$, and consider the function $h(\mathbf{x}) = B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}) - \bar{B}$.
 544 Then we have

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}_k^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \hat{\mathbf{x}}_{k+1}^n)] \\ &= \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \zeta_{2,k}^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x})] \\ &= \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}) - \bar{B}] - \mathbb{E}_{\mathbf{x} \sim \zeta_{2,k}^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}) - \bar{B}] \\ &= \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \zeta_{2,k}^n} [h(\mathbf{x})]. \end{aligned}$$

545 Note that $\mathbb{E}_{\mathbf{x} \sim \zeta_{2,k}^n} [h(\mathbf{x})] = 0$ by the definition of h and thus,

$$\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \zeta_{2,k}^n} [h(\mathbf{x})] = \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h(\mathbf{x})] \leq \sqrt{\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]}. \quad (17)$$

546 In the following, we bound the term above w.r.t. the Chi-squared distance

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}_k^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \hat{\mathbf{x}}_{k+1}^n)] = \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \zeta_{2,k}^n} [h(\mathbf{x})] \\ &= \int_{\mathcal{X}} h(\mathbf{x}) \left(1 - \frac{\zeta_{2,k}^n}{\zeta_{1,k}^n}\right) \zeta_{1,k}^n(d\mathbf{x}) \leq \sqrt{\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]} \sqrt{d_{\chi}(\zeta_{2,k}^n, \zeta_{1,k}^n)} \\ & \hspace{15em} ((\text{Kakade et al., 2020, Lemma C.2.})) \end{aligned}$$

547 With $d_{\chi}(\zeta_{2,k}^n, \zeta_{1,k}^n)$ being the Chi-squared distance.

$$d_{\chi}(\zeta_{2,k}^n, \zeta_{1,k}^n) = \int_{\mathcal{X}} \frac{(\zeta_{1,k}^n - \zeta_{2,k}^n)^2}{\zeta_{1,k}^n} (d\mathbf{x})$$

548 Since both bounds from Equation (17) and bound we got by applying (Kakade et al., 2020, Lemma
 549 C.2.), we can apply minimum and have:

$$\mathbb{E}_{\mathbf{w}_k^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \hat{\mathbf{x}}_{k+1}^n)] \leq \sqrt{\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]} \sqrt{\min \{d_{\chi}(\zeta_{2,k}^n, \zeta_{1,k}^n), 1\}}$$

550 Therefore, following Kakade et al. (2020, Lemma C.2.), we get

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}_k^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \hat{\mathbf{x}}_{k+1}^n)] \\ & \leq \sqrt{\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]} \min \{1/\sigma \|\mathbf{f}^*(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)) - \mathbf{f}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n))\|, 1\} \\ & \leq \sqrt{\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]} (1 + \sqrt{d_x})^{\beta_n/\sigma} \|\boldsymbol{\sigma}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n))\|. \quad ((\text{Sukhija et al., 2024, Cor. 3})) \end{aligned}$$

551 Therefore, we have

$$\begin{aligned} & \sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\mathbb{E}_{\mathbf{w}_k^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \hat{\mathbf{x}}_{k+1}^n)]] \\ & \leq \sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[\sqrt{\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]} (1 + \sqrt{d_x})^{\beta_n/\sigma} \|\boldsymbol{\sigma}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n))\| \right] \\ & \leq \sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} (1 + \sqrt{d_x})^{\beta_n/\sigma} \sqrt{\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]] \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\|\boldsymbol{\sigma}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n))\|^2]} \\ & \leq (1 + \sqrt{d_x})^{\beta_T/\sigma} \sqrt{\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]]} \\ & \times \sqrt{\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\|\boldsymbol{\sigma}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n))\|^2]} \end{aligned}$$

552 Here, for the second and third inequality, we use Cauchy-Schwarz. Now we bound the two terms
 553 above individually.

554 First we bound $\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]$.

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})] &= \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [(B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}) - \bar{B})^2] \\
 &= \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} \left[(B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}) - \mathbb{E}_{\mathbf{x} \sim \zeta_{2,k}^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x})])^2 \right] \\
 &\leq \left(\frac{C_2}{1 - \hat{\lambda}} \right)^2 \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} \left[(2 + V^{\boldsymbol{\pi}_n}(\mathbf{x}) + \mathbb{E}_{\mathbf{x} \sim \zeta_{2,k}^n} [V^{\boldsymbol{\pi}_n}(\mathbf{x})])^2 \right] \quad (\text{Theorem A.3}) \\
 &\leq \left(\frac{C_2}{1 - \hat{\lambda}} \right)^2 \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} \left[(2 + V^{\boldsymbol{\pi}_n}(\mathbf{x}) + \gamma V^{\boldsymbol{\pi}_n}(\mathbf{x}_k^n) + \hat{K})^2 \right] \quad (\text{Lemma A.4}) \\
 &\leq \left(\frac{\sqrt{2}C_2}{1 - \hat{\lambda}} \right)^2 \mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} \left[(V^{\boldsymbol{\pi}_n}(\mathbf{x}))^2 + (2 + \gamma V^{\boldsymbol{\pi}_n}(\mathbf{x}_k^n) + \hat{K})^2 \right] \\
 &\leq \left(\frac{\sqrt{2}C_2}{1 - \hat{\lambda}} \right)^2 \left(\mathbb{E}_{\mathbf{x}_{k+1}^n | \mathbf{x}_k^n} [(V^{\boldsymbol{\pi}_n}(\mathbf{x}_{k+1}))^2] + 2\gamma^2 (V^{\boldsymbol{\pi}_n}(\mathbf{x}_k^n))^2 + 2(2 + \hat{K})^2 \right)
 \end{aligned}$$

555 Furthermore, we have from Lemma A.8.

$$\begin{aligned}
 &\mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} \left[\mathbb{E}_{\mathbf{x}_{k+1}^n | \mathbf{x}_k^n} [(V^{\boldsymbol{\pi}_n}(\mathbf{x}_{k+1}))^2] + 2\gamma^2 (V^{\boldsymbol{\pi}_n}(\mathbf{x}_k^n))^2 \right] \\
 &= \mathbb{E}_{\mathbf{x}_{k+1}^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [(V^{\boldsymbol{\pi}_n}(\mathbf{x}_{k+1}))^2] + 2\gamma^2 \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [(V^{\boldsymbol{\pi}_n}(\mathbf{x}_{k+1}))^2] \leq (1 + 2\gamma^2) D_3(\mathbf{x}_0, K, \gamma, \nu).
 \end{aligned}$$

556 In the end, we get

$$\begin{aligned}
 &\sqrt{\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]]} \\
 &\leq \left(\frac{\sqrt{2}C_2}{1 - \hat{\lambda}} \right) \sqrt{\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} (1 + 2\gamma^2) D_3(\mathbf{x}_0, K, \gamma, \nu) + 2(2 + \hat{K})^2} \\
 &= \left(\frac{\sqrt{2}C_2}{1 - \hat{\lambda}} \right) \sqrt{(1 + 2\gamma^2) D_3(\mathbf{x}_0, K, \gamma, \nu) + 2(2 + \hat{K})^2} \sqrt{\sum_{n=0}^{N-1} H_n} \\
 &= \left(\frac{\sqrt{2}C_2}{1 - \hat{\lambda}} \right) \sqrt{(1 + 2\gamma^2) D_3(\mathbf{x}_0, K, \gamma, \nu) + 2(2 + \hat{K})^2} \sqrt{T}.
 \end{aligned}$$

557 Next, we use the bound from Curi et al. (2020, Lemma 17.) for the second term.

$$\sqrt{\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\|\boldsymbol{\sigma}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n))\|^2]} \leq C' \sqrt{\Gamma_T}$$

558 Here Γ_T is the maximum information gain.

559 If we set $D_4(\mathbf{x}_0, K, \gamma) = \frac{C'(1+\sqrt{d_x})}{\sigma} \left(\frac{\sqrt{2}C_2}{1-\hat{\lambda}} \right) \sqrt{(1 + 2\gamma^2) D_3(\mathbf{x}_0, K, \gamma, \nu) + 2(2 + \hat{K})^2}$, we have

$$\begin{aligned}
 &\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\mathbb{E}_{\mathbf{w}_k^n} [B(\boldsymbol{\pi}_n, \mathbf{f}_n, \mathbf{x}_{k+1}^n) - B(\boldsymbol{\pi}_n, \mathbf{f}_n, \hat{\mathbf{x}}_{k+1}^n)]] \\
 &\leq (1 + \sqrt{d_x})^{\beta_T} / \sigma \sqrt{\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^0 | \mathbf{x}_0} [\mathbb{E}_{\mathbf{x} \sim \zeta_{1,k}^n} [h^2(\mathbf{x})]]}
 \end{aligned}$$

$$\begin{aligned}
& \times \sqrt{\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E}_{\mathbf{x}_k^n, \dots, \mathbf{x}_1^n | \mathbf{x}_0} \left[\|\boldsymbol{\sigma}_n(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n))\|^2 \right]} \\
& \leq (1 + \sqrt{d_x})^{\beta_T/\sigma} \left(\frac{\sqrt{2}C_2}{1 - \hat{\lambda}} \right) \sqrt{(1 + 2\gamma^2)D_3(\mathbf{x}_0, K, \gamma, \nu) + 2(2 + \hat{K})^2 \sqrt{T}C' \sqrt{\Gamma_T}} \\
& \leq D_4(\mathbf{x}_0, K, \gamma) \beta_T \sqrt{T\Gamma_T}
\end{aligned}$$

560 **Proof for (B):**

$$\begin{aligned}
& \sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} \mathbb{E} [B(\boldsymbol{\pi}, \mathbf{f}_n, \mathbf{x}_k^n) - B(\boldsymbol{\pi}, \mathbf{f}_n, \mathbf{x}_{k+1}^n)] = \sum_{n=0}^{N-1} \mathbb{E} [B(\boldsymbol{\pi}, \mathbf{f}_n, \mathbf{x}_0^n) - B(\boldsymbol{\pi}, \mathbf{f}_n, \mathbf{x}_{H_n}^n)] \\
& \leq \frac{C_2}{1 - \hat{\lambda}} \sum_{n=0}^{N-1} (2 + \mathbb{E} [V^\pi(\mathbf{x}_0^n) + V^\pi(\mathbf{x}_{H_n}^n)]) \quad (\text{Theorem A.3}) \\
& \leq \frac{2C_2}{1 - \hat{\lambda}} \sum_{n=0}^{N-1} (1 + D(\mathbf{x}_0, K, \gamma)) \quad (\text{Lemma A.7}) \\
& = \frac{2C_2}{1 - \hat{\lambda}} (1 + D(\mathbf{x}_0, K, \gamma)) N \\
& = D_5(\mathbf{x}_0, K, \gamma) N.
\end{aligned}$$

561 Here $D_5(\mathbf{x}_0, K, \gamma) = \frac{2C_2}{1 - \hat{\lambda}} (1 + D(\mathbf{x}_0, K, \gamma))$. Finally, for our choice, $H_n = H_0 2^n$, we get

$$\sum_{n=0}^{N-1} H_n = H_0 \sum_{n=0}^{N-1} 2^n = H_0 (2^N - 1) = T.$$

562 Therefore, $N = \log_2 \left(\frac{T}{H_0} + 1 \right)$. To this end, we get for our regret

$$\begin{aligned}
R_T &= \mathbb{E} \left[\sum_{n=0}^{N-1} \sum_{k=0}^{H_n-1} c(\mathbf{x}_k^n, \boldsymbol{\pi}_n(\mathbf{x}_k^n)) - A(\boldsymbol{\pi}^*) \right] \\
&\leq D_4(\mathbf{x}_0, K, \gamma) \beta_T \sqrt{T\Gamma_T} + D_5(\mathbf{x}_0, K, \gamma) N \\
&\leq D_4(\mathbf{x}_0, K, \gamma) \beta_T \sqrt{T\Gamma_T} + D_5(\mathbf{x}_0, K, \gamma) \log_2 \left(\frac{T}{H_0} + 1 \right)
\end{aligned}$$

563

□

564 This regret is sublinear for a very rich class of functions. We summarize bounds on Γ_T from
565 [Vakili et al. \(2021\)](#) in Table 1. Furthermore, note that $D_4(\mathbf{x}_0, K, \gamma) \in (0, \infty)$ for all $\mathbf{x}_0 \in \mathcal{X}$ with
566 $\|\mathbf{x}_0\| < \infty$, $K < \infty$, $\gamma \in (0, 1)$. The same holds for $D_5(\mathbf{x}_0, K, \gamma)$. Moreover, since $V^\pi(\mathbf{x})$ is
567 $\Theta(\zeta(\|\mathbf{x}\|))$, both D_4 and D_5 are $\Theta(\zeta(\|\mathbf{x}_0\|))$.

Table 1: Maximum information gain bounds for common choice of kernels.

Kernel	$k(\mathbf{x}, \mathbf{x}')$	Γ_T
Linear	$\mathbf{x}^\top \mathbf{x}'$	$\mathcal{O}(d \log(T))$
RBF	$e^{-\frac{\ \mathbf{x} - \mathbf{x}'\ ^2}{2l^2}}$	$\mathcal{O}(\log^{d+1}(T))$
Matèrn	$\frac{1}{\Gamma(\nu)2^{\nu-1}} \left(\frac{\sqrt{2\nu}\ \mathbf{x} - \mathbf{x}'\ }{l} \right)^\nu B_\nu \left(\frac{\sqrt{2\nu}\ \mathbf{x} - \mathbf{x}'\ }{l} \right)$	$\mathcal{O}\left(T^{\frac{d}{2\nu+d}} \log^{\frac{2\nu}{2\nu+d}}(T)\right)$

568 **A.4 Relaxing Assumption 2.4**

569 Our analysis assumes that Π consists only of policies with bounded energy. This assumption ensures
 570 that during the exploration our system remains stable. The average cost and stability are intertwined
 571 for the LQG case (Anderson & Moore, 2007). Moreover, a bounded average cost of a linear controller
 572 $\pi(\mathbf{x}) = \mathbf{K}\mathbf{x}$ implies stability and vice-versa. This is not necessarily the case for nonlinear systems,
 573 i.e., stability implies a bounded average cost (c.f., Theorem 2.6) but not vice versa. An approach is to
 574 assume this link exists also for the nonlinear case.

575 **Definition A.9** (Stable Policies). We call $\Pi_S(\mathbf{f})$ the set of stable policies for the dynamics \mathbf{f} if
 576 there exists positive constants C_u, C_l with $C_u > C_l$, $\zeta, \kappa \in \mathcal{K}_\infty$, $\gamma \in (0, 1)$ s.t., we have for all
 577 $\pi \in \Pi_S(\mathbf{f})$,

578 Bounded energy; There exists a Lyapunov function $V^\pi : \mathcal{X} \rightarrow [0, \infty)$, $K(\pi, \mathbf{f}) < \infty$ for which

$$\begin{aligned} |V^\pi(\mathbf{x}) - V^\pi(\mathbf{x}')| &\leq \kappa(\|\mathbf{x} - \mathbf{x}'\|) && \text{(uniform continuity)} \\ C_l \xi(\|\mathbf{x}\|) &\leq V^\pi(\mathbf{x}) \leq C_u \xi(\|\mathbf{x}\|) && \text{(positive definiteness)} \\ \mathbb{E}_{\mathbf{x}'|\mathbf{f}, \pi, \mathbf{x}}[V^\pi(\mathbf{x}')] &\leq \gamma V^\pi(\mathbf{x}) + K(\pi, \mathbf{f}) && \text{(drift condition)} \end{aligned}$$

579 Bounded norm of cost;

$$\sup_{\mathbf{x} \in \mathcal{X}} \frac{c(\mathbf{x}, \pi(\mathbf{x}))}{1 + V^\pi(\mathbf{x})} < \infty$$

580 Boundedness of noise with respect to κ

$$\mathbb{E}_{\mathbf{w}}[\kappa(\|\mathbf{w}\|)] < \infty, \mathbb{E}_{\mathbf{w}}[\kappa^2(\|\mathbf{w}\|)] < \infty$$

581 **Assumption A.10** (Bounded average cost implies stability). Consider any dynamics \mathbf{f} , let $\Pi_A(\mathbf{f})$ be
 582 the set of policies with bounded average cost for \mathbf{f} , i.e.,

$$\Pi_A(\mathbf{f}) = \{\pi \in \Pi \mid A(\pi, \mathbf{f}) < \infty\}. \quad (18)$$

583 We assume $\forall n \geq 0$, $\mathbf{f} \in \mathcal{M}_0 \cap \mathcal{M}_n$ that all policies $\pi \in \Pi_A(\mathbf{f})$ are stable, i.e., $\pi \in \Pi_S(\mathbf{f})$.

584 With Assumption A.10 we link the average cost criterion to the stability of our system. A natural
 585 consequence of this link is the following corollary.

586 **Corollary A.11.** *Let Assumption A.10 hold. Then the following two statements are equivalent for all*
 587 *$n \geq 0$, $\mathbf{f} \in \mathcal{M}_0 \cap \mathcal{M}_n$, and $\pi \in \Pi$.*

- 588 1. $\pi \in \Pi_A(\mathbf{f})$
- 589 2. $\pi \in \Pi_S(\mathbf{f})$.

590 *Proof.* 1 \implies 2 follows from Assumption A.10 and 2 \implies 1 from Theorem 2.6. \square

591 **Assumption A.12** (Existence of a stable policy). We assume $\Pi_S(\mathbf{f}^*) \neq \emptyset$.

592 Assumption A.12 assumes that there is at least one stable policy in Π . This is in contrast to
 593 Assumption 2.4, which assumes that all policies in Π are stable. We can relax this requirement
 594 because of Assumption A.10.

595 In the following, we show that $\pi_n \in \Pi_S(\mathbf{f}_n)$ and that this implies $\pi_n \in \Pi_S(\mathbf{f}^*)$. In summary, when
 596 doing optimistic planning, we inherently pick stable policies for the true system.

597 **Lemma A.13.** *Let Assumption 2.1 – 2.2, 2.8, A.10, and Assumption A.12 hold. Let π_n, \mathbf{f}_n denote*
 598 *the solution to Equation (6). Then we have with probability at least $1 - \delta$, $\pi_n \in \Pi_S(\mathbf{f}^*)$.*

599 *Proof.* Since $\Pi_S(\mathbf{f}^*)$ is nonempty, from Corollary A.11, we must have a policy $\pi \in \Pi_A(\mathbf{f}^*)$, and
 600 thus $A(\pi) < \infty$. This implies that $A(\pi^*) < \infty$. Since, Equation (6) is an optimistic estimate of
 601 $A(\pi^*)$, we have $A(\pi_n, \mathbf{f}_n) \leq A(\pi^*) < \infty$. Thus, $\pi_n \in \Pi_A(\mathbf{f}_n)$. Again from Corollary A.11, we
 602 have $\pi_n \in \Pi_S(\mathbf{f}_n)$ and there exists a Lyapunov function V^{π_n} and $K(\pi_n, \mathbf{f}_n)$ such that

$$\mathbb{E}_{\mathbf{x}'|\mathbf{f}_n, \pi_n, \mathbf{x}}[V^{\pi_n}(\mathbf{x}')] \leq \gamma V^{\pi_n}(\mathbf{x}) + K(\pi_n, \mathbf{f}_n)$$

603 Furthermore, due to the uniform continuity of V^{π_n} we have

$$\mathbb{E}_{\mathbf{x}'|\mathbf{x}, \mathbf{f}^*, \pi_n}[V^{\pi_n}(\mathbf{x}')] \leq \mathbb{E}_{\mathbf{x}'|\mathbf{x}, \mathbf{f}_n, \pi_n}[V^{\pi_n}(\mathbf{x}')] + \kappa \left((1 + \sqrt{d_x}) \beta_0 \sqrt{d_x} \sigma_{\max} \right) \quad \text{(c.f., Lemma A.4)}$$

$$\leq \gamma V^{\pi_n}(\mathbf{x}) + \kappa \left((1 + \sqrt{d_x}) \beta_0 \sqrt{d_x} \sigma_{\max} \right) + K(\pi_n, \mathbf{f}_n)$$

604 In summary, we have $\pi_n \in \Pi_S(\mathbf{f}^*)$ with $K(\pi_n, \mathbf{f}^*) = \kappa \left((1 + \sqrt{d_x}) \beta_0 \sqrt{d_x} \sigma_{\max} \right) + K(\pi_n, \mathbf{f}_n)$.
 605 □

606 Lemma A.13 shows that Equation (6) returns policies that are stable for the true system and therefore
 607 with probability at least $1 - \delta$ is optimizing over $\Pi_S(\mathbf{f}^*)$. Thus, even in cases where Π has policies that
 608 do not satisfy Assumption 2.4, these policies are not considered by NEORL. NEORL automatically
 609 optimizes over $\Pi_S(\mathbf{f}^*)$ and the rest of the guarantees follow with $K = \max_{\pi \in \Pi_S(\mathbf{f}^*)} K(\pi, \mathbf{f}^*)$.

610 **B Experimental Details**

611 In the following, we provide all hyperparameters used in our experiments in Table 2 and the cost
 612 function for the environments in Table 3. For NEORL, we use $\beta_n = 2$ for all the experiments, except
 for the Swimmer and the SoftArm environment where we use $\beta_n = 1$.

Table 2: Hyperparameters for results in Section 4.

Environment	iCEM parameters					Model training parameters					H	Action Repeat
	Number of samples	Number of elites	Optimizer steps	Horizon	Particles	Number of ensembles	Network architecture	Learning rate	Batch size	Number of epochs		
Pendulum-GP	500	50	10	20	5	-	-	0.01	64	-	10	1
Pendulum	500	50	10	20	5	10	256×2	0.001	64	50	10	1
MountainCar	1000	100	5	50	5	10	256×2	0.001	64	50	10	2
Reacher	1000	100	10	50	5	10	256×2	0.001	64	50	10	2
CartPole	1000	100	10	50	5	10	256×2	0.001	64	50	10	2
Swimmer	500	50	10	30	5	10	256×4	0.00005	64	100	200	4
SoftArm	500	50	10	20	5	10	256×4	0.00005	64	50	20	1
RaceCar	1000	100	10	50	5	10	256×2	0.001	64	50	10	1

613

Table 3: Cost function for the environments presented in Section 4.

Environment	Cost $c(\mathbf{x}_t, \mathbf{u}_t)$
Pendulum	$\theta_t^2 + 0.1\dot{\theta}_t + 0.1u_t^2$
MountainCar	$0.1u_t^2 - 100(1\{\mathbf{x}_t \in \mathbf{x}_{\text{goal}}\})$
Reacher	$\ \mathbf{x}_t - \mathbf{x}_{\text{target}}\ + 0.1\ u_t\ $
CartPole	$\ \mathbf{x}_t^{\text{pos}} - \mathbf{x}_{\text{target}}^{\text{pos}}\ ^2 + 10(\cos(\theta_t) - 1)^2 + 0.2\ u_t\ ^2$
Swimmer	$\ \mathbf{x}_t - \mathbf{x}_{\text{target}}\ $
SoftArm	$\ \mathbf{x}_t - \mathbf{x}_{\text{target}}\ $
RaceCar	$\ \mathbf{x}_t - \mathbf{x}_{\text{target}}\ $