## TOWARDS RELIABLE OFFLINE REINFORCEMENT LEARNING VIA LYAPUNOV UNCERTAINTY CONTROL

Anonymous authors

Paper under double-blind review

## ABSTRACT

Learning trustworthy and reliable offline policies presents significant challenges due to the inherent uncertainty in pre-collected datasets. In this paper, we propose a novel offline reinforcement learning method to tackle this issue. Inspired by the concepts of Lyapunov stability and control-invariant sets from control theory, the central idea is to introduce a restricted state space for the agent to operate within. This approach allows the learned models to exhibit reduced Bellman uncertainty and make reliable decisions. To achieve this, we regulate the expected Bellman uncertainty associated with the new policy, ensuring that its growth trend in subsequent states remains within acceptable limits. The resulting method, termed Lyapunov Uncertainty Control (LUC), is shown to guarantee that the agent remains within a low-uncertainty state enclosure throughout its entire trajectory. Furthermore, we perform extensive theoretical and experimental analysis to showcase the effectiveness and feasibility of the proposed LUC.

023

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

## 1 INTRODUCTION

025 026

Offline reinforcement learning (RL) allows policy learning from historical data without real-world interaction. However, ensuring reliable sequential decision-making from offline data poses a significant challenge in practical applications. For example, in healthcare (Tang & Wiens, 2021), a reliable diagnostic agent requires avoiding unfamiliar approaches that may introduce errors in subsequent procedures. Similar requirements exist in fields such as autonomous driving (Kiran et al., 2022), robotics control (Lobbezoo et al., 2021), and others.

033 The reliability of offline RL is undermined by the risk of deviating from the scoped safe regions, 034 i.e., stable safe control (Kang et al., 2022) from offline data. More precisely, our aim is to stop the offline-learned agent from entering areas that could cause severe consequences after deployment. 035 Specifically, in practical applications, the safety requirements for decision making are extremely 036 strict (Jiang et al., 2023), demanding that every decision made by the agent at each step be safe. 037 Meanwhile, current methods like the pessimistic and DRRL methods fall short in handling this issue. For instance, Pessimistic methods such as MOPO (Yu et al., 2020), PBRL (Bai et al., 2022) and RORL (Yang et al., 2022) mainly concentrate on making the agent act in accordance with the 040 demonstrations in the dataset. Although these methods are good at quantifying the OOD data, they 041 may not entirely fulfill the previously mentioned safety requirements. On the other hand, Robust 042 RL methods (Panaganti et al., 2022; Shi & Chi, 2024) aim to improve the agent's capacity to deal 043 with distributional shift by establishing the uncertainty set of the transition function and optimizing 044 the lower bound of the policy's long term returns under the worst case scenario within this set. 045 Unfortunately, Robust RL also fails to take into account the safety issue described earlier. In other words, they are unable to prevent the learned agent from straying from the safe region. 046

Inspired by control-invariant sets in control theory (Kerrigan, 2000; Richter & Roy, 2017), where a closed region is delineated in the state space to offer a reliable working environment for the agent, a recent method named Lyapunov Density Model (Kang et al., 2022) defines a region based on data density distribution to ensure adequate data support for the agent during operation. However, setting a reliable region based on a common and pessimistic density criterion overlooks the issue of performance imbalances in complex environments, where data requirements for achieving a certain performance level may vary across regions. For example, in autonomous driving, data needs for learning on a smooth highway differ significantly from those on a rugged mountain road. Instead,



Figure 1: (left) Solely constraining current step's uncertainty is insufficient to identify those boundary states that pose a high risk of the agent deviating into high-uncertainty regions. (right) An
illustration of the failure of traditional pessimism-based agents in accumulating errors by deviating
from reliable regions, as demonstrated on a Halfcheetah robotic agent.

we advocate using a metric linked to model performance (e.g., value functions, policies) - Bellman
 uncertainty - for region definition to meet reliability standards in complex environments.

071 In particular, we aim to define a confined state space for the agent's operation, where the learned 072 models demonstrate reduced Bellman uncertainty and reliable decision-making. To achieve this, 073 we introduce a novel offline reinforcement learning method that regulates the expected Bellman 074 uncertainty associated with the new policy. This regulation ensures that, at subsequent states, the 075 growth trend of uncertainty remains within acceptable limits, allowing the agent to navigate low-076 uncertainty regions that serve as safe zones. Drawing inspiration from the control Lyapunov func-077 tions used in optimal control, we refer to our approach as Lyapunov Uncertainty Control (LUC). 078 We implement our Lyapunov Uncertainty Control (LUC) method using a standard deviation-based 079 uncertainty measure that relies on Q-ensembles, as described in Bai et al. (2022). Theoretically, we demonstrate that our approach can confine the learned agent to operate within a low-uncertainty state enclosure, resulting in secure and reliable trajectories. Furthermore, in certain scenarios, it 081 can enhance the minimum performance bound of the new policy. Finally, we conduct extensive experiments to showcase the effectiveness and feasibility of LUC across various tailored benchmarks, 083 including out-of-distribution (OOD) benchmarks and those with poor demonstrations. 084

The subsequent sections, after a brief review of related works, Section 3 present a concise overview
 of the fundamental concepts in offline RL. Section 4 elaborates on the LUC methodology, providing
 a detailed theoretical analysis of its effectiveness and implementation insights. Section 5 presents
 experimental results evaluating LUC's performance across various scenarios. To conclude, the paper
 summarizes the findings and contributions, along with a brief discussion on limitations.

090 091

092

068

## 2 RELATED WORKS

Offline RL algorithms. Dealing with *distributional shift* poses a significant challenge for offline RL algorithms. Previous works, including CQL (Kumar et al., 2020), BEAR (Kumar et al., 2019), and BRAC (Wu et al., 2019), have aimed to tackle this issue by integrating conservative principles to prevent out-of-distribution (OOD) actions. However, these methods mainly concentrate on reducing the disparity between the new policy and the behavior policy that gathered the dataset. On the other hand, approaches like Implicit Q-Learning (IQL) (Kostrikov et al., 2022) entirely avoid OOD queries for actions during TD target estimation. Nonetheless, these methods heavily depend on the action distribution of the behavior policy, leading to a lack of generalization capability.

101

Pessimistic offline RL. Pessimistic offline RL algorithms introduce Bellman uncertainty quantification (Jin et al., 2021; Xie et al., 2021) to determine reliable actions for generalizing to OOD data. This method has exhibited significant efficacy in model-based offline RL algorithms such as Model-based Offline Policy Optimization (Yu et al., 2020) and MOdel-Bellman Inconsistency penalized offLinE Policy Optimization (MOBILE)(Sun et al., 2023), as well as in model-free algorithms like Pessimistic Bootstrapping for offline RL (PBRL)Bai et al. (2022), Ensemble-Diversified Actor Critic (EDAC)An et al. (2021) and Robust Offline RL (RORL)(Yang et al., 2022). However, this

study reveals that solely managing uncertainty at the current step is insufficient to ensure reliability
 and safety, given the short-sightedness regarding the decision's potential outcomes.

111 **Consequence-driven offline RL.** Consequence-driven offline RL methods have been developed 112 to address the state distributional shift issue utilizing the concept of state recovery(Zhang et al., 113 2022; Jiang et al., 2023). State Deviation Correction(Zhang et al., 2022) involves pre-training a forward dynamics model to facilitate the recovery process, whereas Out-of-sample State Recovery 114 (OSR) (Jiang et al., 2023) employs an inverse dynamics model to implicitly execute the concept. 115 While these techniques assist in rectifying the agent's behavior deviations from reliable regions, 116 their unreliability in OOD states may hinder their success rates. In contrast, our LUC approach can 117 be thought of as a more robust method which not only considers the reliability of the immediate 118 consequence of executing a policy but that of its long term consequence. 119

120

121 122

123

124

125

126

127

**Robust Reinforcement Learning.** Robust RL's pessimism stems from penalizing with uncertainty in the outcome predictions of actions, to deal with the problem of distributional shift(Panaganti et al., 2022; Shi & Chi, 2024). However, if some behaviors that could lead to the deviation from the safe regions are supported well by dataset (as is in the case described in Figure 1), the penalty loses its effectiveness due to that the uncertainty in outcome predictions would be tiny, thereby exacerbating the risk of entering high-uncertainty regions. Therefore, we conclude that Robust RL are unable to prevent the learned agent from straying from the safe region, i.e., the safety requirements.

## 128 129

130 131 3

PRELIMINARIES

In the standard formulation of reinforcement learning, a Markov Decision Process (MDP) is used to 132 model the problem. It is represented by a tuple  $(S, A, P, R, \gamma, \rho_0)$ , where S denotes the state space, 133 A represents the action space, M is the transition function (in a deterministic transitioned MDP, 134 M(s, a) = s', while in a stochastic transitioned MDP, P(s'|s, a) is a distribution of states), R is the 135 reward function with upper bound  $R_{max}$ ,  $\gamma$  is the discount factor, and  $\rho_0$  is the initial state distribu-136 tion. A policy, denoted as  $\pi: S \to A$ , guides the decision-making process in interacting with the 137 environment. To evaluate the expected cumulative rewards, a Q-value function  $Q^{\pi}(s, a)$  is defined 138 as  $(1 - \gamma)\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(a_t|s_t))|s, a]$ . For convenience, the  $\gamma$ -discounted future state distribu-139 tion (stationary state distribution) is defined as  $d^{\pi}(s) = (1-\gamma) \sum t = 0^{\infty} \gamma^t Pr(s_t = s; \pi, \rho_0)$ , with  $\rho_0$  representing the initial state distribution and  $(1 - \gamma)$  is the normalization factor. 140

In the offline setting, Q-Learning (Watkins & Dayan, 1992) is used to learn a Q-value function  $\hat{Q}(s, a)$  and a policy  $\pi$  from a dataset  $\mathcal{D}$  collected by a behavior policy  $\pi_{\beta}$ . The dataset consists of quadruples  $(s, a, r, s') \sim d^{\pi_{\beta}}(s)\pi_{\beta}(a|s)P(r|s, a)P(s'|s, a)$ . The objective is to minimize the Bellman error over the offline dataset (Watkins & Dayan, 1992) and search for a good policy in the policy candidate set  $\Pi \subset (S \to \Delta(A))$  under the supervision of a value-function class  $\mathcal{F} \subset$  $(S \times A \to [0, V_{max}])$  to model the Q-value function,

147

 $Q \longleftarrow \arg\min_{Q} \mathbb{E}_{(s,a,r,s')\sim\mathcal{D}} \left[ r + \gamma \left[ \max_{\pi \in \Pi} \mathbb{E}_{a' \sim \pi(\cdot|s')} Q(s',a') \right] - Q(s,a) \right]^2 \tag{1}$ 

More specifically, in this paper, we denote the optimal Bellman operator over  $\Pi$  as  $\mathcal{T}^{\Pi}f(s,a) = r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)}[\max_{\pi \in \Pi} \mathbb{E}_{a' \sim \pi(\cdot|s')}Q(s',a')]$ , and the empirical Bellman operator as  $\hat{\mathcal{T}}^{\Pi}f(s,a) = r(s,a) + \gamma \mathbb{E}_{s' \sim \hat{P}(\cdot|s,a)}[\max_{\pi \in \Pi} \mathbb{E}_{a' \sim \pi(\cdot|s')}Q(s',a')]$ , where  $\hat{P}$  is the empirical dynamics model based on the dataset. It is worth noting that the TD target in Eq.(1) is estimated by the empirical Bellman operator.

156 4 METHOD

157

155

This section provides a detailed introduction to our work. In Sec.4.1, we formally define our objective mathematically as obtaining a policy that consistently operates within a reliable region. Subsequently, in Sec.4.2, we present a specific algorithm to accomplish this objective. Finally, in Sec.4.3, we analyze the theoretical properties of the algorithm, demonstrating its capability to improve the performance lower bound of the learned policy under specific scenarios.

## 162 4.1 OPERATING WITHIN RELIABLE REGIONS BY LYAPUNOV POLICY

164 In this section, we formally define our conceptual framework for Lyapunov Uncertainty Control. 165 Specifically, analogous to control-invariant sets in control theory, we first define reliable regions in the state space where the policy  $\pi$  can operate effectively. Previous methods defined these regions 166 using density models (Kang et al., 2022); however, as mentioned earlier, a local region may not 167 be reliable even with high density due to the complexity and nonlinearity of the underlying envi-168 ronment. Instead, we introduce a measurement based on epistemic uncertainty, denoted as  $\zeta_f(s, a)$ . The precise computational method for  $\zeta_f(s, a)$  will be presented later; however, it is a positive scalar 170 that is closely related to the agent's current knowledge, reflecting the generalization capability of the 171 learned value function f at the input data (s, a). If f generalizes well,  $\zeta_f(s, a)$  will be small; con-172 versely, if f does not generalize well,  $\zeta_f(s, a)$  will be large. 173

Using  $\zeta_f(s, a)$ , we can evaluate whether an induced policy  $\pi$  can make reliable decisions at a given state s by assessing the uncertainty of the learned value function f. Furthermore, we derive the following definition of the  $f - \pi$  reliable region in Definition 1:

**Definition 1.**  $(f - \pi \text{ reliable region})$ . Given an arbitrary value function f, policy  $\pi$  and a threshold c, we define the  $f - \pi$  reliable region over the state space,

195

$$\mathcal{G}_f(\pi) = \{s | \zeta_f(s,\pi) \le c\}$$
(2)

where 
$$\zeta_f(s,\pi) = \mathbb{E}_{a \sim \pi(\cdot|s)} \zeta_f(s,a).$$

As shown in Figure 1, where 'low-uncertainty region' illustrates the agent's  $f - \pi$  reliable regions. If the agent operates beyond its reliable regions, it would accumulate decision errors, finally failing the task. On the other hand, if an agent with policy  $\pi$  always operate within its  $f - \pi$  reliable regions, we call this policy a reliable policy, defined in Definition 2.

**Definition 2.** (*Reliable policy*). A policy  $\pi$  is reliable if it satisfies that  $\forall s \in \mathcal{D} \cap \mathcal{G}_f(\pi)$ , if  $\forall t, s_t \in supp(P(s_t|s_0 = s, \pi))$ , then  $s_t \in \mathcal{G}_f(\pi)$ .

Furthermore, a Lyapunov policy, as defined in Definition 3, not only manages uncertainty at the current step but also addresses the tendencies of these uncertainty one step ahead. In other words, a Lyapunov policy is capable of controlling current step uncertainty to encompass reliable regions over the state space, while also restricting one-step forward uncertainty to ensure trajectory reliability.

**Definition 3.** (*Lyapunov policy*). Given an arbitrary value function f and offline dataset  $\mathcal{D}$ . A policy  $\pi$  is a Lyapunov policy if it satisfies

$$1.\forall s \in \mathcal{D}, \zeta_f(s,\pi) \le c; \qquad 2.\forall s \in \mathcal{D}, \max_{\hat{a} \in \pi} \zeta_f(M(s,\hat{a}),\pi) \le \zeta_f(s,\pi).$$
(3)

where *M* is the deterministic transition.

Then we have the following results, shown in Theorem 1,

**Theorem 1.** In a deterministic transitioned MDP, a Lyapunov policy  $\pi$  is a reliable policy.

201 The proof of Theorem 1 is given in Appendix A.1. Essentially, this theorem says that a Lyapunov policy is also a Lyapunov reliable policy. In other words, if a policy is a Lyapunov policy, then it will operate within its enclosed reliable region.

**Proposition 1.** (*Existence of reliable policy.*) Suppose the dataset have a sufficient coverage over the optimal policy, i.e.,  $\sup_{s,a} \frac{\pi^*(a|s)}{\pi_{\beta}(a|s)} \leq C^*$ . Then there exists a reliable policy.

207 Proof of Proposition 1 could be seen in Appendix A.1. Proposition 1 shows that there would always
208 exist a reliable policy in the MDP system with sufficient data, which motivates us to learn such a
209 policy for reliable control.

To summary, the main theoretical results in this section are used for problem formulation and func tional requirements of the method - Definition 1 defines a safe region for the agent to operate, and
 Definition 2 defines reliable policy based on it, which is able to verify the safety requirements, i.e.,
 operate within the region defined by Definition 1; Definition 3 and Theorem 1 indicate what kind of
 policies need to be learned to ensure that the agent can operate stably within the safe region with out exceeding it. Proposition 1 demonstrates the existence of policies that meet these functional
 requirements, validating the feasibility of the method proposed in this paper.

## 4.2 IMPLEMENTING LYAPUNOV UNCERTAINTY CONTROL BY VALUE ESTIMATION

In this section, we use the Bellman uncertainty quantifier as in (Jin et al., 2021) to implement the value-epistemic uncertainty in Definition 1,

$$\zeta_f(s,a) = \|\mathcal{T}f - \hat{\mathcal{T}}f\|(s,a) \tag{4}$$

where f is the learned value function.  $\mathcal{T}$  is an arbitrary Bellman operator, while  $\hat{\mathcal{T}}$  is its empirical version according to the dataset. Previous studies (An et al., 2021) suggest that Bellman uncertainty can rely on model predictions to evaluate state-action pairs. High variance in the model's prediction for a particular action implies inadequate data support, leading to low reliability. This property confirms that Bellman uncertainty aligns with the requirement in Definition 1.

Next, our objective is to acquire the Lyapunov reliable policy outlined in Definition 3 from the offline dataset using a model-free approach. Here, we present the Lyapunov value estimation, which straightforwardly penalizes not only the Q-value functions using the uncertainty quantifier from Eq.(4) at the current time step but also the increasing decision uncertainty tendency based on the next time step's situation, as,

$$\mathcal{L}_{LUC}(s, a, s', a', f) = \mathbb{E}_{a \sim \pi(\hat{a}|s)} \zeta_f(s, \hat{a}) + \beta \cdot (\zeta_f(s', a') - \zeta_f(s, a))$$
(5)

Subsequently, we apply regularization using the offline dataset D and the learned models (Q function  $f_k$  and new policy  $\pi$ ), incorporating it as a penalty in the value estimation,

$$f_k(s,a) \leftarrow f_k(s,a) - \hat{\beta} \cdot \mathbb{E}_{a' \sim \pi(\cdot|s')} \mathcal{L}_{LUC}(s,a,s',a',f_k) \tag{6}$$

where  $f_k$  is the Q function learned at the  $k^{th}$  iteration.  $\hat{\beta}$  is the balance coefficient. (s, a, s') is the tuple sampled from the offline dataset  $\mathcal{D}$ , and the  $\pi$  is the currently learned policy.

**Proposition 2.** Suppose the action distribution of new policy  $\pi(a|s)$  is positive correlated to the learnt *Q* function f(s, a), i.e.,  $\pi(a|s) \propto f(s, a)$ . Then the proposed Lyapunov value estimation in Eq.(5) induces a Lyapunov reliable policy as defined in Definition 3.

The proof is available in Appendix A.1. Proposition 2 demonstrates that the policy induced by the value function trained with Eq.(6) may exhibit the traits of the Lyapunov reliable policy described in Definition 3, fulfilling the reliability criteria in our study.

Then like previous pessimistic methods (An et al., 2021; Bai et al., 2022; Yang et al., 2022), we approximate the uncertainty quantifier in Eq.(4) as the standard deviation as,

253 254 255

256

257 258 259

260

261 262 263

220 221 222

233 234 235

236

237 238

 $\Gamma_f(s,a) \approx \beta \cdot Std(f^i(s,a)) = \beta \cdot \sqrt{\frac{1}{K} \sum_{i=1}^K \left( f^i(s,a) - \bar{f}(s,a) \right)}$ (7)

where  $\{f^i\}_{i=1}^K$  is the learned Q-ensembles and  $\overline{f}$  is the average of the K Q-esembles, and  $\beta$  is the balance-coefficient. Then the objective in Eq.(5) is converted to,

$$\mathcal{L}_{LUC}(s, a, s', a', f) = (1 - \beta) \cdot Std(f^i(s, a)) + \beta \cdot \gamma \cdot Std(f^i(s', a'))$$
(8)

where  $U_f(s, a) = std(f^i(s, a))$  and  $\{f^i\}_{i=1}^K$  is the learned K Q-ensembles. In practice, the  $\beta$  is usually selected in (0, 1). Then the regularization of LUC is,

$$\mathcal{L}_{LUC}(f^{i},\pi) = \mathbb{E}_{(\hat{s},\hat{a},\hat{s}'\sim\hat{\mathcal{D}})}(f^{i}(\hat{s},\hat{a}) - \hat{\beta} \cdot \mathbb{E}_{\hat{a}'\sim\pi(\cdot|\hat{s}')}\mathcal{L}_{LUC}(\hat{s},\hat{a},\hat{s}',\hat{a}',f))$$
(9)

where  $\hat{D}$  is the constructed noisy dataset. To be specific, the noised samples in  $\hat{D}$ , as  $\hat{x} = (\hat{s}, \hat{a}, \hat{s'})$ , are the noised version of samples, x = (s, a, s'), in the original offline dataset D, with  $\hat{x} = x + \lambda \cdot \epsilon$ , and  $\epsilon$  is the attached perturbation. Previous studies (Bai et al., 2022; Laskey et al., 2017; Zhang et al., 2022; Jiang et al., 2023) have empirically demonstrated the effectiveness of noise injection in regulating the out-of-distribution (OOD) performance of the trained agent. In the majority of our experiments,  $\epsilon$  is randomly drawn from a standard Gaussian distribution (also, in the OOD observation experiments in Appendix B.1,  $\epsilon$  is generated adversarially as in (Yang et al., 2022)). 270 Then the loss functions of the ensemble critic networks  $(L_c)$  and the actor network  $(L_a)$  are as, 271  $\mathcal{L}_{c} = \mathbb{E}_{(s,a,r,s')\sim\mathcal{D}} \left[ \left( r + \gamma \mathbb{E}_{a'\sim\pi(\cdot|s')} [\min_{i=1...K} f'_{i}(s',a')] - f(s,a) \right]^{2} + \frac{1}{K} \sum_{i=1}^{K} \mathcal{L}_{LUC}(f^{i},\pi)$ (10) 272 273 274  $\mathcal{L}_{a} = \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \pi(\cdot|s)} \big[ \min_{i=1...K} f'_{i}(s', a') \big]$ (11)275 276 277 To sum up, we present our overall approach in Algorithm 1, as follows, 278 279 Algorithm 1 The pseudocode of Lyapunov Uncertainty Control (LUC) algorithm **Input**: The offline dataset  $\mathcal{D}$ . Maximum of episode T. 281 Initialize the policy network, Q-network. 282 Perform the noise injection to generate the noisy dataset  $\mathcal{D}$ . 283 while t < T do 284 Sample mini-batch of transitions  $(s, a, r, s') \sim \mathcal{D}$  and transitions  $(\hat{s}, \hat{a}, \hat{s}') \sim \mathcal{D}$ Update the Q-network minimizing  $\mathcal{L}_q$  according to Eq.(10) 286 Update the policy network minimizing  $\mathcal{L}_{\pi}$  according to Eq.(11) 287 end while 288 **Output:** The learned policy network  $\pi$ . 289

### 4.3 THEORETICAL ANALYSIS

In this section, by Theorem 2, we aim to show that the value function obtained through k-step iterations of the empirical Bellman operator is determined by two factors concerning the true optimal value function: 1) the single-step Bellman uncertainty generated by policies in the policy candidate set, and 2) the growth tendency of Bellman uncertainty along trajectories generated by policies in the policy candidate set. The former has been the focus of previous methods like PBRL (Bai et al., 2022); however, Theorem 2 in this paper indicates that to learn a better value function, attention must be paid to both factors simultaneously.

**Theorem 2.** (*Performance lower bound.*) Given an MDP with max reward  $R_{max}$  and a dataset of size N. Given (s, a) pair, we denote its data density over the dataset is d(s, a). Given an empirical Bellman operator  $\hat{T}^{\Pi}$  and an arbitrary policy candidate set  $\Pi$ , where  $\hat{T}^{\Pi}f(s,a) =$  $r(s,a) + \gamma \mathbb{E}_{s'\sim \hat{P}(s'|s,a)} \max_{\pi \in \Pi} f(s', \pi)$ . Denote the learnt value function as  $f_k$ , with k iterations of  $\hat{f}_k = \hat{T}^{\Pi} \hat{f}_{k-1}$ , and the true optimal value function as  $f^*$ . Then we have,

306 307

308

310 311

290 291

292

$$\|\hat{f}_k - f^*\|_d \le \frac{C^*}{1 - \gamma} \cdot \sup_{\pi \in \Pi} \sum_{s_0} d(s_0) \zeta_{\hat{f}_k}(s_0, \pi) +$$

$$\mathcal{O}\left(\sup_{\substack{\pi \in \Pi, T \ge 0\\s_0, a_0, d(s_0, a_0) > 0}} \mathbb{E}_{P(\tau_T | \pi, s_0, a_0)} (\sum_{t=0}^{T-1} [\gamma^{t+1} \zeta_{\hat{f}_k, t+1} - \gamma^t \zeta_{\hat{f}_k, t}])^2\right)$$
(12)

where  $\zeta_{\hat{f}_k,t}$  is the Bellman uncertainty at time step t, i.e.,  $\zeta_{\hat{f}_k,t} = \|\mathcal{T}^{\Pi}\hat{f}_k - \hat{\mathcal{T}}^{\Pi}\hat{f}_k\|(s_t, a_t)$ .  $\tau_T$  is the trajectory with length of T. And  $C^*$  is assumed by  $\sup_{s,a} \frac{\pi^*(a|s)}{\pi_{\beta}(a|s)} \leq C^*$ .

<sup>316</sup> Proof of Theorem 2 is found in Appendix A.2. Theorem 2 primarily shows the following points:

1) The role of LUC term in enhancing the agent's performance - it helps optimize the lower bound of the agent's performance. Specifically, Theorem 2 illustrates that when iterated using the empirical Bellman operator, the difference between the learned value function  $\hat{f}_k$  and the true optimal value function  $f^*$ , which is also known as the performance lower bound of the offline algorithm, can be controlled through the LUC method.

2) One intuitive way to understand Theorem 2 is that we control the right term in Eq.(12) by adjusting the policy candidate set, thereby enhancing the lower bound of the method's output policy

324 performance. Specifically, we align the policy candidate set with the definition of Lyapunov reliable 325 policy (Definition 3) through the loss function in Eq. (5). 326

3) Such operation can control the right term of Eq.(12). And then consequently improve the lower 327 bound of the algorithm's performance. The proposed method would not conflict with the objective 328 of approaching the optimal policy, under the assumption of optimal coverage. 329

Then to further simplify the calculation complexity, Proposition 3 indicates that one-step Lyapunov 330 Uncertainty-penalization could bound the second term in Eq.12. 331

332 **Proposition 3.** If the first term of Eq.(12) is bounded, i.e.,  $\forall \pi \in \Pi$ , we have  $\mathbb{E}_{d(s_0)}\zeta_{\hat{f}_{i_k}}(s_0,\pi) \leq c$ , 333 then we can bound the second term with one-step Lyapunov Uncertainty-penalization, i.e.,  $\forall s \sim D$ , 334

$$\min_{\pi} [\gamma \mathbb{E}_{P(s'|s,\pi)} \zeta_{\hat{f}_{k},t+1}(s',\pi) - \zeta_{\hat{f}_{k},t+1}(s,\pi)] \Rightarrow \min_{\pi} \mathbb{E}_{P(\tau_{T}|\pi,s_{0}=s)} (\sum_{t=0}^{T-1} [\gamma^{t+1} \zeta_{\hat{f}_{k},t+1} - \gamma^{t} \zeta_{\hat{f}_{k},t}]$$

340

341

343

344

345

347 348

349 350

351

352 353

354

355

356

357

335

Furthermore, assume the dataset fully covers dynamics modes, i.e.,  $\forall s, a \in \mathcal{D}, supp(P(s'|s, a)) \subseteq$  $supp(\hat{P}(s'|s,a))$ , then the left part is controlled by Lyapunov value estimation in Eq.(5).

Then the left term in Proposition 3 could be empirically estimated by the Lyapunov value estimation as in Definition 3. Theorem 2 and Proposition 3 demonstrate that to control the performance lower 342 bound of the learned agent, despite controlling the current step's Bellman uncertainty, it is also important to control the one-step forward growth tendency of the Bellman uncertainty along with the whole trajectory, which is the main contribution of this paper. This helps the learned value function to be more likely to converge to the fixed point of the empirical Bellman operator, which is hence for controlling the performance lower bound of the learned agent. 346

- 5 **EXPERIMENTAL RESULTS**
- Our experiments primarily aim to address three key questions:
  - 1. Can LUC enhance the state-of-the-art performance on standard MuJoCo benchmarks?
  - 2. Is LUC capable of consistently learning reliable operation regions from noisy datasets with poor demonstrations?
  - 3. Does LUC exhibit superior generalization ability in avoiding deviations from reliable regions under various types of OOD perturbations?

358 Our experimental section includes the following components: first, we validate the performance of the method proposed in this paper on standard D4RL benchmarks, particularly on non-expert 359 datasets, demonstrating our method's superiority over others, addressing question 1. Next, to ad-360 dress question 2, we design noise data at different levels - where noise represents the discrepancy 361 between the policy and the optimal policy, resulting in varying degrees of poor demonstrations. We 362 then evaluate the performance of different algorithms on such noisy data. Subsequently, we introduce Out-of-distribution (OOD) MuJoCo benchmarks with various perturbations to increase the 364 likelihood of entering high-uncertainty states, assessing the agent's OOD generalization capability, 365 answering question 3. Finally, we conduct ablation experiments to verify the effectiveness of the 366 LUC method. A brief introduction of our code is provided in Appendix B.2.

367 368

369

## 5.1 LEARNING ON STANDARD MUJOCO BENCHMARKS

We assess our method using the D4RL benchmark (Fu et al., 2020) across various continuous-control 370 tasks and datasets. We compare LUC with several offline RL algorithms, including CQL (Kumar 371 et al., 2020), PBRL (Bai et al., 2022), MOPO (Yu et al., 2020), RORL (Yang et al., 2022), and 372 MOBILE (Sun et al., 2023). Among these, PBRL (Bai et al., 2022), MOPO (Yu et al., 2020), and 373 MOBILE (Sun et al., 2023) are most closely related to LUC as they are all based on uncertainty 374 penalization techniques<sup>1</sup>. 375

<sup>376</sup> <sup>1</sup>Unfortunately, as the LDM method (Kang et al., 2022) is mainly used in model-based RL as a constraint 377 on the model optimizer, it is unclear how this method could be properly used for the task of offline RL, so we did not make any comparison with this method at the current stage.

random ratios

Table 1: Normalized scores on standard MuJoCo tasks, averaged over 4 random seeds. Part of the results are reported in the RORL and MOBILE papers. Top two scores for each task are highlighted.  $(\cdot)$  indicates the average without 'expert' datasets.

	-	-					
	Task Name	CQL	PBRL	MOPO	RORL	MOBILE	LUC (Ours)
halfcheeta	ah random	31.3±3.5	11.0±5.8	<b>35.4</b> ±2.5	28.5±0.8	<b>39.3</b> ±3.0	$31.3 \pm 1.4$
	medium	46.9±0.4	57.9±1.5	42.3±1.6	66.8±0.7	<b>74.6</b> ±1.2	<b>68.2</b> ± 1.1
	medium-expert	95.0±1.4	92.3±1.1	63.3±38.0	107.8±1.1	108.2±2.5	<b>111.6</b> ± 1.2
	medium-replay	45.3±0.3	45.1±8.0	53.1±2.0	61.9±1.5	<b>71.7</b> ±1.2	<b>65.9</b> ± 1.9
	expert	97.3±1.1	92.4±1.7	-	105.2±0.7	-	$108.3 \pm 0.5$
hopper	random	5.3±0.6	26.8±9.3	11.7±0.4	31.4±0.1	<b>31.9±</b> 0.6	<b>31.9</b> ± 1.4
	medium	61.9±6.4	75.3±31.2	28.0±12.4	$104.8 \pm 0.1$	<b>106.6±</b> 0.6	<b>106.9</b> ± 0.4
	medium-expert	96.9±15.1	$110.8\pm0.8$	23.7±6.0	<b>112.7±</b> 0.2	112.6±0.2	<b>114.3</b> ± 1.1
	medium-replay	86.3±7.3	100.6±1.0	67.5±24.7	102.8±0.5	<b>103.9±</b> 1.0	$103.6 \pm 0.7$
	expert	106.5±9.1	110.5±0.4	-	112.8±0.2	-	$114.2 \pm 0.4$
walker2d	random	5.4±1.7	8.1±4.4	13.6±2.6	<b>21.4</b> ±0.2	17.9±3.0	<b>25.6</b> ± 1.2
	medium	79.5±3.2	89.6±0.7	17.8±19.3	<b>102.4</b> ±1.4	87.7±1.1	<b>103.6</b> ± 1.3
	medium-expert	109.1±0.2	110.1±0.3	44.6±12.9	121.2±1.5	115.2±0.7	$124.1 \pm 0.9$
	medium-replay	76.8±10.0	77.7±14.5	39.0±9.6	<b>90.4</b> ±0.5	89.9±1.5	<b>92.8</b> ± 1.5
	expert	109.3±0.1	108.3±0.3	-	<b>115.4</b> ±0.5	-	<b>116.2</b> ± 0.4
Average s	score	70.2	74.4	(36.7)	85.7	(80.0)	87.9 (81.7)
120	Halfcheetah-noisy	115	Hopper-no	bisy	120	Walker2d-n	oisy
Se 100	al the	S 110			9 100	1.1.	
260		- 0, 105 -			- D 60 -		
		aliz					

Figure 2: Results of CQL, PBRL SDC and LUC on tasks with different levels of non-expert data.

CQL PBRL SDC LUC(ours)

The results are presented in Table 1. It is evident that our method, LUC, outperforms other methods in most tasks and achieves a higher overall score. Particularly, LUC performs significantly better on non-expert datasets, notably the 'medium-expert' dataset, showcasing its robustness against the reward shift due to the conservative regularization term on non-expert data. Additionally, LUC's performance in the 'hopper' and 'walker2d' environments surpasses the State-of-the-art (SOTA), possibly due to these environments being more susceptible to noise from non-expert data.

5.2 LEARNING ON TASKS WITH DIFFERENT LEVELS OF NON-EXPERT NOISE

In this section, we modify the discrepancy between the behavioral policy and the optimal policy by blending datasets produced from expert and random policies at various proportions to generate noisy datasets at different levels. Evaluating these datasets not only validates the influence of non-expert problems on conventional conservative constraints but also confirms the robustness of the proposed LUC method against noise stemming from non-expert data. We compare several representative methods: CQL (Kumar et al., 2020), PBRL (Bai et al., 2022), and SDC (Zhang et al., 2022).

The results are depicted in Figure 2. It is evident that most constraint-based methods are influenced by non-expert issues, wherein an increase in the discrepancy between the behavioral policy and the optimal policy leads to a substantial performance decline. However, as the randomness level rises, the performance degradation of LUC is comparatively minor, suggesting that LUC demonstrates greater resilience to noise from non-expert data under suboptimal behavioral policy settings. Partic-ularly at elevated randomness levels (e.g., 0.9), LUC can maintain effective performance on these benchmarks.

#### TESTING ON OUT-OF-DISTRIBUTION MUJOCO BENCHMARKS 5.3

To evaluate the agent's capacity to avoid straying from reliable regions, we introduce three types of OOD perturbations, applying varying intensities of perturbations at different intervals to the agent

432 employed for a higher risk of deviation. This investigation encompasses three combinations of noise 433 intensities and intervals<sup>2</sup>. It is important to highlight that our approach in this study differs from the 434 perturbation noise discussed in Yang et al. (2022); our method modifies the actual state in which the 435 agent operates, rather than solely perturbing state observations. The research is centered on three 436 MuJoCo environments: Halfcheetah, Hopper, and Walker2d; with the agent trained on 'mediumexpert' datasets. 437

438 439 440

451

460 461

Table 2: The results of OSR, RORL and LUC (ours) on Out-of-distribution MuJoCo benchmarks. The highest scores for each task are highlighted.

	Ha	lfcheetah-o	od	]	Hopper-ood	[	V	Valker2d-oc	od
	small	medium	large	small	medium	large	small	medium	large
OSR	93.8	90.4	88.7	111.3	103.4	85.7	112.7	110.5	105.8
RORL	102.7	94.3	82.4	111.5	92.8	72.7	117.4	107.3	86.9
LUC(ours)	104.8	102.9	99.0	110.9	106.6	83.3	119.4	114.5	111.9

448 We have chosen two key algorithms, OSR and RORL, tailored for managing OOD states and obser-449 vations, to contrast with the proposed LUC in these OOD benchmarks. The outcomes are detailed in 450 Table 2. Analysis reveals that LUC surpasses the other two methods across the majority of tasks, notably demonstrating substantial benefits in extensive OOD perturbation assignments like Halfcheetah and Walker2d. This implies that these settings might be more sensitive to OOD perturbations, 452 necessitating advanced the agent to tackle OOD scenarios. 453

454 To delve deeper into the factors contributing to the superior performance of LUC in Halfcheetah-455 ood and Walker2d-ood tasks, we present the visualized results of LUC in Figure 3. The analysis 456 reveals that each perturbation event leads the agent into high-uncertainty regions; however, LUC 457 effectively prevents error accumulation and guides the agent back to lower-error (reliable) regions. This underscores the robustness of LUC in handling OOD scenarios by constraining the agent to 458 operate within the reliable regions. 459



Figure 3: Visualized results of LUC on OOD MuJoCo benchmarks - 'Halfcheetah' and 'Walker2d', with large scales of perturbations. The Bellman uncertainty (error) is estimated by the standard deviation uncertainty based on the learned Q-ensembles. The interval between every two frames is about five steps.

471 472 473

474 475

476

477

478

479

480

481

482

468

469

470

#### 5.4 MORE COMPLICATED ENVIRONMENTS - ANTMAZE AND ADROIT

Compared to the MuJoCo environment, the AntMaze and Adroit environments require the agent to have the ability of multi-step dynamic planning, making them considered a more complex scenario. The results are shown in Table 3 and 4. In the AntMaze environment, based on the size and shape of the maze, it can be categorized into 'umaze,' 'medium,' and 'large'; and based on different tasks, it can be classified as 'diverse' and 'play'. While the Adroit domain features three types of datasets: demonstration data from humans ("human"), expert data from a reinforcement learning policy ("expert"), and mixed data combining human demonstrations with an imitation policy ("cloned"). The tasks in Adroit are more complex than those in the Gym domain, and the inclusion of human demonstrations adds an additional layer of difficulty.

<sup>&</sup>lt;sup>2</sup>Details regarding the construction of Out-of-distribution MuJoCo benchmarks are outlined in Appendix B.3.

Here, we compare CQL (Kumar et al., 2020), IQL (Kostrikov et al., 2022), SPOT (Wu et al., 2022),
ATAC (Cheng et al., 2022), SDC (Zhang et al., 2022), and OSR-10 (Jiang et al., 2023) in AntMaze,
while CQL (Kumar et al., 2020) and PBRL (Bai et al., 2022) in Adroid.

Table 3: Results of LUC(ours), CQL, IQL, SPOT, ATAC, SDC and OSR-10 on offline AntMaze tasks averaged over 4 seeds. We bold the highest scores in each task.

		CQL	IQL	SPOT	ATAC	SDC	<b>OSR-10</b>	LUC(Ours)
	umaze	82.6	87.5	93.5	70.6	89.0	89.9	94.3±1.3
	umaze-diverse	10.2	62.2	40.7	54.3	57.3	<b>74.0</b>	62.1±3.7
AntMaze	medium-play	59.0	71.2	74.7	72.3	71.9	66.0	80.1±2.2
	medium-diverse	46.6	70.0	79.1	68.7	78.7	80.0	78.7±3.9
	large-play	16.4	39.6	35.3	38.5	37.2	37.9	44.3±7.2
	large-diverse	3.2	47.5	36.3	43.1	33.2	37.9	41.7±6.1
average		36.3	63.0	59.9	57.9	61.2	64.3	66.9

Table 4: Results of **LUC(ours**), CQL and PBRL on offline Adroit tasks averaged over 4 seeds. We bold the highest scores in each task.

	pen-hu.	pen-cl.	pen-ex.	hammer-hu.	hammer-cl.	hammer-ex.	door-hu.	door-cl.	door-ex.	avg.
CQL	37.5	39.2	107.0	4.4	2.1	86.7	9.9	0.4	101.5	43.2
PBRL	35.4	74.9	137.7	0.4	0.8	127.5	0.1	4.6	95.7	53.0
LUC(ours)	41.6	68.2	142.5	7.9	4.7	130.4	7.8	8.4	103.4	57.2

From these benchmarks, we observe that our method outperforms the other methods in most benchmarks. In particular, our method performs much better in the "hammer," "door," and "large" maze benchmarks than other methods. This could be attributed to our method's ability to scope a safe region for the agent to stably operate within, thus addressing more complex environmental dynamics.

## 514 5.5 ABLATION STUDY

516 iments on the proposed LUC method to con-517 firm its contribution to the overall framework. 518 The findings, illustrated in Figure 4, demon-519 strate a substantial performance enhancement in 520 the LUC module compared to LUC without re-521 ward shaping, particularly on certain non-expert 522 datasets. This highlights the pivotal role of the 523 LUC method in improving the performance of 524 pessimistic offline RL approaches.

Furthermore, we conducted experiments on Outof-Distribution (OOD) observation benchmarks to further validate the robustness of the proposed method, as detailed in Appendix B.1. More in-

In this section, we performed ablation experiments on the proposed LUC method to confirm its contribution to the overall framework. *B*.



Figure 4: Ablation study.

## 531 6 CONCLUSION

This paper aims to identify a reliable operational region for the agent based on offline data. To achieve this, we introduce the Lyapunov Uncertainty Control (LUC) algorithm in an offline, model-free manner. Theoretically, in deterministic MDPs or when the dataset fully covers all dynamic modes, LUC can confine the agent's operations within low-uncertainty areas, thereby enhancing decision-making reliability. Empirically, LUC-trained agents demonstrate superior robustness and reliability in high-risk scenarios compared to various other methods. In future works, LUC can serve as a versatile tool in diverse offline reinforcement learning frameworks, including model-based approaches, potentially paving the way for new research opportunities.

490

491

501 502

510

511

512

513

515

525

526

527

528

529 530

## 540 REFERENCES

567

- Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline reinforcement learning with diversified q-ensemble. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 7436-7447, 2021. URL https://proceedings.neurips.cc/paper/2021/hash/ 3d3d286a8d153a4a58156d0e02d8570c-Abstract.html.
- Chenjia Bai, Lingxiao Wang, Zhuoran Yang, Zhi-Hong Deng, Animesh Garg, Peng Liu, and Zhao ran Wang. Pessimistic bootstrapping for uncertainty-driven offline reinforcement learning. In
   *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022. URL https://openreview.net/forum?id=
   Y4cs1Z3HnqL.
- Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor critic for offline reinforcement learning. In *International Conference on Machine Learning*, pp. 3852–3878. PMLR, 2022.
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4RL: datasets for deep
   data-driven reinforcement learning. *CoRR*, abs/2004.07219, 2020. URL https://arxiv.
   org/abs/2004.07219.
- Mohammad Ghavamzadeh, Marek Petrik, and Yinlam Chow. Safe policy improvement by minimizing robust baseline regret. *Advances in Neural Information Processing Systems*, 29, 2016.
- Ke Jiang, Jia-Yu Yao, and Xiaoyang Tan. Recovering from out-of-sample states via inverse dynamics
   in offline reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Ying Jin, Zhuoran Yang, and Zhaoran Wang. Is pessimism provably efficient for offline rl? In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 5084–5096. PMLR, 2021. URL http://proceedings.mlr. press/v139/jin21e.html.
- Katie Kang, Paula Gradu, Jason J Choi, Michael Janner, Claire Tomlin, and Sergey Levine. Lya punov density models: Constraining distribution shift in learning-based control. In *International Conference on Machine Learning*, pp. 10708–10733. PMLR, 2022.
- Eric C Kerrigan. *Robust constraint satisfaction: Invariant sets and predictive control.* University of London, 2000.
- B. Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A. Al Sallab, Senthil Kumar Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE Trans. Intell. Transp. Syst.*, 23(6):4909–4926, 2022. doi: 10.1109/TITS.2021.3054625.
  URL https://doi.org/10.1109/TITS.2021.3054625.
- Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Vir- tual Event, April 25-29, 2022.* OpenReview.net, 2022. URL https://openreview.net/
   forum?id=68n2s9ZJWF8.
- Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 11761–11771, 2019. URL https://proceedings.neurips.cc/paper/2019/hash/c2073ffa77b5357a498057413bb09d3a-Abstract.html.

594 595 596 597 598 599	Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: An- nual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 0d2b2061826a5df3221116a5085a6052-Abstract.html.
600 601 602 603 604 605	Michael Laskey, Jonathan Lee, Roy Fox, Anca D. Dragan, and Ken Goldberg. DART: noise injection for robust imitation learning. In <i>1st Annual Conference on Robot Learning, CoRL 2017, Mountain</i> <i>View, California, USA, November 13-15, 2017, Proceedings</i> , volume 78 of <i>Proceedings of Ma- chine Learning Research</i> , pp. 143–156. PMLR, 2017. URL http://proceedings.mlr. press/v78/laskey17a.html.
606 607 608	Andrew Lobbezoo, Yanjun Qian, and Hyock-Ju Kwon. Reinforcement learning for pick and place operations in robotics: A survey. <i>Robotics</i> , 10(3):105, 2021. doi: 10.3390/robotics10030105. URL https://doi.org/10.3390/robotics10030105.
609 610	Yixiu Mao, Hongchang Zhang, Chen Chen, Yi Xu, and Xiangyang Ji. Supported value regularization for offline reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
611 612 613 614	Rémi Munos. Error bounds for approximate value iteration. In <i>Proceedings of the National Confer-</i> <i>ence on Artificial Intelligence</i> , volume 20, pp. 1006. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2005.
615 616 617	Kishan Panaganti, Zaiyan Xu, Dileep Kalathil, and Mohammad Ghavamzadeh. Robust reinforce- ment learning using offline data. <i>Advances in neural information processing systems</i> , 35:32211– 32224, 2022.
618 619 620	Charles Richter and Nicholas Roy. Safe visual navigation via deep learning and novelty detection. <i>Robotics: Science and Systems Foundation</i> , 2017.
621 622	Laixi Shi and Yuejie Chi. Distributionally robust model-based offline reinforcement learning with near-optimal sample complexity. <i>Journal of Machine Learning Research</i> , 25(200):1–91, 2024.
623 624 625 626	Yihao Sun, Jiaji Zhang, Chengxing Jia, Haoxin Lin, Junyin Ye, and Yang Yu. Model-bellman inconsistency for model-based offline reinforcement learning. In <i>International Conference on Machine Learning</i> , pp. 33177–33194. PMLR, 2023.
627 628 629	Shengpu Tang and Jenna Wiens. Model selection for offline reinforcement learning: Practical con- siderations for healthcare settings. In <i>Machine Learning for Healthcare Conference</i> , pp. 2–35. PMLR, 2021.
630 631	Christopher J. C. H. Watkins and Peter Dayan. Technical note q-learning. <i>Mach. Learn.</i> , 8:279–292, 1992. doi: 10.1007/BF00992698. URL https://doi.org/10.1007/BF00992698.
633 634 635 636 637 638 639	Jialong Wu, Haixu Wu, Zihan Qiu, Jianmin Wang, and Mingsheng Long. Supported pol- icy optimization for offline reinforcement learning. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural In- formation Processing Systems 35: Annual Conference on Neural Information Process- ing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/ caa934a507a952698d54efb24845fc4b-Abstract-Conference.html.
640 641	Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. <i>CoRR</i> , abs/1911.11361, 2019. URL http://arxiv.org/abs/1911.11361.
642 643 644 645	Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-consistent pessimism for offline reinforcement learning. <i>Advances in neural information processing systems</i> , 34:6683–6694, 2021.
646 647	Rui Yang, Chenjia Bai, Xiaoteng Ma, Zhaoran Wang, Chongjie Zhang, and Lei Han. RORL: robust offline reinforcement learning via conservative smoothing. <i>CoRR</i> , abs/2206.02829, 2022. doi: 10.48550/arXiv.2206.02829. URL https://doi.org/10.48550/arXiv.2206.02829.

Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y. Zou, Sergey Levine, Chelsea Finn, and Tengyu Ma. MOPO: model-based offline policy optimization. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Con-ference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ a322852ce0df73e204b7e67cbbef0d0a-Abstract.html. 

Hongchang Zhang, Jianzhun Shao, Yuhang Jiang, Shuncheng He, Guanwen Zhang, and Xiangyang Ji. State deviation correction for offline reinforcement learning. In Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, pp. 9022–9030. AAAI Press, 2022. URL https://ojs. aaai.org/index.php/AAAI/article/view/20886.

#### 702 APPENDIX А 703

721 722

738

739

745

753

704 A.1 PROOFS OF MAIN THEOREMS IN SEC.4.1 705

Proposition 1. (Existence of reliable policy.) Suppose the dataset have a sufficient coverage over 706 the optimal policy, i.e.,  $\sup_{s,a} \frac{\pi^*(a|s)}{\pi_{\beta}(a|s)} \leq C^*$ ). Then there exists a reliable policy. 708

*Proof sketch.* From the assumption that the dataset have a sufficient over the optimal policy, we 709 have all the behaviors of the optimal policy would be supported by the dataset, i.e.,  $\pi^* \subseteq \mathcal{D}$ . Then 710 the optimal stationary state distribution is also supported by the dataset  $d^{\pi^*}(s) \subseteq \mathcal{D}$ . Due to the fact 711 that all the states over the dataset are contained in the  $f - \pi^*$  reliable region, and  $d^{\pi^*}(s) \subset \mathcal{D}$ , so the 712 optimal policy would never operate beyond its  $f - \pi^*$  reliable region, hence it is a reliable policy. 713

714 **Theorem 1.** In an MDP with deterministic transition, Lyapunov policy  $\pi$  is a reliable policy as defined in Definition 2. 715

716 Proof of Theorem 1. We prove this theorem in a contradiction way. Denote the deterministic 717 transition as M(s, a) = s'. First, we suppose the Lyapunov policy  $\pi$  is not a reliable policy, then 718 we have:  $\forall s_0 \in \mathcal{D} \cap \mathcal{G}_f(\pi), \exists t, s_t \in supp(P(s_t|s_0, \pi))$ , such that  $s_t \notin \mathcal{G}_f(\pi)$ . Then we have 719  $\zeta_f(s_t,\pi) > c$ . Then we aim to find the contradiction. From the definition of Lyapunov policy, 720  $\forall s \in \mathcal{D}$ , we have,

$$c < \zeta_f(s_t, \pi) \le \max_{\hat{a} \in \pi} \zeta_f(P(s_{t-1}, \hat{a}), \pi) \le \epsilon_f(s_{t-1}, \pi)$$
(13)

Then we have  $\epsilon_f(s_{t-1}, \pi) > c$ . By recurrently applying the above derivation by t times, we would 723 have  $\zeta_f(s_0,\pi) > c$ . Then we have  $s_0 \notin \mathcal{G}_f(\pi)$ , which is conflicted with  $s_0 \in \mathcal{D} \cap \mathcal{G}_f(\pi)$ . Completing 724 the proof, and we can conclude that the Lyapunov policy is reliable. 725

726 **Proposition 2.** Suppose the action distribution of new policy  $\pi(a|s)$  is positive correlated to the 727 learnt Q function f(s, a), i.e.,  $\pi(a|s) \propto f(s, a)$ . Then the proposed Lyapunov value estimation 728 induces a Lyapunov policy as defined in Definition 3.

729 *Proof of Proposition 2.* Denote the empirical behavior of the dataset  $\mathcal{D}$  as  $\pi_{\beta}$ , whose actions are 730 always supported by the dataset. The Lyapunov value estimation implicitly penalize the new policy 731  $\pi$  at an given state s with two aspects: 732

1)  $\min_{\pi} \mathbb{E}_{a \sim \pi(a|s)} \zeta_f(s, a)$ ; This constrains the new policy would not generated OOD actions beyond 733 the demonstration of the offline data, i.e.,  $supp(\pi(a|s)) \subset supp(\pi_{\beta}(a|s))$ . In previous works Wu 734 et al. (2022); Mao et al. (2024), such supported constraint is often achieved in a data density based 735 way as  $\min_{\pi} \sum_{a \notin \pi_{\beta}} \pi(a|s)$ . Then we will show that the current step's error controlling minimizes 736 the upper bound of the above supported constraint, 737

$$\min_{\pi} \sum_{a \notin \pi_{\beta}} \pi(a|s) \Leftrightarrow \max_{\pi} \mathbb{E}_{a \in \pi(a|s)} d(s, a)$$
(14)

740  
741 
$$\leq^{(a)} \max_{\pi} C \cdot \mathbb{E}_{a \in \pi(a|s)} \frac{1}{\zeta_f^2(s,a)}$$
 (15)  
742

$$\Leftrightarrow \min_{\pi} \mathbb{E}_{a \in \pi(a|s)} \zeta_f^2(s, a) \tag{16}$$

$$\Leftrightarrow^{(b)} \min \mathbb{E}_{a \in \pi(a|s)} \zeta_f(s, a) \tag{17}$$

746 The inequality (a) holds because of the Lemma 1 in Appendix A.2. The equivalence (b) holds 747 because the Bellman uncertainty is always non-negative.

748 Then with this constraint, we can have the new policy would reject the actions that is not supported 749 by the behavior policy. This means, the new policy would only select the data-supported actions. 750

2) 
$$\forall s, a \in \mathcal{D}$$
, we have  $\zeta_f(M(s, a), \pi) \leq \zeta_f(s, a)$ . This could also be consider as,  

$$\min_{a \in \pi_\beta} \zeta_f(s, a) - \zeta_f(M(s, a), \pi)$$
(18)

754 
$$\leq^{(a)} \min_{a \in \pi} \zeta_f(s, a) - \zeta_f(M(s, a), \pi)$$
(19)

$$\leq \zeta_f(s, a_m) - \zeta_f(M(s, a_m), \pi) \tag{20}$$

where  $a_m = \arg \max_{a_m \in \pi} \zeta_f(M(s, a_m), \pi)$ . The (a) holds because of the assumption that the condition 1) is perfectly confirmed, i.e.,  $supp(\pi(a|s)) \subset supp(\pi_{\beta}(a|s))$ .

Then we have  $\forall s \in \mathcal{D}, \max_{\hat{a} \in \pi} \zeta_f(M(s, \hat{a}), \pi) \leq \zeta_f(s, \pi)$ . And we can conclude that the Lyapunov value estimation would help to induce a Lyapunov policy.

A.2 PROOFS OF MAIN THEOREM WITH STOCHASTIC TRANSITION SETTING

First we define the policy candidate set  $\Pi$  based on a version space of all the functions  $f \in \mathcal{F}$ , where we have  $\forall \pi \in \Pi, \exists f \in \mathcal{F}, \pi(a|s) \propto f(s, a)$ , and  $\forall f \in \mathcal{F}, \exists \pi \in \Pi, \pi$  is the greedy policy according to f. Then we define a corresponding Bellman operator  $\mathcal{T}^{\Pi}f(s,a) = r + \gamma \mathbb{E}_{s'\sim P} \max_{\pi \in \Pi} f(s',\pi)$ and its empirical version is  $\hat{\mathcal{T}}^{\Pi}f(s,a) = r + \gamma \mathbb{E}_{s'\sim \hat{P}} \max_{\pi \in \Pi} f(s',\pi)$ . Then before the introduction of theoretical results, a basic assumption should be made.

Assumption 1. (Optimal coverage.) (Xie et al., 2021) We assume the dataset have sufficient coverage over the optimal policy's visitation, i.e.,  $\sup_{s,a} \frac{\pi^*(a|s)}{\pi_{\beta}(a|s)} \leq C^*$ , and  $\pi^* \in \Pi$ .

Similar assumptions has been utilized in theoretical analysis for offline RL (Xie et al., 2021). Compared with the more common assumption - Concentrability assumption (Munos, 2005; Kumar et al., 2019) that the dataset should fully cover the whole state space, Assumption 1 is much looser and more feasible in practice.

**Definition 4.** (*Recoverability*) Define the recoverability of a given policy  $\pi$  from the given  $(s_0, a_0)$ pair,

$$R(\pi)\big|_{s_0,a_0} = \inf_{T \ge 0} \mathbb{E}_{s_T,a_T \sim P(s_T,a_T | \pi, s_0, a_0)} \frac{d(s_T, a_T)}{d(s_0, a_0)}$$
(21)

where d(s, a) is the data density at (s, a).

761

762

778

779

785 786 787

789

793 794

796 797

**Definition 5.** (*Recoverability risk*) We define the most risk of a policy  $\pi$  that the agent is able to recover to the regions with low Bellman uncertainty, i.e., its familiar regions, from the given  $(s_0, a_0)$ pair,

$$Risk(\pi)\big|_{s_0,a_0} = \sup_{T \ge 0} \mathbb{E}_{s_T,a_T \sim P(s_T,a_T|\pi,s_0,a_0)} \frac{\|\mathcal{T}^{\pi}f - \hat{\mathcal{T}}^{\pi}f\|(s_T,a_T)}{\|\mathcal{T}^{\pi}f - \hat{\mathcal{T}}^{\pi}f\|(s_0,a_0)} = \sup_{T \ge 0} \mathbb{E}_{s_T,a_T} \frac{\zeta_T}{\zeta_0}$$
(22)

where  $\mathcal{T}^{\pi}$  is the true Bellman operator and  $\hat{\mathcal{T}}^{\pi}$  is the empirical Bellman.

**Lemma 1.** Given an MDP with max reward  $R_{max}$  and a dataset of size N. The dimension of state space is |S| and that of action space is |A|. Given (s, a) pair, we denote its data density over the dataset is d(s, a). Then with probability  $1 - \delta$ , we have,

$$d(s,a) \le \gamma^2 \cdot R_{max}^2 \frac{2}{N \cdot \|\mathcal{T}^{\pi} f - \hat{\mathcal{T}}^{\pi} f\|^2(s,a)} \log(\frac{|S||A| \cdot 2^{|S|}}{\delta})$$
(23)

where  $\|\mathcal{T}^{\pi}f - \hat{\mathcal{T}}^{\pi}f\|(s,a)$  is Bellman uncertainty.

Proof of Lemma 1.

$$\|\mathcal{T}^{\pi}f - \hat{\mathcal{T}}^{\pi}f\|(s,a) = \gamma\|\sum_{s'} (\hat{P}(s'|s,a) - P(s'|s,a)) \cdot f(s',\pi)\|$$
(24)

$$\leq \gamma \| \sum_{s'} (\hat{P}(s'|s,a) - P(s'|s,a)) \| \cdot \| f(s',\pi) \|$$
(25)

$$\leq \gamma \cdot R_{max} \cdot \|\hat{P}(s'|s,a) - P(s'|s,a)\|_1$$
(26)

$$\leq^{(a)} \gamma \cdot R_{max} \cdot \sqrt{\frac{2}{N(s,a)} \log(\frac{|S||A| \cdot 2^{|S|}}{\delta})} \tag{27}$$

808  
809 
$$\Rightarrow d(s,a) \le \gamma^2 \cdot R_{max}^2 \frac{2}{N \cdot \|\mathcal{T}^{\pi} f - \hat{\mathcal{T}}^{\pi} f\|^2(s,a)} \log(\frac{|S||A| \cdot 2^{|S|}}{\delta})$$
(28)

The inequality (a) holds because of the **Proposition 9** in Ghavamzadeh et al. (2016). And N(s, a) is the number of (s, a) samples in the dataset, so the density  $d(s, a) = \frac{N(s, a)}{N}$ . Completing the proof.

Then we give Lemma 2, which is a general formulation of Lemma ?? in the main text.

**Lemma 2.** For any policy  $\pi$  and  $(s_0, a_0)$  pair, we have,

$$Risk(\pi)\big|_{s_0,a_0} = \sup_{T \ge 0} \mathbb{E}_{\pi} \left[ \frac{\sum_{t=0}^{T-1} [\gamma^{t+1} \zeta_{t+1} - \gamma^t \zeta_t] + \zeta_0}{\zeta_0 \cdot \gamma^T} \big| s_0, a_0 \right]$$
(29)

Proof of Lemma 2.

$$Risk(\pi)\big|_{s_0,a_0} = \sup_{T \ge 0} \mathbb{E}_{\pi} \left[ \frac{\zeta_T}{\zeta_0} \big| s_0, a_0 \right]$$
(30)

$$= \sup_{T \ge 0} \mathbb{E}_{\pi} \left[ \frac{\gamma^{T} \zeta_{T} - \gamma^{T-1} \zeta_{T-1} + \gamma^{T-1} \zeta_{T-1} - \dots + \gamma \zeta_{1} - \zeta_{0} + \zeta_{0}}{\zeta_{0} \cdot \gamma^{T}} \Big| s_{0}, a_{0} \right]$$
(31)

$$= \sup_{T \ge 0} \mathbb{E}_{\pi} \left[ \frac{\sum_{t=0}^{T-1} [\gamma^{t+1} \zeta_{t+1} - \gamma^{t} \zeta_{t}] + \zeta_{0}}{\zeta_{0} \cdot \gamma^{T}} | s_{0}, a_{0} \right]$$
(32)

## 827 Completing the proof.

**Lemma 3.** Given an arbitrary Bellman operator (maybe empirical Bellman operator)  $\mathcal{T}$  and an arbitrary policy candidate set  $\Pi$ . We have  $\mathcal{T}f(s, a) = r + \gamma \mathbb{E}_{s'} \max_{\pi \in \Pi} f(s', \pi)$ . Then for any value function  $f_1, f_2 \in \mathcal{F}$  and  $t \ge 0$ , we have,

$$\|\mathcal{T}^{(t)}f_1 - \mathcal{T}^{(t)}f_2\|_d \le \sup_{s_0, a_0, d(s_0, a_0) > 0} \frac{\gamma^t}{R(\hat{\pi})|_{s_0, a_0}} \|f_1 - f_2\|_d$$
(33)

We denote the greedy policy induced by  $f_1$  as  $\pi_1$  and the greedy policy induced by  $f_2$  as  $\pi_2$ . Then  $\hat{\pi}$ is the pessimistic policy of  $f_1$  and  $f_2$ , i.e.,  $\hat{\pi}(a|s) = \pi_1(a|s)$  if  $f_1(s,\pi_1) \le f_2(s,\pi_2)$  and  $\hat{\pi}(a|s) = \pi_2(a|s)$  if  $f_2(s,\pi_2) \le f_1(s,\pi_1)$ . And  $||x||_d = \sum_x d(x)|x|$  is a distributional weighted norm, where d here is the density of the dataset.

*Proof of Lemma 3.* First we denote  $\pi_1(a|s) = \arg \max_{\pi \in \Pi} f_1(s,\pi)$  and  $\pi_2(a|s) = \arg \max_{\pi \in \Pi} f_2(s,\pi)$ . Then,

$$-\mathcal{T}f_2(s,a) = \gamma \mathbb{E}_{P(s'|s,a)}[f_1(s',\pi_1) - f_2(s',\pi_2)]$$
(34)

$$\leq \gamma \mathbb{E}_{P(s'|s,a)}[f_1(s',\pi_1) - f_2(s',\pi_1)]$$
(35)

On the other hand,

$$(\mathcal{T}f_1 - \mathcal{T}f_2)(s, a) \ge \gamma \mathbb{E}_{P(s'|s, a)}[f_1(s', \pi_2) - f_2(s', \pi_2)]$$
(36)

846 Then we construct  $\hat{\pi}(a|s) = \pi_1(a|s)$  if  $f_1(s,\pi_1) \le f_2(s,\pi_2)$  and  $\hat{\pi}(a|s) = \pi_2(a|s)$  if  $f_2(s,\pi_2) \le f_1(s,\pi_1)$ . So we have,

$$|\mathcal{T}f_1 - \mathcal{T}f_2|(s,a) \le \gamma |\mathbb{E}_{P(s'|s,a)}[f_1(s',\hat{\pi}) - f_2(s',\hat{\pi})]|$$
(37)

Then if we recursively apply the  $\hat{\pi}$ , we would have,

 $(\mathcal{T}f_1)$ 

$$|\mathcal{T}^{(t)}f_1 - \mathcal{T}^{(t)}f_2|(s,a) \le \gamma^t \cdot |\mathbb{E}_{P(s_t,a_t|\hat{\pi},s_0=s,a_0=a)}[f_1(s_t,a_t) - f_2(s_t,a_t)]|$$
(38)  
Then we aim to bound the  $\|\mathcal{T}^{(t)}f_1 - \mathcal{T}^{(t)}f_2\|_d$ ,

$$\|\mathcal{T}^{(t)}f_1 - \mathcal{T}^{(t)}f_2\|_d = \sum_{s_0, a_0} |\mathcal{T}^{(t)}f_1 - \mathcal{T}^{(t)}f_2|(s_0, a_0)d(s_0, a_0)$$
(39)

$$\leq \gamma^{t} \cdot \sum_{s_{0}, a_{0}} |\mathbb{E}_{P(s_{t}, a_{t} | \hat{\pi}, s_{0}, a_{0})}[f_{1}(s_{t}, a_{t}) - f_{2}(s_{t}, a_{t})]|d(s_{0}, a_{0})$$
(40)

$$\leq \gamma^{t} \cdot \sum_{s_{0}, a_{0}} \sum_{s_{t}, a_{t}} d(s_{0}, a_{0}) P(s_{t}, a_{t} | \hat{\pi}, s_{0}, a_{0}) |[f_{1}(s_{t}, a_{t}) - f_{2}(s_{t}, a_{t})]|$$
(41)

$$\leq \sup_{s_0, a_0, d(s_0, a_0) > 0} \frac{\gamma^t}{R(\hat{\pi})\big|_{s_0, a_0}} \cdot \sum_{s_t, a_t} d(s_t, a_t) |[f_1(s_t, a_t) - f_2(s_t, a_t)]|$$
(42)

$$= \sup_{s_0, a_0, d(s_0, a_0) > 0} \frac{\gamma^t}{R(\hat{\pi})|_{s_0, a_0}} \cdot \|[f_1(s_t, a_t) - f_2(s_t, a_t)]\|_d$$
(43)

Completing the proof.

**Corollary 1.** *Especially, if the*  $f_2$  *in Lemma 3 is the fix point of*  $\mathcal{T}$ *, as*  $f^*$ *, then we have,* 

$$\|\mathcal{T}^{(t)}f_1 - f^*\|_d \le \sum_{s_0, a_0} \frac{\gamma^t}{R(\pi_1)|_{s_0, a_0}} \|f_1 - f^*\|_d \tag{44}$$

This is easily obtained by the fact that  $\forall s, f^*(s, \pi^*) \ge f_1(s, \pi_1)$ .

**Lemma 4.** Denote the learnt value function as  $f_k$ , with k iterations of  $f_k = \mathcal{T} f_{k-1}$ , and the fixed point of  $\mathcal{T}$  is  $f^*$ . Then,

$$\|f_k - f^*\|_d \le R(\pi_0) \cdot \gamma^k \cdot \|\triangle_{(0)}\|_d + \epsilon_{max} \cdot \sum_{t=1}^k [R(\hat{\pi}_{t-1}) \cdot \gamma^t] + \epsilon_{max}$$
(45)

where  $R(\pi_k) = \sum_{s_0, a_0} \frac{1}{R(\pi_k) \Big|_{s_0, a_0}}$ ,  $\epsilon_{max} = \max_{t \le k-1} \|f_{t+1} - \mathcal{T}f_t\|_d$  and  $\triangle_{(0)} = \|f_0 - f^*\|_{\infty}$ .

Proof of Lemma 4.

$$||f_k - f^*||_d \le ||\mathcal{T}f_{k-1} - f^*||_d + ||f_k - \mathcal{T}f_{k-1}||_d$$
(46)

$$\leq \|\mathcal{T}^{(2)}f_{k-2} - f^*\|_d + \|\mathcal{T}f_{k-1} - \mathcal{T}^{(2)}f_{k-2}\|_d + \epsilon_{max}$$
(47)

$$\leq \|\mathcal{T}^{(2)}f_{k-2} - f^*\|_d + R(\pi_{k-1}) \cdot \gamma \cdot \epsilon_{max} + \epsilon_{max}$$
(48)

$$\leq \|\mathcal{T}^{(k)}f_0 - f^*\|_d + \epsilon_{max} \cdot \sum_{t=1}^k [R(\hat{\pi}_{t-1}) \cdot \gamma^t] + \epsilon_{max}$$

$$\tag{50}$$

$$\leq^{(a)} R(\pi_0) \cdot \gamma^k \cdot \|\triangle_{(0)}\|_d + \epsilon_{max} \cdot \sum_{t=1}^k [R(\hat{\pi}_{t-1}) \cdot \gamma^t] + \epsilon_{max}$$
(51)

The inequality (a) holds because of Corollary 1. Completing the proof.

. . . . . . . . . . . . .

**Lemma 5.** Given an MDP with max reward  $R_{max}$  and a dataset of size N. The dimension of state space is |S| and that of action space is |A|. Given (s, a) pair, we denote its data density over the dataset is d(s, a). Given an empirical Bellman operator  $\hat{T}$  and an arbitrary policy candidate set  $\Pi$ , where  $\hat{T}f(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}(s'|s, a)} \max_{\pi \in \Pi} f(s', \pi)$ . Denote the learnt value function as  $f_k$ , with k iterations of  $\hat{f}_k = \hat{T}\hat{f}_{k-1}$ , and the fixed point of  $\hat{T}$  is  $\hat{f}^*$ . Then we have,

$$\|\hat{f}_{k} - \hat{f}^{*}\|_{d} \leq \mathcal{O}\left(\sup_{s_{0}, a_{0}, d(s_{0}, a_{0}) > 0} \sup_{\pi \in \Pi, T \geq 0} \mathbb{E}_{s_{T}, a_{T} \sim P(s_{T}, a_{T} | \pi, s_{0}, a_{0})} (\sum_{t=0}^{T-1} \gamma^{t+1} \zeta_{t+1} - \gamma^{t} \zeta_{t})^{2}\right)$$
(52)

where  $\zeta_t$  is the Bellman uncertainty at time step t, i.e.,  $\zeta_t = \|\mathcal{T}^{\pi}f - \hat{\mathcal{T}}^{\pi}f\|(s_t, a_t)$ .

Proof of Lemma 5. From Lemma 4 we have known that if we want to bound  $\|\hat{f}_k - \hat{f}^*\|_d$ , we should bound  $\frac{1}{R(\pi)|_{s_0,a_0}}$  at each time steps.

$$\frac{1}{R(\pi)\big|_{s_0,a_0}} = \sup_{T \ge 0} \mathbb{E}_{s_T,a_T \sim P(s_T,a_T \mid \pi, s_0, a_0)} \frac{d(s_0, a_0)}{d(s_T, a_T)}$$
(53)

With Lemma 1, we have,

$$\frac{1}{R(\pi)\big|_{s_0,a_0}} = \sup_{T \ge 0} \mathbb{E}_{s_T,a_T \sim P(s_T,a_T \mid \pi, s_0, a_0)} \frac{d(s_0, a_0)}{d(s_T, a_T)}$$
(54)

915 
$$((x_1)_{s_0,a_0} \to (x_1,x_1))$$

916  
917 
$$\leq^{(a)} \sup_{T \ge 0} \left( \mathbb{E}_{s_T, a_T \sim P(s_T, a_T | \pi, s_0, a_0)} \frac{\|\mathcal{T}^{\pi} f - \mathcal{T}^{\pi} f\|(s_T, a_T)}{\|\mathcal{T}^{\pi} f - \hat{\mathcal{T}}^{\pi} f\|(s_0, a_0)} \right)$$
(55)

The inequality (a) holds because of  $d(s_0, a_0) \cdot \zeta_0 \leq \sup_T \mathbb{E}_{s_T, a_T} d(s_T, a_T) \cdot \zeta_T$ . Then following Lemma ??, we have,

$$\sup_{T \ge 0} \mathbb{E}_{s_T, a_T \sim P(s_T, a_T | \pi, s_0, a_0)} \frac{\|\mathcal{T}^\pi f - \hat{\mathcal{T}}^\pi f\|^2(s_T, a_T)}{\|\mathcal{T}^\pi f - \hat{\mathcal{T}}^\pi f\|^2(s_0, a_0)}$$
(56)

$$= \sup_{T \ge 0} \mathbb{E}_{s_T, a_T \sim P(s_T, a_T | \pi, s_0, a_0)} \left( \frac{\sum_{t=0}^{T-1} [\gamma^{t+1} \zeta_{t+1} - \gamma^t \zeta_t]}{\zeta_0 \cdot \gamma^T} + \frac{1}{\gamma^T} \right)^2$$
(57)

$$\leq \mathcal{O}\left(\sup_{\pi \in \Pi, T \geq 0} \mathbb{E}_{s_T, a_T \sim P(s_T, a_T | \pi, s_0, a_0)} (\sum_{t=0}^{T-1} \gamma^{t+1} \zeta_{t+1} - \gamma^t \zeta_t)^2\right)$$
(58)

Completing the proof.

932 Please note that Lemma 5 holds for any estimation value function f. And we can utilize the learned 933  $\hat{f}_k$  at the  $k^{th}$  iteration.

# Furthermore, by plugging Eq.(58) in Corollary 1, we would have the proof for Lemma ??. Then we give the proof for Theorem 2 in the main text.

**Theorem 2.** Given an MDP with max reward  $R_{max}$  and a dataset of size N. The dimension of state space is |S| and that of action space is |A|. Given (s, a) pair, we denote its data density over the dataset is d(s, a). Given an empirical Bellman operator  $\hat{\mathcal{T}}^{\Pi}$  and an arbitrary policy candidate set  $\Pi$ , where  $\hat{\mathcal{T}}^{\Pi}f(s,a) = r(s,a) + \gamma \mathbb{E}_{s' \sim \hat{P}(s'|s,a)} \max_{\pi \in \Pi} f(s', \pi)$ . Denote the learnt value function as  $f_k$ , with k iterations of  $\hat{f}_k = \hat{\mathcal{T}}^{\Pi}\hat{f}_{k-1}$ , and the true optimal value function as  $f^*$ . Then we have,

$$\|\hat{f}_k - f^*\|_d \le \frac{C^*}{1 - \gamma} \cdot \sup_{\pi \in \Pi} \sum_{s_0} d(s_0) \zeta(s_0, \pi) +$$
(59)

$$\mathcal{O}\left(\sup_{\substack{\pi\in\Pi,T\geq0\\s_{0},a_{0},d(s_{0},a_{0})>0}} \mathbb{E}_{s_{T},a_{T}\sim P(s_{T},a_{T}|\pi,s_{0},a_{0})} (\sum_{t=0}^{T-1} \gamma^{t+1} \zeta_{t+1} - \gamma^{t} \zeta_{t})^{2}\right)$$
(60)

where  $\zeta_t$  is the Bellman uncertainty at time step t, i.e.,  $\zeta_t = \|\mathcal{T}^\Pi \hat{f}_k - \hat{\mathcal{T}}^\Pi \hat{f}_k\|(s_t, a_t)$ . Proof of Theorem 2.

$$\|\hat{f}_k - f^*\|_d \le \|\hat{f}_k - \hat{f}^*\|_d + \|\hat{f}^* - f^*\|_d$$
(61)

where  $\hat{f}^*$  is the fixed point of  $\hat{\mathcal{T}}^{\Pi}$ . Then Lemma 5 bounds  $\|\hat{f}_k - \hat{f}^*\|_d$ , i.e.,

$$\|\hat{f}_{k} - \hat{f}^{*}\|_{d} \leq \mathcal{O}\left(\sup_{\substack{\pi \in \Pi, T \ge 0\\ s_{0}, a_{0}, d(s_{0}, a_{0}) > 0}} \mathbb{E}_{s_{T}, a_{T} \sim P(s_{T}, a_{T} | \pi, s_{0}, a_{0})} (\sum_{t=0}^{T-1} \gamma^{t+1} \zeta_{t+1} - \gamma^{t} \zeta_{t})^{2}\right)$$
(62)

On the other hand,

$$\|\hat{f}^* - f^*\|_d \le \|\hat{\mathcal{T}}^{\Pi}\hat{f}^* - \hat{\mathcal{T}}^{\Pi}f^*\|_d + \|\hat{\mathcal{T}}^{\Pi}f^* - \mathcal{T}^{\Pi}f^*\|_d$$
(63)

$$\leq \|\hat{\mathcal{T}}^{\Pi}f^{*} - \mathcal{T}^{\Pi}f^{*}\|_{d} + \gamma \|\hat{f}^{*} - f^{*}\|_{d}$$
(64)

$$\Rightarrow \|\hat{f}^* - f^*\|_d \le \frac{\|\hat{\mathcal{T}}^{\Pi} f^* - \mathcal{T}^{\Pi} f^*\|_d}{1 - \gamma}$$
(65)

Then due to the optimal coverage assumption that  $\sup_{s,a} \frac{\pi^*(a|s)}{\pi_{\beta}(a|s)} \leq C^*$ , and  $\pi^* \in \Pi$ , we have,

 $\|\hat{\mathcal{T}}^{\Pi}f^* - \mathcal{T}^{\Pi}f^*\|_d = \|\sum_{s'} (P(s'|s,a) - \hat{P}(s'|s,a)) \max_{\pi \in \Pi} f^*(s',\pi)\|_d$ (66)

$$= \|\sum_{s'} (P(s'|s,a) - \hat{P}(s'|s,a))f^*(s',\pi^*)\|_d$$
(67)

$$\leq C^* \sup_{\pi \in \Pi} \sum_{s_0, a_0} d(s_0, a_0) \mathbb{E}_{a \sim \pi(a|s_0)} \|\hat{\mathcal{T}}^{\Pi} f^* - \mathcal{T} f^* \| (s_0, a)$$
(68)

$$\leq C^* \sup_{\pi \in \Pi} \sum_{s_0, a_0} d(s_0, a_0) \mathbb{E}_{a \sim \pi(a|s_0)} \| \hat{\mathcal{T}}^{\Pi} \hat{f}_k - \mathcal{T} \hat{f}_k \| (s_0, a)$$
(69)

The last inequality holds because the assumption that the optimal policy is reliable, so its value function would have an ideally low uncertainty.

Therefore, we have,

$$\|\hat{f}_k - f^*\|_d \le \sup_{\pi \in \Pi} \sum_{s_0} d(s_0) \zeta(s_0, \pi) +$$
(70)

$$\mathcal{O}\left(\sup_{\substack{\pi \in \Pi, T \ge 0\\s_0, a_0, d(s_0, a_0) > 0}} \mathbb{E}_{s_T, a_T \sim P(s_T, a_T \mid \pi, s_0, a_0)} (\sum_{t=0}^{T-1} \gamma^{t+1} \zeta_{t+1} - \gamma^t \zeta_t)^2\right)$$
(71)

Completing the proof. 

Theorem 2 inspires us that when restricting the policy candidate set, it is essential not only to constrain the uncertainty of the new policy's actions, as in traditional pessimistic algorithms, but also to limit the growth tendency of the Bellman uncertainty caused by the new policy. By controlling this tendency to be as minimal as possible, even ensuring that the Bellman uncertainty monoton-ically decreases over time steps, we can guarantee that the learned value function exhibits better performance and consequently induces a policy with superior performance. 

**Proposition 3.** If the first term of Eq.(12) is bounded, i.e.,  $\forall \pi \in \Pi$ , we have  $\mathbb{E}_{d(s_0)} \zeta_{f_1}(s_0, \pi) \leq c$ , then we can bound the second term with one-step Lyapunov Uncertainty-penalization, i.e.,  $\forall s \sim D$ , 

$$\min_{\pi} [\gamma \mathbb{E}_{P(s'|s,a)} \zeta_{\hat{f}_k,t+1}(s',\pi) - \zeta_{\hat{f}_k,t+1}(s,\pi)]$$

$$\Rightarrow \min_{\pi} \mathbb{E}_{P(\tau_T|\pi,s_0=s)} (\sum_{t=0}^{T-1} [\gamma^{t+1} \zeta_{\hat{f}_k,t+1} - \gamma^t \zeta_{\hat{f}_k,t}]$$
(72)
(73)

(73)

Furthermore, if we assume the dataset fully covers dynamics modes, i.e.,  $\forall s, a \in \mathcal{D}, P(s'|s, a) \subseteq$  $\hat{P}(s'|s, a)$ , then the left part could be controlled by Lyapunov value estimation.

*Proof of Proposition 3.*  $\min_{\pi} \mathbb{E}_{a \sim \pi(a|s)} \zeta_f(s, a)$ ; This constrains the new policy would not generated OOD actions beyond the demonstration of the offline data, i.e.,  $supp(\pi(a|s)) \subset supp(\pi_{\beta}(a|s))$ . In previous works Wu et al. (2022); Mao et al. (2024), such supported constraint is often achieved in a data density based way as  $\min_{\pi} \sum_{a \notin \pi_{\beta}} \pi(a|s)$ . Then we will show that the current step's error controlling minimizes the upper bound of the above supported constraint, 

$$\min_{\pi} \sum_{a \notin \pi_{\beta}} \pi(a|s) \Leftrightarrow \max_{\pi} \mathbb{E}_{a \in \pi(a|s)} d(s, a)$$
(74)

$$\leq^{(a)} \max_{\pi} C \cdot \mathbb{E}_{a \in \pi(a|s)} \frac{1}{\zeta_f^2(s,a)} \tag{75}$$

$$\lim_{\pi} \mathbb{E}_{a \in \pi(a|s)} \zeta_f^2(s, a) \tag{76}$$

$$\Leftrightarrow^{(b)} \min_{\pi} \mathbb{E}_{a \in \pi(a|s)} \zeta_f(s, a) \tag{77}$$

The inequality (a) holds because of the Lemma 1 in Appendix A.2. The equivalence (b) holds because the Bellman uncertainty is always non-negative.

Then with this constraint, we can have the new policy would reject the actions that is not supported by the behavior policy. This means, the new policy would only select the data-supported actions. Then,

$$\min_{\pi} \mathbb{E}_{P(\tau_T \mid \pi, s_0 = s, a_0 = a)} \left( \sum_{t=0}^{T-1} [\gamma^{t+1} \zeta_{\hat{f}_k, t+1} - \gamma^t \zeta_{\hat{f}_k, t}] \right)$$
(78)

$$\leq^{(a)} \min_{\pi} \frac{1 - \gamma^{T}}{1 - \gamma} \max_{s_{t} \in P(s_{t}|\pi, s_{0}=s)} (\gamma \mathbb{E}_{P(s_{t+1}|s_{t}, \pi)} \zeta_{\hat{f}_{k}}(s_{t+1}, \pi) - \zeta_{\hat{f}_{k}}(s_{t}, \pi))$$
(79)

$$\leq^{(b)} \min_{\pi} \frac{1 - \gamma^{I}}{1 - \gamma} \max_{s_{t} \in P(s_{t} \mid \pi_{\beta}, s_{0} = s)} (\gamma \mathbb{E}_{P(s_{t+1} \mid s_{t}, \pi)} \zeta_{\hat{f}_{k}}(s_{t+1}, \pi) - \zeta_{\hat{f}_{k}}(s_{t}, \pi)) \tag{80}$$

$$\leq^{(c)} \min_{\pi} \frac{1 - \gamma^{T}}{1 - \gamma} \max_{s \in \mathcal{D}} (\gamma \mathbb{E}_{P(s'|s,\pi)} \zeta_{\hat{f}_{k}}(s',\pi) - \zeta_{\hat{f}_{k}}(s,\pi))$$
(81)

1041 The inequality (a) is obtained using the formula for the sum of a geometric series. The inequal-1042 ity (b) is due to  $\pi \subseteq \pi_{\beta}$ , where  $\pi_{\beta}$  is the behavior policy. Finally, inequality (c) holds because 1043  $P(s_t | \pi_{\beta}, s_0 = s) \subseteq \mathcal{D}$ . Then we have,  $\forall s \in \mathcal{D}$ ,

$$\min_{\pi} (\gamma \mathbb{E}_{P(s'|s,\pi)} \zeta_{\hat{f}_k}(s',\pi) - \zeta_{\hat{f}_k}(s,\pi)) \Rightarrow \min_{\pi} \mathbb{E}_{P(\tau_T|\pi,s_0=s)} (\sum_{t=0}^{T-1} [\gamma^{t+1} \zeta_{\hat{f}_k,t+1} - \gamma^t \zeta_{\hat{f}_k,t}]$$
(82)

Then we would bound the left part through the Lyapunov value estimation. Due to  $\pi(a|s) \propto \hat{f}_k(s,a)$ , then the penalization to the s, a pairs is equivalent to minimize the preference of  $\pi$  to the action a at state s. In this way, the Lyapunov Uncertainty-penalization could be converted to,

$$\min_{\pi} (\gamma \max_{s' \in \hat{P}(s'|s,\pi)} \zeta_{\hat{f}_k}(s',\pi) - \zeta_{\hat{f}_k}(s,\pi)) \ge \min_{\pi} (\gamma \max_{s' \in P(s'|s,\pi)} \zeta_{\hat{f}_k}(s',\pi) - \zeta_{\hat{f}_k}(s,\pi))$$
(83)

 $\geq \min_{\pi} (\gamma \mathbb{E}_{P(s'|s,\pi)} \zeta_{\hat{f}_k}(s',\pi) - \zeta_{\hat{f}_k}(s,\pi))$ 

(84)

1054 1055 1056

1058

1061

1062

1052 1053

1039 1040

The first inequality holds because of  $P(s'|s, a) \subseteq \hat{P}(s'|s, a)$ . Completing the proof.

### **B** ADDITIONAL EXPERIMENTAL DETAILS

### B.1 TESTING ON ENVIRONMENTS WITH OOD OBSERVATIONS

1063 In this section, we conducted tests on benchmarks with OOD observations. This type of testing is 1064 primarily aimed at evaluating the generalization ability of conservative/pessimistic methods in offline reinforcement learning when there may be some noise present in the environment to make the agent deviate from the in-distributional regions. To address this challenge, we compared the perfor-1067 mance of the proposed LUC method with RORL (Yang et al., 2022) and OSR (Jiang et al., 2023) 1068 methods. These two methods are also improvements upon traditional conservative methods, incor-1069 porating enhancements such as smoothness constraints (RORL) and recovery constraints (OSR) to 1070 enhance the generalization performance of models on unknown states. We utilize the models trained on the 'medium' datasets of the three benchmarks - 'Halfcheetah', 'Hopper' and 'Walker2d', and 1071 three kinds of OOD noises - 'random', 'action\_diff' and 'min\_Q' like in (Yang et al., 2022). 1072

The results are shown in the Figure 5. We can observe that the proposed LUC method achieved good results on most benchmarks, particularly on the two types of adversarial attacks related to action differences. This may be because the learning approach in this paper has the weakest reliance on the behavioral policy, as the conservatism of LUC mainly stems from the evaluation of consequential reliability rather than a specific action distribution. In the context of min\_Q, this tests the robustness of these methods in maintaining the optimality of the value function. In this type of attack, LUC also exhibited superior performance compared to other methods, indicating that LUC has better robustness in dealing with the OOD situations than other methods.



Figure 5: Results of RORL, OSR and LUC on environments with OOD observations.

## B.2 CODE

We constructed our approach based on the RORL project available on GitHub<sup>3</sup>. The rationale behind selecting YangRui2015's project is as follows: 1) The RORL framework serves as a fundamental benchmark for conservative offline reinforcement learning, built on the PBRL implementation (Bai et al., 2022). 2) Implementing conservative Q functions is straightforward with the RORL framework. 3) As far as we are aware, the RORL framework stands out as the leading baseline in MuJoCo benchmarks. The code for our approach is included in the supplementary material.

1111 1112

1113

1100

1101 1102 1103

1104

## B.3 CONSTRUCTION OF OUT-OF-DISTRIBUTION MUJOCO BENCHMARKS

In this section, we introduce how to construct the testing environments for Out-of-distribution MuJoCo benchmarks in detail. First, we set three kinds of perturbations (different scales and intervals)
over three kinds of MuJoCo environments as shown in Table 5. The perturbation is randomly sampled from the Uniform distribution.

1118 1119

1120 1121 1122

Table 5	5: Param	eters for th	e constr	uction of	f Out-of-dis	tribution	n MuJoC	o benchma	rks.
	1	Halfcheetah	ı		Hopper			Walker2d	
	small	medium	large	small	medium	large	small	medium	large
scales	0.05	0.15	0.3	0.01	0.03	0.05	0.03	0.05	0.07
intervals	10	50	100	100	100	100	10	50	100

- 1123 1124 1125
- 1126

1127 1128 Then we visualize some of the perturbed situations, as is shown in Figure 6.

B.4 HYPERPARAMETERS OF LUC

1130 In Table 6 and Table 7, we give the hyperparameters used by LUC to generate Table 1 results. The  $\lambda_{LUC}$  is the weight of the reward shaping.

<sup>&</sup>lt;sup>3</sup>Project of RORL: https://github.com/YangRui2015/RORL

		1		1	1	
	0			1	Î	1
	5 4	-				
	T					
				1	1	
		-				
	U.	- 28				
Figure 6: Some visu	alized sample	s of perturbation	ons. The same	nples in blu	e box is the	normal
while the red box is	the correspond	ling perturbed	states. First	line is the sr	nall scales o	f perturb
second line is mediu	m; third line i	s large.				
	<					
Table	e 6: Hyperpara	imeters of LUC	in standard	l MuJoCo be	enchmarks.	
		Halfcheetah	Hopper	Walker2d		
	$\lambda_{LUC}$	0.1	0.1	0.05	_	
					_	
Table	e 7: Hyperpara	imeters of LUC	in adversa	rial attack be	enchmarks.	
		Halfcheetah	Hopper	Walker2d	_	
	$\lambda_{LUC}$	0.1	0.1	0.1		
					-	
		T T T	a			
B.5 NEURAL NET	WORK STRUC	TURES OF LU	С			
B.5 NEURAL NET In this section, we in	WORK STRUC	TURES OF LU	C networks we	e use in this	paper: polic	v netwo
B.5 NEURAL NET In this section, we ir Q network.	WORK STRUC	TURES OF LU	C networks we	e use in this	paper: polic	y netwo
B.5 NEURAL NET In this section, we ir Q network.	WORK STRUC	TURES OF LU ructure of the r	C networks we	e use in this	paper: polic	y netwo
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a	WORK STRUC ntroduce the st policy netwo and 'a_dim' is	TURES OF LU ructure of the r rk and Q netw the dimension	C networks we orks is as sl of actions.	e use in this nown in Tab 'h_dim' is th	paper: polic ble 8, where dimension	y networ 's_dim'
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our	TURES OF LU ructure of the r rk and Q netw the dimension experiments.	C networks we orks is as sl of actions. The policy r	e use in this nown in Tab 'h_dim' is th tetwork is a	paper: polic ble 8, where ble dimension Guassian po	y netwo 's_dim' of the l
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our en Q function 1	TURES OF LU ructure of the r rk and Q netw the dimension experiments. 7 networks and te	C networks we orks is as sl of actions. The policy r en target Q f	e use in this nown in Tab 'h_dim' is th tetwork is a function netw	paper: polic ble 8, where he dimensior Guassian po vorks.	y netwo 's_dim' of the l blicy and
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our en Q function r	TURES OF LU ructure of the r rk and Q netw the dimension experiments. ' networks and te	C networks we orks is as sl of actions. The policy r en target Q f	e use in this nown in Tab 'h_dim' is th tetwork is a function netw	paper: polic ble 8, where le dimensior Guassian po vorks.	y networ 's_dim' of the h licy and
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our on Q function 1 Table 8: The s	TURES OF LU ructure of the r rk and Q netw the dimension experiments. ' networks and te	C networks we orks is as sl of actions. The policy r en target Q f	e use in this nown in Tab 'h_dim' is th tetwork is a unction network	paper: polic ble 8, where le dimensior Guassian po vorks. works	y netwoi 's_dim' of the h olicy and
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our en Q function n Table 8: The s	TURES OF LU ructure of the r rk and Q netw the dimension experiments. ' networks and te tructure of the p	C networks we orks is as sl of actions. The policy r on target Q f	e use in this nown in Tab 'h_dim' is th thetwork is a function network and the Q net	paper: polic ble 8, where he dimension Guassian po vorks. works.	y netwo 's_dim' of the h blicy and
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our en Q function n Table 8: The s policy n	TURES OF LU ructure of the r rk and Q netw the dimension experiments. ' networks and te tructure of the p et	C networks we orks is as sl of actions. The policy r n target Q f policy net an Q net	e use in this nown in Tab 'h_dim' is th tetwork is a function network and the Q net	paper: polic ble 8, where le dimensior Guassian po vorks. works.	y netwo 's_dim' n of the h licy and
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our en Q function n Table 8: The s policy n Linear(s	TURES OF LU ructure of the r rk and Q netw the dimension experiments. T networks and te tructure of the p et adim, 256)	C networks we orks is as sl of actions. The policy r en target Q f policy net an Q net Linear(s_C	e use in this nown in Tab 'h_dim' is th tetwork is a unction netw nd the Q net	paper: polic ble 8, where le dimensior Guassian po vorks. works.	y netwo 's_dim' n of the l licy and
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our en Q function n Table 8: The s policy n Linear(s Relu() Linear(s	TURES OF LU ructure of the r rk and Q netw the dimension experiments. T networks and te tructure of the p et s_dim, 256)	C networks we orks is as sl of actions. The policy r on target Q f policy net an Q net Linear(s_c Relu() Linear(b)	e use in this nown in Tab 'h_dim' is th thetwork is a function network and the Q net lim, h_dim)	paper: polic ble 8, where le dimensior Guassian po vorks. works.	y netwo 's_dim' n of the l licy and
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	TWORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our en Q function n Table 8: The s policy n Linear(s Relu() Linear(h Relu()	TURES OF LU ructure of the r rk and Q netw the dimension experiments. T networks and te tructure of the p et s_dim, 256) n_dim, h_dim)	C networks we orks is as sl of actions. The policy r en target Q f policy net an Q net Linear(s_c Relu() Linear(h_c Relu()	e use in this nown in Tab 'h_dim' is th thetwork is a function network ad the Q net lim, h_dim)	paper: polic ole 8, where le dimension Guassian po vorks. works.	y netwo 's_dim' of the h olicy and
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our en Q function n Table 8: The s policy n Linear(s Relu() Linear(f Relu() Linear(f	TURES OF LU ructure of the n rk and Q netw the dimension experiments. ' networks and te tructure of the p et s_dim, 256) n_dim, h_dim) n_dim, a_dim)	C networks we orks is as sl of actions. The policy r en target Q f policy net ar Q net Linear(s_C Relu() Linear(h_C Relu()	e use in this hown in Tab 'h_dim' is th tetwork is a function netw hd the Q net lim, h_dim) lim, h_dim) lim, 1)	paper: polic ble 8, where dimension Guassian po vorks. works.	y netwo 's_dim' n of the h licy and
B.5 NEURAL NET In this section, we ir Q network. The structure of the dimension of states a layers, which is usua networks includes te	WORK STRUC ntroduce the st policy netwo and 'a_dim' is ally 256 in our en Q function n Table 8: The s policy n Linear(s Relu() Linear(f Relu() Linear(f	TURES OF LU ructure of the n rk and Q netw the dimension experiments. ' networks and te tructure of the p et s_dim, 256) n_dim, h_dim) n_dim, a_dim)	C networks we orks is as sl of actions. The policy r n target Q f policy net at Q net Linear(s_C Relu() Linear(h_C Relu()	e use in this hown in Tab 'h_dim' is th tetwork is a function network hd the Q net lim, h_dim) lim, h_dim) lim, 1)	paper: polic ble 8, where le dimension Guassian po vorks. 	y netwo 's_dim' of the f blicy and

We conducted all our experiments using a server equipped with one Intel Xeon Gold 5218 CPU, with 32 cores and 64 threads, and 256GB of DDR4 memory. We used a NVIDIA RTX3090 GPU

with 24GB of memory for our deep learning experiments. All computations were performed using
 Python 3.8 and the PyTorch deep learning framework.

1191 1192

1193

C DISCUSSION

1194 C.1 LIMITATIONS

In highly stochastic MDP environments with incomplete dataset coverage of transition outcomes, this paper's method may increase the likelihood of the agent straying from low-uncertainty regions, compromising decision reliability. Nevertheless, experimental results in Section 5.3 highlight the superior reliability of our LUC method over alternative approaches, showcasing enhanced generalization abilities in addressing previously unseen OOD scenarios.

- 1201 1202
- 1203

### 1204 C.2 DIFFERENCES BETWEEN LUC AND ROBUST RL 1205

Robust RL methods, although they seem to share a similar form to our method (Lyapunov uncertainty control), it is crucial to note the fundamental differences between them and some technical tips:

1209 1) Objective and motivation. Robust RL still falls short of addressing the safety issue we mentioned 1210 before. Specifically, the pessimism of Robust RL stems from penalizing with uncertainty in the 1211 outcome predictions of actions, to deal with the problem of distributional shift. However, if some 1212 behaviors that could lead to the deviation from the safe regions are supported well by dataset, the 1213 penalty loses its effectiveness due to that the uncertainty in outcome predictions would be tiny, thereby exacerbating the risk of entering high-uncertainty regions. While in our method, we aim to 1214 learn the policy to satisfy the Lyapunov reliable properties, of which we has shown the effectiveness 1215 in stable safe control in both theoretical and experimental ways. 1216

2) Penalization mechanism. Robust RL's penalty is reward-driven, aiming to maximize the expected cumulative return under the worst case to deal with transitioned distributional shift. On the other hand, our method places greater emphasis on the safety of the agent's decisions, aiming to minimize the risk of deviation from the safe regions under the worst case to achieve stable safety control for the agent. From this view, we can conclude that the Robust methods focus on the agent's generalization performance in solving offline RL problems, while our method prioritizes stable control of the agent to meet safety requirements in practical applications.

1224

### 1225 1226 C.3 DIFFERENCES BETWEEN LUC AND EDAC

There are several key differences between our method, LUC, and the EDAC (An et al., 2021) method in terms of motivation, implementation, and effectiveness, summarized as follows:

1) Motivation: The motivation behind the EDAC method lies in enhancing the sensitivity of Q-1230 ensembles to out-of-distribution (OOD) data. Therefore, the EDAC could be seen as a traditional 1231 pessimistic method. Pessimistic methods that focus only on having the agent behave like the demon-1232 strations present in the dataset - while helpful in the quantification for the OOD data, do not fully 1233 meet the safety requirements as mentioned before. This is because these methods often focus on 1234 an average effect, biasing the agent towards regions with higher data coverage or better model per-1235 formance, without strictly constraining the agent's activities within the safety region. This can be 1236 observed from Theorem 1 in this paper, where both rules defined in Definition 3 are essential for 1237 Lyapunov reliability: controlling the action uncertainty of current step and the growth of uncertainty 1238 in the future. Unfortunately, previous pessimistic methods only control the former, while ignoring 1239 the problem of uncertainty accumulation, failing to satisfy the Lyapunov reliability. While in our method, Q-ensembles serve as a tool to help us measure the safety of the model at certain data 1240 points, scoping a safe region for the agent to operate within. These two methods address completely 1241 different issues.

1242 2) Implementation: While the EDAC method, similar to traditional pessimistic methods, primarily penalizes the agent for selecting OOD actions, our method not only penalizes such actions but also considers whether choosing a specific action would lead to an increase in uncertainty, thereby preventing the agent from entering regions of high uncertainty (unreliability).

3) Effectiveness: Experimental results indicate that compared to traditional pessimistic methods like
 EDAC (such as PBRL, RORL, etc.), our method demonstrates greater reliability when facing highly stochastic real-world environments (see Section 5 of our paper).

- 1251 C.4 DIFFERENCES BETWEEN LUC AND CQL

When comparing our method to conservative approaches like CQL (Kumar et al., 2020), several key differences can be identified:

1255 1) OOD quantification: Conservative methods like CQL lack the ability to quantitatively measure the degree of out-of-distribution (OOD) data, leading to the rejection of all unseen data and reducing their generalization capabilities. In contrast, our method evaluates OOD data based on uncertainty, assessing the reliability of the model at that specific OOD data point, thus enabling it to generalize effectively on OOD data.

2) Safety assessment of consequences/long-term trends: CQL still focuses on penalizing OOD ac-tions without considering the outcomes' safety of these actions, making it challenging to ensure the agent's operation regions in practical deployment and potentially resulting in trajectory devia-tions (Zhang et al., 2022; Jiang et al., 2023; Kang et al., 2022). On the other hand, our method defines a closed region where the agent can operate stably, suitable for scenarios requiring higher reliability and safety in decision-making, such as autonomous driving and healthcare applications. This could also be seen in the experimental results in our paper, where we have compared our method LUC with CQL from various benchmarks.