

---

# RAMBO-RL: Robust Adversarial Model-Based Offline Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Offline reinforcement learning (RL) aims to find performant policies from logged  
2 data without further environment interaction. Model-based algorithms, which learn  
3 a model of the environment from the dataset and perform conservative policy opti-  
4 misation within that model, have emerged as a promising approach to this problem.  
5 In this work, we present Robust Adversarial Model-Based Offline RL (RAMBO), a  
6 novel approach to model-based offline RL. To enforce conservatism, we formulate  
7 the problem as a two-player zero sum game against an adversarial environment  
8 model. The model is trained to minimise the value function while still accurately  
9 predicting the transitions in the dataset, forcing the policy to act conservatively in  
10 areas not covered by the dataset. To approximately solve the two-player game, we  
11 alternate between optimising the policy and adversarially optimising the model.  
12 The problem formulation that we address is theoretically grounded, resulting in a  
13 probably approximately correct (PAC) performance guarantee and a pessimistic  
14 value function which lower bounds the value function in the true environment. We  
15 evaluate our approach on widely studied offline RL benchmarks, and demonstrate  
16 that it outperforms existing state-of-the-art baselines.

## 17 1 Introduction

18 Reinforcement learning (RL) [61] has achieved state-of-the-art performance on many sequential  
19 decision-making problems [40, 58, 44]. However, the need for extensive exploration prohibits the  
20 application of RL to many real world domains where such exploration is costly or dangerous. Offline  
21 RL [32, 34] overcomes this limitation by learning policies from static, pre-recorded datasets.

22 Online RL algorithms perform poorly in the offline setting due to the distribution shift between the  
23 state-action pairs in the dataset and those taken by the learnt policy. Thus, an important aspect of  
24 offline RL is to introduce *conservatism* to prevent the learnt policy from executing state-action pairs  
25 which are out of distribution. Model-free offline RL algorithms [28, 30, 72, 22, 26, 14] train a policy  
26 from only the data present in the fixed dataset, and incorporate conservatism either into the value  
27 function or by directly constraining the policy.

28 On the other hand, model-based offline RL algorithms [76, 75, 24, 63, 39] use the dataset to learn  
29 a model of the environment, and train a policy using additional synthetic data generated from that  
30 model. By training on additional synthetic data, model-based algorithms can potentially generalise  
31 better to states not present in the dataset, or to solving new tasks. Previous approaches to model-based  
32 offline RL incorporate conservatism by estimating the uncertainty in the model and applying reward  
33 penalties for state-action pairs that have high uncertainty [76, 24]. However, uncertainty estimation  
34 can be unreliable for neural network models [75, 36]. Like recent work [75], we propose an approach  
35 for offline model-based RL which *does not require uncertainty estimation*.

36 In this work we present Robust Adversarial Model-Based Offline (RAMBO) RL, a new algorithm for  
37 model-based offline RL. RAMBO incorporates conservatism by modifying the *transition dynamics* of  
38 the learnt Markov decision process (MDP) model in an adversarial manner. We formulate the problem  
39 of offline RL as a zero-sum game against an adversarial environment. To solve the resulting maximin  
40 optimisation problem, we alternate between optimising the agent and optimising the adversary in  
41 the style of Robust Adversarial RL (RARL) [49]. Unlike existing RARL approaches, our model-  
42 based approach forgoes the need to define and train an adversary policy, and instead only learns an  
43 adversarial model of the MDP. We train the agent policy with an actor-critic algorithm using synthetic  
44 data generated from the model in addition to data sampled from the dataset, similar to Dyna [60] and  
45 a number of recent methods [21, 76, 24, 75]. We update the environment model so that it reduces  
46 the value function for the agent policy, whilst still accurately predicting the transitions in the dataset.  
47 As a result, our approach introduces conservatism by generating *pessimistic synthetic transitions*  
48 for state-action pairs which are out-of-distribution. The theoretical formulation of offline RL that  
49 our algorithm addresses yields a PAC bound for the performance gap with respect to any policy  
50 covered by the dataset, and a pessimistic value function that lower bounds the value function in the  
51 true environment.

52 In summary, the main contributions of this work are:

- 53 • RAMBO, a novel and theoretically-grounded model-based offline RL algorithm which enforces  
54 conservatism by training an adversarial dynamics model.
- 55 • Adapting the Robust Adversarial RL approach to model-based offline RL by proposing a new  
56 formulation of RARL, where instead of defining and training an adversary policy, we directly train  
57 the model adversarially.

58 In our experiments we demonstrate that RAMBO outperforms current state-of-the-art algorithms  
59 on the D4RL benchmarks [13], and we provide ablation results which show that if the model is not  
60 trained adversarially, the performance of RAMBO degrades significantly.

## 61 2 Related Work

62 **Offline RL:** Offline RL addresses the problem of learning policies from fixed datasets, and has  
63 been applied to domains such as healthcare [42, 56], natural language processing [23, 22], and  
64 robotics [29, 37, 50]. Model-free offline RL algorithms do not require a learnt model. Approaches for  
65 model-free offline RL include importance sampling algorithms [35, 41], constraining the learnt policy  
66 to be similar to the behaviour policy [15, 27, 72, 22, 57, 14], incorporating conservatism into the value  
67 function during training [9, 26, 30, 73], using uncertainty quantification to generate more robust value  
68 estimates [1, 2, 28], or applying only a single iteration of policy iteration [7, 48]. In contrast, model-  
69 based approaches learn a model of the environment and generate synthetic data from that model [60]  
70 to optimise a policy using either planning [4] or RL algorithms [24, 76, 75]. By training a policy on  
71 additional synthetic data, model-based approaches have the potential for broader generalisation and  
72 for solving new tasks [6, 76]. A simple approach to ensuring conservatism is to constrain the policy  
73 to be similar to the behaviour policy in the same fashion as some model-free approaches [8, 39, 63].  
74 Another approach is to apply reward penalties for executing state-action pairs with high uncertainty  
75 in the environment model [24, 74, 76]. However, this requires explicit uncertainty estimates which  
76 may be unreliable for neural network models [16, 36, 45, 75]. COMBO [75] obviates the need for  
77 uncertainty estimation in model-based offline RL by adapting model-free techniques [30] to regularise  
78 the value function for out-of-distribution samples. Like COMBO, our approach does not require  
79 uncertainty estimation.

80 Most approaches to model-based offline RL use maximum likelihood estimates (MLE) of the MDP  
81 trained using standard supervised learning [4, 39, 63, 76, 75]. However, other methods have been  
82 proposed to learn models which are more suitable for offline policy optimisation. One approach  
83 is to reweight the loss function to ensure the model is accurate under the state-action distribution  
84 generated by the policy [33, 52, 19]. A *reverse* model is used by [67] to generate rollouts to goal  
85 states which are in-distribution. In contrast, our approach produces pessimistic synthetic transitions  
86 when out-of-distribution.

87 Most related to our work is a recent paper [66] which introduces the maximin formulation of offline  
88 RL that we address. This existing work motivates our approach theoretically by showing that the  
89 problem formulation obtains probably approximately correct (PAC) guarantees. However, [66]  
90 only addresses the theoretical aspects of the problem formulation and does not propose a practical

91 algorithm. In this work, we propose a practical RL algorithm to solve the maximin formulation of  
 92 model-based offline RL.

93 **Robust RL:** Algorithms for Robust MDPs [5, 20, 43, 53, 64, 70] find the policy with the best worst-  
 94 case performance over a set of possible MDPs. Typically, it is assumed that the uncertainty set of  
 95 MDPs is specified a priori. To eliminate the need to specify the set of possible MDPs, model-free  
 96 approaches to Robust MDPs [68, 54] instead assume that samples can be drawn from a misspecified  
 97 MDP which is similar to the true MDP. As our work addresses offline RL, we assume that we have a  
 98 fixed dataset from the true MDP.

99 Our approach is conceptually similar to Robust Adversarial RL (RARL) [49], a method proposed to  
 100 improve the robustness of RL policies in the *online* setting. RARL is posed as a two-player zero-sum  
 101 game where the agent plays against an adversary which perturbs the environment. Formulations  
 102 of model-free RARL differ in how they define the action space of the adversary. Options include  
 103 allowing the adversary to apply perturbation forces to the simulator [49], add noise to the agent’s  
 104 actions [65], or periodically take over control [47]. A model-based approach to RARL is proposed  
 105 in [12], which learns an optimistic and pessimistic model to encourage online exploration. However,  
 106 this existing approach requires uncertainty estimation as well as an adversarial policy to be learnt  
 107 in *addition to the model*. Our work follows the paradigm of RARL and alternates between agent and  
 108 adversarial updates in a maximin formulation. We adapt model-based RARL to the *offline* setting  
 109 and propose an alternative formulation: we eliminate the need to learn an adversary policy and  
 110 instead *directly modify the MDP model adversarially*. Related to our work is [38] which adversarially  
 111 modifies trajectories while training a policy. However, [38] assumes that the dynamics are known in  
 112 the form of a differential equation, whereas we assume that we only have access to sampled data.

### 113 3 Preliminaries

114 **MDPs and Offline RL:** An MDP is defined by the tuple,  $M = (S, A, T, R, \mu_0, \gamma)$ .  $S$  and  $A$   
 115 denote the state and action spaces respectively,  $R(s, a)$  is the reward function,  $T(s'|s, a)$  is the  
 116 transition function,  $\mu_0$  is the initial state distribution, and  $\gamma \in (0, 1)$  is the discount factor. In  
 117 this work we consider Markovian policies,  $\pi \in \Pi$ , which map from each state to a distribution  
 118 over actions. We denote the (improper) discounted state visitation distribution of a policy by  
 119  $d_M^\pi(s) := \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s | \pi, M)$ , where  $\Pr(s_t = s | \pi, M)$  is the probability of reaching state  $s$   
 120 at time  $t$  by executing policy  $\pi$  in  $M$ . The improper state-action visitation distribution is  $d_M^\pi(s, a) =$   
 121  $\pi(a|s) \cdot d_M^\pi(s)$ . We also denote the normalised state-action visitation distribution by  $\tilde{d}_M^\pi(s, a) =$   
 122  $(1 - \gamma) \cdot d_M^\pi(s, a)$ .

123 The value function,  $V_M^\pi(s)$ , represents the expected discounted return from executing  $\pi$  from state  
 124  $s$  in  $M$ :  $V_M^\pi(s) = \mathbb{E}_{\pi, M} [\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$ . We write  $V_M^\pi$  to indicate the value function under  
 125 the initial state distribution, i.e.  $V_M^\pi = \sum_{s \in S} \mu_0(s) V_M^\pi(s)$ . The standard objective for MDPs is to  
 126 find the policy which maximises  $V_M^\pi$ . The state-action value function,  $Q_M^\pi(s, a)$ , is the expected  
 127 discounted cumulative reward from taking action  $a$  at state  $s$  and then executing  $\pi$  thereafter.

128 In offline RL we only have access to a fixed dataset of transitions from the MDP,  $\mathcal{D} =$   
 129  $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^{|\mathcal{D}|}$ . The goal of offline RL is to find the best possible policy using the fixed dataset.

130 **Model-Based Offline RL Algorithms:** Model-based approaches to offline RL use a model of the  
 131 MDP to help train a policy. The dataset is used to learn a dynamics model,  $\hat{T}$ , which is typically trained  
 132 via maximum likelihood estimation:  $\min_{\hat{T}} \mathbb{E}_{(s, a, s') \sim \mathcal{D}} [-\log \hat{T}(s'|s, a)]$ . A model of the reward  
 133 function,  $\hat{R}(s, a)$ , can also be learnt if it is unknown. The estimated MDP,  $\hat{M} = (S, A, \hat{T}, \hat{R}, \mu_0, \gamma)$ ,  
 134 has the same state and action space as the true MDP but uses the learnt transition and reward functions.  
 135 Thereafter, any planning or RL algorithm can be used to recover optimal policy in the learnt model,  
 136  $\hat{\pi} = \arg \max_{\pi \in \Pi} V_{\hat{M}}^\pi$ .

137 Unfortunately, directly applying this approach to the offline RL setting does not perform well due  
 138 to distribution shift. In particular, if the dataset does not cover the entire state-action space, the  
 139 model will inevitably be inaccurate for some state-action pairs. Thus, naive policy optimisation on a  
 140 learnt model in the offline setting can result in *model exploitation* [21, 31, 52]. To mitigate this issue,  
 141 we propose the novel approach of enforcing conservatism by adversarially modifying the transition  
 142 dynamics of  $\hat{M}$ .

143 In line with existing works [76, 75, 8], we use model-based policy optimisation (MBPO) [21] to  
 144 learn the optimal policy for  $\widehat{M}$ . MBPO utilises a standard actor-critic RL algorithm. However, the  
 145 value function is trained using an augmented dataset  $\mathcal{D} \cup \mathcal{D}_{\widehat{M}}$ , where  $\mathcal{D}_{\widehat{M}}$  is synthetic data generated  
 146 by simulating rollouts in the learnt model. To generate the synthetic data, MBPO performs  $k$ -step  
 147 rollouts in  $\widehat{M}$  starting from states  $s \in \mathcal{D}$ , and adds this data to  $\mathcal{D}_{\widehat{M}}$ . To train the policy, minibatches  
 148 of data are drawn from  $\mathcal{D} \cup \mathcal{D}_{\widehat{M}}$ , where each datapoint is sampled from the real data,  $\mathcal{D}$ , with  
 149 probability  $f$ , and from  $\mathcal{D}_{\widehat{M}}$  with probability  $1 - f$ .

150 **Robust Adversarial Reinforcement Learning:** RARL addresses the problem of finding a robust  
 151 agent policy,  $\pi$ , in the online RL setting by posing the problem as a two-player zero sum game against  
 152 adversary policy,  $\bar{\pi}$ :

$$\pi = \arg \max_{\pi \in \Pi} \min_{\bar{\pi} \in \bar{\Pi}} V_M^{\pi, \bar{\pi}} \quad (1)$$

153 where  $V_M^{\pi, \bar{\pi}}$  is the expected value from executing  $\pi$  and  $\bar{\pi}$  in environment  $M$ . Different approaches  
 154 define the action space for  $\bar{\pi}$  in different ways, as discussed in Section 2. For a scalable approximation  
 155 to the optimisation problem in Equation 1, algorithms for RARL alternate between applying steps of  
 156 stochastic gradient ascent to the agent’s policy to increase the expected value, and stochastic gradient  
 157 descent to the adversary’s policy to decrease the expected value. In our work, we follow the RARL  
 158 paradigm of alternating between agent and adversarial updates and adapt it to the model-based offline  
 159 setting. Instead of defining a separate adversary policy, we treat the *model itself* as the policy to be  
 160 adversarially trained.

## 161 4 Problem Formulation

162 For the sake of generality, we assume that both the transition function and reward function are  
 163 unknown. Hereafter, we write  $\widehat{T}$  to denote both the learnt dynamics and reward function, where  
 164  $\widehat{T}(s', r | s, a)$  is the probability of receiving reward  $r$  and transitioning to  $s'$  after executing  $(s, a)$ . We  
 165 address the maximin formulation of offline RL recently proposed by [66]:

166 **Problem 1.** For some dataset,  $\mathcal{D}$ , and some fixed constant  $\xi > 0$ , find the policy  $\pi$  defined by

$$\pi = \arg \max_{\pi \in \Pi} \min_{\widehat{T} \in \mathcal{M}_{\mathcal{D}}} V_{\widehat{T}}^{\pi}, \text{ where} \quad (2)$$

$$\mathcal{M}_{\mathcal{D}} = \left\{ \widehat{T} \mid \mathbb{E}_{\mathcal{D}} [\text{TV}(\widehat{T}_{\text{MLE}}(\cdot | s, a), \widehat{T}(\cdot | s, a))^2] \leq \xi \right\}, \quad (3)$$

168 where  $\text{TV}(P_1, P_2)$  is the total variation distance between distributions  $P_1$  and  $P_2$ , and  $\widehat{T}_{\text{MLE}}$  denotes  
 169 the maximum likelihood estimate of the MDP given the offline dataset,  $\mathcal{D}$ .

170 Thus, the set defined in Equation 3 contains MDPs which are similar to the maximum likelihood  
 171 estimate under state-action pairs in  $\mathcal{D}$ . However, because the expectation in Equation 3 is taken  
 172 under  $\mathcal{D}$ , there is no restriction on  $\widehat{T}$  for regions of the state-action space not covered by  $\mathcal{D}$ . We  
 173 present a brief overview of the theoretical guarantees from [66] in the following subsection.

174 **Remark 1.** Note that Problem 1 differs from the pessimistic MDP formulations introduced by  
 175 MOPO [76] and MOREL [24]. Problem 1 considers the worst-case transition dynamics, while the  
 176 pessimistic MDPs constructed by MOPO and MOREL only modify the reward function by applying  
 177 reward penalties for state-action pairs with high uncertainty.

### 178 4.1 Theoretical Motivation

179 The theoretical analysis from [66] shows that solving Problem 1 outputs a policy that with high  
 180 probability is approximately as good as any policy with a state-action distribution that is covered by  
 181 the dataset. This is formally stated in the following theorem.

**Theorem 1** (PAC guarantee from [66], Theorem 5). Denote the true MDP transition function by  $T$ ,  
 and let  $\mathcal{M}$  denote a hypothesis class of MDP models such that  $T \in \mathcal{M}$ . Let  $\pi$  denote the solution to  
 Problem 1 for dataset  $\mathcal{D}$ . Then with probability  $1 - \delta$  for any policy,  $\pi^* \in \Pi$ , we have

$$V_T^{\pi^*} - V_T^{\pi} \leq (1 - \gamma)^{-2} c_1 \sqrt{C_{\pi^*}} \sqrt{G_{\mathcal{M}_1} + G_{\mathcal{M}_2} + \xi_n^2 + \frac{\ln(c/\delta)}{|\mathcal{D}|}}, \text{ where}$$

$$C_{\pi^*} = \max_{T' \in \mathcal{M}} \frac{\mathbb{E}_{(s,a) \sim \tilde{d}_{T'}^*} [\text{TV}(T'(\cdot|s,a), T(\cdot|s,a))^2]}{\mathbb{E}_{(s,a) \sim \rho} [\text{TV}(T'(\cdot|s,a), T(\cdot|s,a))^2]},$$

182 where  $\rho$  is the state-action distribution from which  $\mathcal{D}$  was sampled, and  $c$  and  $c_1$  are universal  
183 constants. We refer the reader to Appendix A of [66] for the definitions of  $G_{\mathcal{M}_1}$ ,  $G_{\mathcal{M}_2}$ , and  $\xi_n$ .

184 The quantity  $C_{\pi^*}$  is upper bounded by the maximum density ratio between the comparator policy,  
185  $\pi^*$ , and the offline distribution, i.e.  $C_{\pi^*} \leq \max_{(s,a)} \tilde{d}_{T'}^*(s,a)/\rho(s,a)$ . It represents the discrepancy  
186 between the distribution of data in the dataset compared to the visitation distribution of policy  $\pi^*$ .  
187 Theorem 1 shows that if we find a policy by solving Problem 1, the performance gap of that policy is  
188 bounded with respect to any other policy  $\pi^*$  which has a state-action distribution which is covered by  
189 the dataset.

190 Furthermore, the value function under the worst-case model in the set defined by Problem 1 is a lower  
191 bound on the value function in the true environment, as stated by Proposition 1.

**Proposition 1** (Pessimistic value function). *Let  $T$  denote the true transition function for some MDP, and let  $\mathcal{M}_{\mathcal{D}}$  be the set of MDP models defined in Equation 3. Then for any policy  $\pi$ , with probability  $1 - \delta$  we have that*

$$\min_{\hat{T} \in \mathcal{M}_{\mathcal{D}}} V_{\hat{T}}^{\pi} \leq V_T^{\pi}.$$

192 Proposition 1 follows from the fact that  $T \in \mathcal{M}_{\mathcal{D}}$  with high probability, which is proven in [66]  
193 (Appendix E.2). Proposition 1 shows that we can expect the performance of any policy in the true  
194 MDP to be at least as good as the value in the worst-case model defined in Problem 1.

195 While [66] provides the theoretical motivation for solving Problem 1, it does not propose a practical  
196 algorithm. In this work, we focus on developing a practical approach to solving Problem 1.

## 197 5 RAMBO-RL

198 In this section, we present Robust Adversarial Model-Based Offline RL (RAMBO), our algorithm  
199 for solving Problem 1. The main difficulty with solving Problem 1 is that it is unclear how to find  
200 the worst-case MDP in the set defined in Equation 3. To arrive at a scalable solution, we propose a  
201 novel approach which is in the spirit of RARL. We alternate between optimising the agent policy to  
202 increase the expected value, and adversarially optimising the model to decrease the expected value.  
203 In this section, we first describe how we compute the gradient to adversarially train the model. Then,  
204 we discuss how to ensure that the model remains approximately within the constraint set defined in  
205 Problem 1. Finally, we present our overall algorithm.

### 206 5.1 Model Gradient

207 We propose a policy gradient-inspired approach to adversarially optimise the model. Typically,  
208 policy gradient algorithms are used to modify the distribution over actions taken by a policy at each  
209 state [61, 62, 71]. In contrast, the update that we propose modifies the likelihood of the successor  
210 states and rewards in the MDP model.

211 We assume that the MDP model is defined by parameters,  $\phi$ , and we write  $\hat{T}_{\phi}$  to indicate this. We  
212 denote by  $V_{\phi}^{\pi}$  the value function for policy  $\pi$  in model  $\hat{T}_{\phi}$ . To approximately find  $\min_{\hat{T}_{\phi} \in \mathcal{M}_{\mathcal{D}}} V_{\phi}^{\pi}$  as  
213 required by Problem 1 via gradient descent, we wish to compute the gradient of the model parameters  
214 which reduces the value of the policy within the model, i.e.  $\nabla_{\phi} V_{\phi}^{\pi}$ .

215 **Proposition 2** (Model Gradient). *Let  $\phi$  denote the parameters of a parametric MDP model  $\hat{T}_{\phi}$ , and  
216 let  $V_{\phi}^{\pi}$  denote the value function for policy  $\pi$  in  $\hat{T}_{\phi}$ . Then:*

$$\nabla_{\phi} V_{\phi}^{\pi} = \mathbb{E}_{s \sim d_{\phi}^{\pi}, a \sim \pi, (s', r) \sim \hat{T}_{\phi}} [(r + \gamma V_{\phi}^{\pi}(s')) \cdot \nabla_{\phi} \log \hat{T}_{\phi}(s', r | s, a)] \quad (4)$$

217 The proof of Proposition 2 is given in Appendix A. We can subtract the baseline  $Q_{\phi}^{\pi}(s, a)$  without  
218 biasing the gradient estimate (see Appendix A.1 for details):

$$\nabla_{\phi} V_{\phi}^{\pi} = \mathbb{E}_{s \sim d_{\phi}^{\pi}, a \sim \pi, (s', r) \sim \hat{T}_{\phi}} [(r + \gamma V_{\phi}^{\pi}(s') - Q_{\phi}^{\pi}(s, a)) \cdot \nabla_{\phi} \log \hat{T}_{\phi}(s', r | s, a)] \quad (5)$$

219 The Model Gradient differs from the standard policy gradient in that a) it is used to update the  
 220 likelihood of successor states in the model, rather than actions in a policy, and b) the “advantage”  
 221 term  $r + \gamma V_\phi^\pi(s') - Q_\phi^\pi(s, a)$  compares the utility of receiving reward  $r$  and transitioning to  $s'$  to the  
 222 expected value of state action pair  $(s, a)$ . In contrast, the standard advantage term,  $Q(s, a) - V(s)$  [55],  
 223 compares the utility of executing state-action pair  $(s, a)$  to the expected value at  $s$ . To estimate  $V_\phi^\pi$   
 224 and  $Q_\phi^\pi$ , we use the critic learnt by the actor-critic algorithm used for policy optimisation. Thus, we  
 225 use the critic *both* for training the policy and adversarially training the model.

226 **Remark 2.** The Model Gradient can be thought of as a specific instantiation of the policy gradient if  
 227 we view  $\hat{T}_\phi$  as a adversarial “policy” on an augmented MDP,  $M^+$ , i.e.  $\hat{T}_\phi : S^+ \rightarrow \text{Dist}(A^+)$ . The  
 228 augmented state space,  $S^+ = S \times A$  includes the state in the original MDP augmented by the action  
 229 taken by the agent. The augmented action space consists of the reward applied and the successor  
 230 state,  $A^+ = S \times [R_{\min}, R_{\max}]$ . Thus, we can think of this approach as an instantiation of RARL in  
 231 which the adversary policy ( $\bar{\pi}$  in Equation 1) that we train is the *model itself*.

232 **Remark 3.** In principle, we could have instead adapted backpropagation through time (BPTT) [11,  
 233 69] to adversarially train the model. However, BPTT is known to have issues with numerical  
 234 instability [11].

## 235 5.2 Adversarial Model Training

236 If we were to update the model using Equation 5 alone, this would allow the model to be modified  
 237 arbitrarily such that the value function in the model is reduced. However, the set of plausible MDPs  
 238 given by Equation 3 states that over the dataset,  $\mathcal{D}$ , the model  $\hat{T}_\phi$  should be close to the maximum  
 239 likelihood estimate,  $\hat{T}_{\text{MLE}}$ . Specifically, in the inner optimisation of Problem 1 we wish to find a  
 240 solution to the constrained optimisation problem

$$\min_{\hat{T}_\phi} V_\phi^\pi, \quad \text{s.t.} \quad \mathbb{E}_{\mathcal{D}} [\text{TV}(\hat{T}_{\text{MLE}}(\cdot|s, a), \hat{T}_\phi(\cdot|s, a))^2] \leq \xi. \quad (6)$$

241 The Lagrangian relaxation leads to the unconstrained problem

$$\max_{\lambda \geq 0} \min_{\hat{T}_\phi} \left( L(\hat{T}, \lambda) := V_\phi^\pi + \lambda (\mathbb{E}_{\mathcal{D}} [\text{TV}(\hat{T}_{\text{MLE}}(\cdot|s, a), \hat{T}_\phi(\cdot|s, a))^2] - \xi) \right), \quad (7)$$

242 where  $\lambda$  is the Lagrange multiplier. Rather than optimising the Lagrange multiplier, we find that  
 243 in practice fixing  $\lambda$  to apply a constant weighting between the two terms works well with minimal  
 244 tuning. To facilitate easier tuning of the learning rate, in our implementation we apply the weighting  
 245 constant to the value function term rather than the model term, which is equivalent up to a scaling  
 246 factor. This leads to

$$\min_{\hat{T}_\phi} \left( \lambda V_\phi^\pi + \mathbb{E}_{\mathcal{D}} [\text{TV}(\hat{T}_{\text{MLE}}(\cdot|s, a), \hat{T}_\phi(\cdot|s, a))^2] \right). \quad (8)$$

247 Rather than explicitly minimising the TV distance between the model and MLE model, we apply the  
 248 standard MLE loss for the model term. This leads to the final loss function:

$$\mathcal{L}_\phi = \lambda V_\phi^\pi - \mathbb{E}_{(s, a, s') \sim \mathcal{D}} [\log \hat{T}_\phi(s'|s, a)]. \quad (9)$$

249 By minimising the loss function in Equation 9, the model is trained to a) predict the transitions within  
 250 the dataset, and b) reduce the value function of the policy, with  $\lambda$  determining the tradeoff between  
 251 these two objectives. Choosing  $\lambda$  to be small ensures that the MLE term dominates for transitions  
 252 within  $\mathcal{D}$ , ensuring that the model fits the dataset accurately. Because the MLE term is only computed  
 253 over  $\mathcal{D}$ , the value function term dominates outside of the dataset meaning that the model is modified  
 254 adversarially for transitions outside of the dataset.

255 To estimate the gradient of the loss function in Equation 9 for stochastic gradient descent, we sample  
 256 a minibatch of transitions from  $\mathcal{D}$  to estimate the MLE term. The gradient for the value function term  
 257 is computed using the Model Gradient. The transitions used to compute the Model Gradient term  
 258 must be sampled under the current policy and model (Equation 5). Therefore, to estimate the Model  
 259 Gradient term, we generate a minibatch of transitions by simulating the current policy in  $\hat{T}_\phi$ .

260 Like previous works [10, 76, 75, 8], we represent the dynamics model using an ensemble of neural  
 261 networks. Each neural network produces a Gaussian distribution over the next state and reward:

262  $\hat{T}_\phi(s', r|s, a) = \mathcal{N}(\mu_\phi(s, a), \Sigma_\phi(s, a))$ . A visualisation of the result of adversarially training the  
263 dynamics model can be found in Appendix C.4.

264 **Normalisation** The composite loss function in Equation 9 comprises two terms which may have  
265 different magnitudes across domains depending on the scale of the states and rewards. To enable  
266 easier tuning of the adversarial loss weighting,  $\lambda$ , across different domains we perform the following  
267 normalisation procedure. Prior to training, we normalise the states in  $\mathcal{D}$  in the manner proposed  
268 by [14], by subtracting the mean and dividing by the standard deviation of each state dimension in  
269 the dataset. Additionally, when computing the gradient in Equation 5 we normalise the advantage  
270 terms,  $r + \gamma V_\phi^\pi(s') - Q_\phi^\pi(s, a)$ , according to the mean and standard deviation across each minibatch.  
271 Advantage normalisation is already common practice in policy gradient RL implementations [3, 51].

## 272 Algorithm

273 We are now ready to present our overall approach in Algorithm 1. The first step of RAMBO is to  
274 pretrain the environment dynamics model using standard MLE (Line 1). Thereafter, the algorithm  
275 follows the format of RARL. At each iteration, we apply gradient updates to the agent to increase the  
276 expected value, followed by gradient updates to the model to decrease the expected value.

277 Prior to each agent update, we generate synthetic  $k$ -step rollouts starting from states in  $\mathcal{D}$  by simulating  
278 rollouts in the current MDP model  $\hat{T}_\phi$ . This data is added to the synthetic dataset  $\mathcal{D}_{\hat{T}_\phi}$  (Line 3).  
279 Following previous approaches [76, 75, 21] we store data in  $\mathcal{D}_{\hat{T}_\phi}$  in a first in, first out manner so that  
280 only data generated from recent iterations is stored in  $\mathcal{D}_{\hat{T}_\phi}$ . The agent’s policy and value functions  
281 are trained using an off-policy actor critic algorithm using samples from  $\mathcal{D} \cup \mathcal{D}_{\hat{T}_\phi}$  (Line 4). In our  
282 implementation, we use soft actor-critic (SAC) [18] for agent training.

283 To update the model to minimise the loss in Equation 9 in Line 5, we sample data from  $\mathcal{D}$  to estimate  
284 the gradient for the MLE component. To compute the adversarial component, we generate samples by  
285 simulating the current policy and model and utilise the value function learnt by the agent to compute  
286 the gradient according to Equation 5.

---

### Algorithm 1 RAMBO-RL

---

**Require:** Normalised dataset,  $\mathcal{D}$ ;

- 1:  $\hat{T}_\phi \leftarrow$  MLE dynamics model.
  - 2: **for**  $i = 1, 2, \dots, n_{\text{iter}}$  **do**
  - 3:   Generate synthetic  $k$ -step rollouts. Add transition data to  $\mathcal{D}_{\hat{T}_\phi}$ .
  - 4:   *Agent update:* Update  $\pi$  and  $Q_\phi^\pi$  with an actor critic algorithm, using samples from  $\mathcal{D} \cup \mathcal{D}_{\hat{T}_\phi}$ .
  - 5:   *Adversarial model update:* Update  $\hat{T}_\phi$  according to Eq. 9, using samples from  $\mathcal{D}$  for the MLE component and the current critic  $Q_\phi^\pi$  and synthetic data sampled from  $\pi$  and  $\hat{T}_\phi$  for the adversarial component.
- 

## 287 6 Experiments

288 In our experiments, we examine how well RAMBO performs compared to state-of-the-art baselines,  
289 and explore the importance of adversarial training to the performance of the algorithm. The code  
290 for our experiments is included in the supplementary material. We evaluate our approach on the  
291 following domains from the D4RL benchmarks [13].

292 **MuJoCo** There are three different environments representing different robots (*HalfCheetah*, *Hopper*,  
293 *Walker2D*), each with 4 datasets (*Random*, *Medium*, *Medium-Replay*, *Medium-Expert*). *Random*  
294 contains transitions collected by a random policy. *Medium* consists of the replay buffer generated by an  
295 early-stopped SAC policy. *Medium-Replay* consists of the replay buffer generated while training the  
296 *Medium* policy. The *Medium-Expert* dataset contains a mixture of suboptimal and expert data.

297 **AntMaze** The agent controls a robot and navigates to reach a goal, receiving a sparse reward only if  
298 the goal is reached. There are three different layouts of maze (*Umaze*, *Medium*, *Large*), and different  
299 dataset types (*Fixed*, *Play*, *Diverse*) which differ in terms of the variety of start and goal locations  
300 used to collect the dataset.

301 **Single Transition Example** This domain has a one-dimensional state and action space, and several  
302 distinct regions of the action space are covered by the dataset. We use this domain to illustrate the  
303 difference between RAMBO and the most similar prior algorithm, COMBO. See Appendix C.1 for a  
304 more detailed description.

305 **Hyperparameter Selection and Evaluation** The base hyperparameters that we use for RAMBO  
306 mostly follow those used in SAC [18] and COMBO [75]. We find that the performance of RAMBO  
307 is sensitive to the choice of rollout length,  $k$ , consistent with findings in previous works [21, 36]. The  
308 other critical parameter for RAMBO is the choice of the adversarial weighting,  $\lambda$ .

309 For each dataset, we choose the rollout length and the adversarial weighting from one of three possible  
310 configurations:  $(k, \lambda) \in \{(2, 3e-4), (5, 3e-4), (5, 0)\}$ . We included  $(k, \lambda) = (5, 0)$  as we found that  
311 an adversarial weighting of 0 worked well for some datasets. For all MuJoCo datasets we performed  
312 model rollouts using the current policy, but for some AntMaze datasets we used a random rollout  
313 policy as we found that this performed better. Further details about the hyperparameters are in  
314 Appendix B.

315 To evaluate RAMBO, we ran each of the three hyperparameter configurations for five seeds each, and  
316 report the best performance across the three configurations. Thus, our evaluation of RAMBO utilises  
317 limited online hyperparameter tuning to choose hyperparameters for each dataset. Note that this is the  
318 most common practice among existing model-based offline RL algorithms [8, 24, 36, 39, 76]. The  
319 performance obtained for each of the hyperparameter configurations is included in Appendix C.2.

320 Offline hyperparameter selection is an important topic in offline RL [46, 77]. This is because  
321 in some applications selecting hyperparameters based on online performance may be infeasible.  
322 Therefore, we present additional results in Appendix C.3 where we select between the three choices  
323 of hyperparameters offline using a simple heuristic based on the magnitude and stability of the  
324  $Q$ -values during training.

325 **Computational Resources** During our experiments, each run of RAMBO has access to 2 CPUs  
326 of an Intel Xeon Platinum 8259CL processor at 3.1GHz, and half of an Nvidia T4 GPU. With this  
327 hardware, each run of RAMBO takes 24-30 hours.

328 **Baselines** We compare RAMBO against state-of-the-art model-based (COMBO [75], MOReL [24],  
329 and MOPO [76]) and model-free (CQL [30], IQL [27], and TD3+BC [14]) offline RL algorithms. To  
330 facilitate a fair comparison, we provide results for all algorithms for the MuJoCo-v2 D4RL datasets  
331 which contain more performant data than v0 [6] (details in Appendix B.6). We report results for the  
332 AntMaze-v0 datasets.

333 **Results** The results in Table 1 show that RAMBO outperforms existing state-of-the-art methods  
334 on the MuJoCo locomotion domains. Relative to existing algorithms, RAMBO obtains the best  
335 performance on the Medium-Replay datasets, and performs less well on the Medium-Expert datasets.  
336 For the high-quality data in the Medium-Expert datasets, simpler approaches such as performing  
337 behaviour cloning on the best 10% of trajectories can be used to achieve stronger performance than  
338 offline RL methods [27]. Therefore, the suboptimal performance of RAMBO on the Medium-Expert  
339 datasets is less of a concern, as applying offline RL algorithms may not be the most suitable approach  
340 for these datasets.

341 For AntMaze, RAMBO performs considerably less well than the model-free baselines. This echoes  
342 the findings of [67] that model-based approaches that use forward rollouts struggle to perform well in  
343 the AntMaze domains, potentially because model-based algorithms are too aggressive and collide  
344 with walls. The authors of [67] show that *reverse* rollouts can lead to stronger performance for  
345 model-based methods in these domains. In future work, we wish to investigate whether combining  
346 RAMBO with reverse rollouts can lead to stronger performance on these domains.

347 In Table 2 we present ablations comparing the performance of RAMBO to the same approach without  
348 any adversarial updates to the model. These results demonstrate that overall performance degrades  
349 if the adversarial training is removed. This parallels previous findings that conservatism is crucial  
350 to obtaining strong performance in offline RL. Interestingly however, for some specific datasets  
351 we obtain the best performance with no adversarial training (Appendix C.2). This suggests that a  
352 potential direction for future work could be trying to identify which types of problems do not require  
353 conservatism for a successful policy to be trained offline with model-based RL.

Table 1: Results for the D4RL benchmark using the normalisation procedure proposed by [13]. We report the normalised performance during the last 10 iterations of training averaged over 5 seeds.  $\pm$  captures the standard deviation over seeds. Highlighted numbers indicate results within 1% of the most performant algorithm. \* indicates the total without random datasets.

		Ours	Model-based baselines			Model-free baselines			
		RAMBO	COMBO	MOPO	MOReL	CQL	IQL	TD3+BC	BC
Random	HalfCheetah	39.5 $\pm$ 3.5	38.8	35.4	25.6	19.6	-	11.0	2.1
	Hopper	25.4 $\pm$ 7.5	17.9	4.1	53.6	6.7	-	8.5	9.8
	Walker2D	0.0 $\pm$ 0.3	7.0	4.2	37.3	2.4	-	1.6	1.6
Medium	HalfCheetah	77.9 $\pm$ 4.0	54.2	69.5	42.1	49.0	47.4	48.3	36.1
	Hopper	87.0 $\pm$ 15.4	94.9	48.0	95.4	66.6	66.3	59.3	29.0
	Walker2D	84.9 $\pm$ 2.6	75.5	-0.2	77.8	83.8	78.3	83.7	6.6
Medium Replay	HalfCheetah	68.7 $\pm$ 5.3	55.1	68.2	40.2	47.1	44.2	44.6	38.4
	Hopper	99.5 $\pm$ 4.8	73.1	39.1	93.6	97.0	94.7	60.9	11.8
	Walker2D	89.2 $\pm$ 6.7	56.0	69.4	49.8	88.2	73.9	81.8	11.3
Medium Expert	HalfCheetah	95.4 $\pm$ 5.4	90.0	72.7	53.3	90.8	86.7	90.7	35.8
	Hopper	88.2 $\pm$ 20.5	111.1	3.3	108.7	106.8	91.5	98.0	111.9
	Walker2D	56.7 $\pm$ 39.0	96.1	-0.3	95.6	109.4	109.6	110.1	6.4
<b>MuJoCo-v2 Total:</b>		812.4 $\pm$ 47.8	769.7	413.4	773.0	767.4	692.6*	698.5	300.8
AntMaze	Umaze	25.0 $\pm$ 12.0	80.3	0.0	0.0	74.0	87.5	78.6	65.0
	Medium-Play	16.4 $\pm$ 17.9	0.0	0.0	0.0	61.2	71.2	3.0	0.0
	Large-Play	0.0 $\pm$ 0.0	0.0	0.0	0.0	15.8	39.6	0.0	0.0
	Umaze-Diverse	0.0 $\pm$ 0.0	57.3	0.0	0.0	84.0	62.2	71.4	55.0
	Medium-Diverse	23.2 $\pm$ 14.2	0.0	0.0	0.0	53.7	70.0	10.6	0.0
	Large-Diverse	2.4 $\pm$ 3.3	0.0	0.0	0.0	14.9	47.5	0.2	0.0
<b>AntMaze-v0 Total:</b>		67.0 $\pm$ 14.9	137.6	0.0	0.0	303.6	378.0	163.8	120.0

Table 2: Ablation of the adversarial updates for RAMBO. These results use the same rollout length for each dataset as RAMBO but with no adversarial updates (i.e.  $\lambda = 0$ ). The scores are averaged over 5 seeds. To produce these results, we gave runs where the value function diverged a score of 0.

RAMBO (No Adversarial Training)	
<b>MuJoCo-v2 Total:</b>	671.5 $\pm$ 43.5   <b>AntMaze-v0 Total:</b> 45.8 $\pm$ 21.8

Table 3: Comparison between RAMBO and COMBO for the Single Transition Example. The results are averaged over 10 seeds.  $\pm$  captures the standard deviation over seeds.

<b>RAMBO:</b>	1.49 $\pm$ 0.04	<b>COMBO:</b>	1.34 $\pm$ 0.14
---------------	-----------------	---------------	-----------------

354 Table 3 compares the performance of RAMBO and COMBO on the Single Transition Example.  
355 Further analysis in Appendix C.1 shows that for this problem, the pessimistic value function updates  
356 used by COMBO create local maxima in the  $Q$ -function throughout training. Policy optimisation can  
357 become stuck in these local maxima. On the other hand, the value function produced by RAMBO is  
358 initially optimistic, and pessimism is introduced into the value function *gradually* as the transition  
359 function is modified adversarially. As a result, the likelihood of the policy becoming stuck in poor  
360 local maxima is reduced. This may explain why RAMBO is able to achieve strong performance  
361 relative to existing algorithms in the MuJoCo domains. Gradually increasing the level of pessimism  
362 could be a useful modification for existing offline RL algorithms.

## 363 7 Conclusion and Future Directions

364 RAMBO is a promising new approach to offline RL which imposes conservatism by adversarially  
365 modifying the transition dynamics of a learnt model. Our approach is theoretically justified, and  
366 achieves state-of-the-art performance on standard benchmarks.

367 There are a number of possible extensions to RAMBO, some of which we have already discussed.  
368 In addition, we would like to apply RAMBO to image-space domains by using deep latent variable  
369 models to compress the state space [17, 50] and adversarially perturbing the transition dynamics in  
370 the latent space representation. Furthermore, we wish to improve the computational efficiency of  
371 RAMBO. RAMBO is more computationally demanding than many prior offline RL methods due to  
372 the need to simultaneously train the model as well as the policy. Potential directions to improve the  
373 efficiency include warmstarting the policy with behaviour cloning [24, 59] or reducing the frequency  
374 of model updates. Finally, another potential research direction is to apply the idea of adversarially  
375 trained dynamics models to improve robustness in the online setting.

376 **References**

- 377 [1] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective  
378 on offline reinforcement learning. In *International Conference on Machine Learning*, pages  
379 104–114. PMLR, 2020.
- 380 [2] Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline  
381 reinforcement learning with diversified q-ensemble. *Advances in neural information processing*  
382 *systems*, 34:7436–7447, 2021.
- 383 [3] Marcin Andrychowicz, Anton Raichuk, Piotr Stańczyk, Manu Orsini, Sertan Girgin, Raphaël  
384 Marinier, Leonard Hussenot, Matthieu Geist, Olivier Pietquin, Marcin Michalski, Sylvain Gelly,  
385 and Olivier Bachem. What matters for on-policy deep actor-critic methods? A large-scale study.  
386 In *International Conference on Learning Representations*, 2021.
- 387 [4] Arthur Argenson and Gabriel Dulac-Arnold. Model-based offline planning. In *International*  
388 *Conference on Learning Representations*, 2021.
- 389 [5] J Andrew Bagnell, Andrew Y Ng, and Jeff G Schneider. Solving uncertain Markov decision  
390 processes. Technical report, 2001.
- 391 [6] Philip J Ball, Cong Lu, Jack Parker-Holder, and Stephen Roberts. Augmented world models  
392 facilitate zero-shot dynamics generalization from a single offline environment. In *International*  
393 *Conference on Machine Learning*, pages 619–629. PMLR, 2021.
- 394 [7] David Brandfonbrener, Will Whitney, Rajesh Ranganath, and Joan Bruna. Offline RL without  
395 off-policy evaluation. *Advances in Neural Information Processing Systems*, 34:4933–4946,  
396 2021.
- 397 [8] Catherine Cang, Aravind Rajeswaran, Pieter Abbeel, and Michael Laskin. Behavioral priors  
398 and dynamics models: Improving performance and domain transfer in offline RL. In *Deep RL*  
399 *Workshop NeurIPS 2021*, 2021.
- 400 [9] Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor  
401 critic for offline reinforcement learning. *arXiv preprint arXiv:2202.02446*, 2022.
- 402 [10] Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement  
403 learning in a handful of trials using probabilistic dynamics models. *Advances in neural*  
404 *information processing systems*, 31, 2018.
- 405 [11] Ignasi Clavera, Violet Fu, and Pieter Abbeel. Model-augmented actor-critic: Backpropagating  
406 through paths. *International Conference on Learning Representations*, 2020.
- 407 [12] Sebastian Curi, Ilija Bogunovic, and Andreas Krause. Combining pessimism with optimism for  
408 robust and efficient model-based deep reinforcement learning. In *International Conference on*  
409 *Machine Learning*, pages 2254–2264. PMLR, 2021.
- 410 [13] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4RL: Datasets for  
411 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 412 [14] Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning.  
413 *Advances in Neural Information Processing Systems*, 34, 2021.
- 414 [15] Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning  
415 without exploration. In *International Conference on Machine Learning*, pages 2052–2062.  
416 PMLR, 2019.
- 417 [16] Jakob Gawlikowski, Cedrique Rovile Njieutcheu Tassi, Mohsin Ali, Jongseok Lee, Matthias  
418 Humt, Jianxiang Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, et al. A  
419 survey of uncertainty in deep neural networks. *arXiv preprint arXiv:2107.03342*, 2021.
- 420 [17] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. *Ad-*  
421 *vances in Neural Information Processing Systems*, 31, 2018.
- 422 [18] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-  
423 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International*  
424 *Conference on Machine Learning*, pages 1861–1870. PMLR, 2018.
- 425 [19] Toru Hishinuma and Kei Senda. Weighted model estimation for offline model-based reinforc-  
426 ment learning. *Advances in Neural Information Processing Systems*, 34, 2021.

- 427 [20] Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*,  
428 30(2):257–280, 2005.
- 429 [21] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model:  
430 Model-based policy optimization. *Advances in Neural Information Processing Systems*, 32,  
431 2019.
- 432 [22] Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza,  
433 Noah Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch deep reinforcement  
434 learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- 435 [23] Kirthevasan Kandasamy, Yoram Bachrach, Ryota Tomioka, Daniel Tarlow, and David Carter.  
436 Batch policy gradient methods for improving neural conversation models. *International Confer-*  
437 *ence on Learning Representations*, 2017.
- 438 [24] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. MOREL:  
439 Model-based offline reinforcement learning. *Advances in neural information processing systems*,  
440 33:21810–21823, 2020.
- 441 [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International*  
442 *Conference on Learning Representations*, 2015.
- 443 [26] Ilya Kostrikov, Rob Fergus, Jonathan Tompson, and Ofir Nachum. Offline reinforcement  
444 learning with Fisher divergence critic regularization. In *International Conference on Machine*  
445 *Learning*, pages 5774–5783. PMLR, 2021.
- 446 [27] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit  
447 Q-learning. In *International Conference on Learning Representations*, 2022.
- 448 [28] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-  
449 policy Q-learning via bootstrapping error reduction. *Advances in Neural Information Processing*  
450 *Systems*, 32, 2019.
- 451 [29] Aviral Kumar, Anikait Singh, Stephen Tian, Chelsea Finn, and Sergey Levine. A workflow for  
452 offline model-free robotic reinforcement learning. *Conference on Robot Learning*, 2021.
- 453 [30] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative Q-learning for  
454 offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–  
455 1191, 2020.
- 456 [31] Thanard Kurutach, Ignasi Clavera, Yan Duan, Aviv Tamar, and Pieter Abbeel. Model-ensemble  
457 trust-region policy optimization. In *International Conference on Learning Representations*,  
458 2018.
- 459 [32] Sascha Lange, Thomas Gabel, and Martin Riedmiller. Batch reinforcement learning. In  
460 *Reinforcement learning*, pages 45–73. Springer, 2012.
- 461 [33] Byung-Jun Lee, Jongmin Lee, and Kee-Eung Kim. Representation balancing offline model-  
462 based reinforcement learning. In *International Conference on Learning Representations*, 2020.
- 463 [34] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:  
464 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 465 [35] Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Off-policy policy gradient  
466 with state distribution correction. *International Conference on Machine Learning RL4RealLife*  
467 *Workshop*, 2019.
- 468 [36] Cong Lu, Philip Ball, Jack Parker-Holder, Michael Osborne, and Stephen J Roberts. Revisiting  
469 design choices in offline model based reinforcement learning. In *International Conference on*  
470 *Learning Representations*, 2022.
- 471 [37] Ajay Mandlekar, Fabio Ramos, Byron Boots, Silvio Savarese, Li Fei-Fei, Animesh Garg, and  
472 Dieter Fox. Iris: Implicit reinforcement without interaction at scale for learning control from  
473 offline robot manipulation data. In *2020 IEEE International Conference on Robotics and*  
474 *Automation (ICRA)*, pages 4414–4420. IEEE, 2020.
- 475 [38] Ajay Mandlekar, Yuke Zhu, Animesh Garg, Li Fei-Fei, and Silvio Savarese. Adversarially robust  
476 policy learning: Active construction of physically-plausible perturbations. In *2017 IEEE/RSJ*  
477 *International Conference on Intelligent Robots and Systems (IROS)*, pages 3932–3939. IEEE,  
478 2017.

- 479 [39] Tatsuya Matsushima, Hiroki Furuta, Yutaka Matsuo, Ofir Nachum, and Shixiang Gu.  
480 Deployment-efficient reinforcement learning via model-based offline optimization. In *In-*  
481 *ternational Conference on Learning Representations*, 2021.
- 482 [40] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G  
483 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.  
484 Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- 485 [41] Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Al-  
486 gaeDICE: Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019.
- 487 [42] Xinkun Nie, Emma Brunskill, and Stefan Wager. Learning when-to-treat policies. *Journal of*  
488 *the American Statistical Association*, 116(533):392–409, 2021.
- 489 [43] Arnab Nilim and Laurent El Ghaoui. Robust control of Markov decision processes with  
490 uncertain transition matrices. *Operations Research*, 53(5):780–798, 2005.
- 491 [44] OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew,  
492 Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, Jonas Schneider,  
493 Nikolas Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba,  
494 and Lei Zhang. Solving rubik’s cube with a robot hand. *arXiv preprint*, 2019.
- 495 [45] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua  
496 Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty?  
497 Evaluating predictive uncertainty under dataset shift. *Advances in Neural Information Process-*  
498 *ing Systems*, 32, 2019.
- 499 [46] Tom Le Paine, Cosmin Paduraru, Andrea Michi, Caglar Gulcehre, Konrad Zolna, Alexander  
500 Novikov, Ziyu Wang, and Nando de Freitas. Hyperparameter selection for offline reinforcement  
501 learning. *arXiv preprint arXiv:2007.09055*, 2020.
- 502 [47] Xinlei Pan, Daniel Seita, Yang Gao, and John Canny. Risk averse robust adversarial reinforce-  
503 ment learning. In *2019 International Conference on Robotics and Automation (ICRA)*, pages  
504 8522–8528. IEEE, 2019.
- 505 [48] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression:  
506 Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- 507 [49] Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial  
508 reinforcement learning. In *International Conference on Machine Learning*, pages 2817–2826.  
509 PMLR, 2017.
- 510 [50] Rafael Rafailov, Tianhe Yu, Aravind Rajeswaran, and Chelsea Finn. Offline reinforcement  
511 learning from images with latent space models. In *Learning for Dynamics and Control*, pages  
512 1154–1168. PMLR, 2021.
- 513 [51] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah  
514 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of*  
515 *Machine Learning Research*, 2021.
- 516 [52] Aravind Rajeswaran, Igor Mordatch, and Vikash Kumar. A game theoretic framework for  
517 model based reinforcement learning. In *International Conference on Machine Learning*, pages  
518 7953–7963. PMLR, 2020.
- 519 [53] Marc Rigter, Bruno Lacerda, and Nick Hawes. Minimax regret optimisation for robust planning  
520 in uncertain Markov decision processes. In *Proceedings of the AAAI Conference on Artificial*  
521 *Intelligence*, volume 35, pages 11930–11938. Association for the Advancement of Artificial  
522 Intelligence, 2021.
- 523 [54] Aurko Roy, Huan Xu, and Sebastian Pokutta. Reinforcement learning under model mismatch.  
524 *Advances in Neural Information Processing Systems*, 30, 2017.
- 525 [55] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-  
526 dimensional continuous control using generalized advantage estimation. *International Confer-*  
527 *ence on Learning Representations*, 2016.
- 528 [56] Susan M Shortreed, Eric Laber, Daniel J Lizotte, T Scott Stroup, Joelle Pineau, and Susan A  
529 Murphy. Informing sequential clinical decision-making through reinforcement learning: an  
530 empirical study. *Machine learning*, 84(1):109–136, 2011.

- 531 [57] Noah Siegel, Jost Tobias Springenberg, Felix Berkenkamp, Abbas Abdolmaleki, Michael  
532 Neunert, Thomas Lampe, Roland Hafner, Nicolas Heess, and Martin Riedmiller. Keep doing  
533 what worked: Behavior modelling priors for offline reinforcement learning. In *International*  
534 *Conference on Learning Representations*, 2020.
- 535 [58] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur  
536 Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of  
537 Go without human knowledge. *Nature*, 550(7676):354–359, 2017.
- 538 [59] Avi Singh, Albert Yu, Jonathan Yang, Jesse Zhang, Aviral Kumar, and Sergey Levine. Cog:  
539 Connecting new skills to past experience with offline reinforcement learning. *Conference on*  
540 *Robot Learning*, 2020.
- 541 [60] Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM*  
542 *Sigart Bulletin*, 2(4):160–163, 1991.
- 543 [61] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press,  
544 2018.
- 545 [62] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient meth-  
546 ods for reinforcement learning with function approximation. *Advances in Neural Information*  
547 *Processing Systems*, 12, 1999.
- 548 [63] Phillip Swazinna, Steffen Udluft, and Thomas Runkler. Overcoming model bias for robust offline  
549 deep reinforcement learning. *Engineering Applications of Artificial Intelligence*, 104:104366,  
550 2021.
- 551 [64] Aviv Tamar, Shie Mannor, and Huan Xu. Scaling up robust MDPs using function approximation.  
552 In *International Conference on Machine Learning*, pages 181–189. PMLR, 2014.
- 553 [65] Chen Tessler, Yonathan Efroni, and Shie Mannor. Action robust reinforcement learning and  
554 applications in continuous control. In *International Conference on Machine Learning*, pages  
555 6215–6224. PMLR, 2019.
- 556 [66] Masatoshi Uehara and Wen Sun. Pessimistic model-based offline reinforcement learning under  
557 partial coverage. *International Conference on Learning Representations*, 2022.
- 558 [67] Jianhao Wang, Wenzhe Li, Haozhe Jiang, Guangxiang Zhu, Siyuan Li, and Chongjie Zhang.  
559 Offline reinforcement learning with reverse model-based imagination. *Advances in Neural*  
560 *Information Processing Systems*, 34, 2021.
- 561 [68] Yue Wang and Shaofeng Zou. Online robust reinforcement learning with model uncertainty.  
562 *Advances in Neural Information Processing Systems*, 34, 2021.
- 563 [69] Paul J Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of*  
564 *the IEEE*, 78(10):1550–1560, 1990.
- 565 [70] Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust Markov decision processes.  
566 *Mathematics of Operations Research*, 38(1):153–183, 2013.
- 567 [71] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforce-  
568 ment learning. *Machine learning*, 8(3):229–256, 1992.
- 569 [72] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement  
570 learning. *arXiv preprint arXiv:1911.11361*, 2019.
- 571 [73] Tengyang Xie, Ching-An Cheng, Nan Jiang, Paul Mineiro, and Alekh Agarwal. Bellman-  
572 consistent pessimism for offline reinforcement learning. *Advances in Neural Information*  
573 *Processing Systems*, 34, 2021.
- 574 [74] Yijun Yang, Jing Jiang, Tianyi Zhou, Jie Ma, and Yuhui Shi. Pareto policy pool for model-based  
575 offline reinforcement learning. In *International Conference on Learning Representations*, 2021.
- 576 [75] Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea  
577 Finn. COMBO: Conservative offline model-based policy optimization. *Advances in Neural*  
578 *Information Processing Systems*, 34, 2021.
- 579 [76] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea  
580 Finn, and Tengyu Ma. MOPO: Model-based offline policy optimization. *Advances in Neural*  
581 *Information Processing Systems*, 33:14129–14142, 2020.
- 582 [77] Siyuan Zhang and Nan Jiang. Towards hyperparameter-free policy selection for offline rein-  
583 forcement learning. *Advances in Neural Information Processing Systems*, 34, 2021.

584 **Checklist**

585 The checklist follows the references. Please read the checklist guidelines carefully for information on  
586 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or  
587 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing  
588 the appropriate section of your paper or providing a brief inline description. For example:

- 589 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 590 • Did you include the license to the code and datasets? **[No]** The code and the data are  
591 proprietary.
- 592 • Did you include the license to the code and datasets? **[N/A]**

593 Please do not modify the questions and only use the provided macros for your answers. Note that the  
594 Checklist section does not count towards the page limit. In your paper, please delete this instructions  
595 block and only keep the Checklist section heading above along with the questions/answers below.

596 1. For all authors...

- 597 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
598 contributions and scope? **[Yes]**
- 599 (b) Did you describe the limitations of your work? **[Yes]** We described two main limi-  
600 tations: a) our approach performs less well on the AntMaze datasets (discussed in  
601 experiments), and b) requires more computation than many other approaches (discussed  
602 in conclusion).
- 603 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 604 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
605 them? **[Yes]**

606 2. If you are including theoretical results...

- 607 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
- 608 (b) Did you include complete proofs of all theoretical results? **[Yes]** We include a complete  
609 proof of Proposition 2. For the theoretical results from [66], we provide references to  
610 the appropriate parts of that paper.

611 3. If you ran experiments...

- 612 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
613 perimental results (either in the supplemental material or as a URL)? **[Yes]** See the  
614 supplementary material.
- 615 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were  
616 chosen)? **[Yes]** Detailed information about implementation details and hyperparameters  
617 is in the appendix.
- 618 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
619 ments multiple times)? **[Yes]** Standard deviation with respect to seeds is reported for  
620 all results.
- 621 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
622 of GPUs, internal cluster, or cloud provider)? **[Yes]** Details of compute used per run is  
623 in the experiments section. Details of total compute used for full evaluation is in the  
624 appendix.

625 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 626 (a) If your work uses existing assets, did you cite the creators? **[Yes]** We cite the D4RL  
627 benchmark datasets [13].
- 628 (b) Did you mention the license of the assets? **[Yes]** In Appendix B.6.
- 629 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**  
630 Code for our work is in the supplementary material.
- 631 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
632 using/curating? **[N/A]**
- 633 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
634 information or offensive content? **[N/A]**

635

5. If you used crowdsourcing or conducted research with human subjects...

636

(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]

637

638

(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]

639

640

(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

641