

---

# Less Tuning, Better Planning: Simplifying Offline Model-Based Planning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Offline reinforcement learning enables policies to be learned from previously  
2 collected experiences without requiring online interaction. However, these policies  
3 are typically deployed as fixed, zero-shot agents and lack the ability to adapt  
4 their behavior at test time. Offline model-based planning offers a promising way  
5 to enable flexible test-time adaptation, but its performance is highly sensitive to  
6 critical design choices, particularly the planning horizon and the action proposer.  
7 In practice, these choices are often tuned through online evaluation, contradicting  
8 the premise of offline RL. In this work, we introduce Soft Horizon Aggregation for  
9 Planning (SHARP), an offline plug-and-play planning method that eliminates the  
10 need for an online-tuned planning horizon. Instead of using a fixed horizon across  
11 all states, SHARP performs soft horizon aggregation by dynamically weighting  
12 returns according to model uncertainty estimated from an ensemble of dynamics  
13 models. We further investigate the role of the action proposer and find that stronger  
14 offline policies do not necessarily lead to better planning performance. Instead,  
15 a simple behavior cloning (BC) policy is often sufficient as an action proposer  
16 while avoiding the effort required for extensive policy extraction. Combining  
17 these insights, we propose SHARP-BC, which consistently outperforms existing  
18 baselines while reducing reliance on extensive online hyperparameter tuning.

## 19 1 Introduction

20 Offline reinforcement learning (RL) [1] enables the conversion of previously collected experiences  
21 into policies for control problems without requiring online interaction during training. It is particularly  
22 appealing in safety-critical or resource-constrained domains where online exploration is expensive,  
23 risky, or impractical, such as robotic manipulation [2, 3], autonomous driving [4, 5], and healthcare  
24 decision-making [6]. A central challenge in offline RL is extrapolation error: at test time, the learned  
25 policy may encounter state-action pairs that are poorly covered by the offline dataset, leading to  
26 unreliable value estimates. Prior work typically addresses this challenge by learning policies that  
27 favor safe behaviors observed in the dataset [7–10], or by training adaptive policies that condition on  
28 histories, or uncertainty estimates to infer environment variation [11, 12]. However, these policies  
29 are usually optimized only during training and then deployed as fixed policies in a zero-shot manner,  
30 without the ability to adapt their behavior at test time.

31 To enable adaptive decision-making at deployment, offline model-based planning methods [13, 14]  
32 learn dynamic models from offline data and perform model-predictive control [15] during evaluation.  
33 A representative line of work follows a plug-and-play paradigm [16, 17], where components learned  
34 by offline RL, such as the policy and value function, are used to guide planning at test time. The  
35 procedure typically involves (1) generating candidate actions using an action proposer (usually the  
36 policy), (2) evaluating their expected returns with the dynamic models and value function, and (3)  
37 selecting the best action to be executed in the environment through an optimization procedure. By

38 leveraging knowledge from dynamic models, current plug-and-play methods enable flexible test-time  
39 adaptation and effectively improve policy performance [16, 17].

40 In plug-and-play offline planning, two design choices are especially critical: (1) **planning horizon**:  
41 how far the planner should roll out the learned dynamics, and (2) **action proposer**: which action  
42 distribution to use to generate candidate trajectories. Existing methods typically expose these  
43 components as hyperparameters or design choices, whose best settings can only be identified through  
44 online evaluation, contradicting the premise of “offline” RL. In this work, we aim to eliminate such  
45 online-tuned choices by designing a planning algorithm that avoids explicit horizon selection and  
46 does not require separately tuning the action proposer.

47 First, we examine the most influential hyperparameters in offline model-based planning: the planning  
48 horizon. It controls how far into the future we should roll out, aiming to balance the information  
49 gained and errors produced by the dynamic model. Existing methods typically employ a fixed  
50 planning horizon across all states and tune it via online evaluation. However, this practice can be  
51 highly sensitive and is often impractical. We empirically find that the optimal planning horizon can  
52 vary substantially across tasks and even across states, as illustrated in Figure 1. Therefore, using a  
53 single fixed horizon may be suboptimal, motivating the need for adaptive horizon selection.

54 To this end, we propose **Soft Horizon Aggregation for Planning (SHARP)**. Instead of using a fixed  
55 horizon across all states, we use soft horizon aggregation to achieve stable performance across various  
56 planning horizons. Specifically, we first generate multiple rollout branches from an ensemble of  
57 dynamics models, then estimate the mean and variance of the expected returns at every horizon across  
58 all branches. The estimated variance is then used to soft-aggregate the return across all horizons:  
59 horizons with larger variance receive smaller weights, reflecting greater model uncertainty at those  
60 horizons. On D4RL locomotion datasets [18], we show that our method achieves substantially greater  
61 stability across different planning horizons. Moreover, it delivers superior performance compared to  
62 existing baselines without requiring online hyperparameter tuning.

63 Second, we investigate the role of the action proposer, an underexplored component in planning.  
64 Intuitively, using a higher-performance action proposer should yield better planning performance.  
65 Therefore, existing works tune and leverage policies with higher zero-shot performance as their action  
66 proposer [16, 17]. However, policy extraction introduces an additional design space that requires  
67 online hyperparameter tuning. By conducting a systematic study of different policy extraction  
68 methods, including advantage-weighted regression (AWR) [19] and behavior-constrained policy  
69 gradients (DDPG+BC) [9], we find that stronger offline policies do not necessarily contribute to  
70 better planning performance. Surprisingly, simple **behavior cloning (BC)**, despite its low zero-  
71 shot performance, can still serve as a sufficient action proposer to achieve comparable planning  
72 performance. Therefore, we suggest using BC as the default action proposer in offline model-based  
73 planning, as it reduces the effort required for extensive policy extraction and hyperparameter tuning  
74 while maintaining strong planning performance.

75 Together, we propose **SHARP-BC**, a simple yet effective offline model-based planning framework  
76 that combines soft horizon aggregation with a simple BC policy as an action proposer. By eliminating  
77 the need for sensitive horizon tuning and reducing reliance on complex policy extraction, our  
78 approach simplifies the design of plug-and-play planning methods. Empirically, SHARP-BC achieves  
79 competitive or superior performance compared to other existing planning methods, without requiring  
80 any extensive online hyperparameter tuning.

## 81 2 Preliminaries

82 This work studies under the framework of Markov decision processes (MDPs), characterized by a tuple  
83  $M = (\mathcal{S}, \mathcal{A}, T, r, d_0, \gamma)$ , where  $\mathcal{S}$  and  $\mathcal{A}$  denote the continuous state and action spaces,  $T(s' | s, a)$   
84 is the transition probability,  $r(s, a)$  is the reward function,  $d_0$  is the initial state distribution, and  
85  $\gamma \in [0, 1]$  is the discount factor.

### 86 2.1 Offline reinforcement learning

87 In offline reinforcement learning (offline RL), an agent is given access to a fixed dataset  $\mathcal{D} =$   
88  $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$  collected by some unknown behavior policy, and is not permitted to interact  
89 with the environment during training. The objective is to learn a policy  $\pi(a | s)$  that maximizes the  
90 expected discounted return:  $J(\pi) = \mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$ .

91 A common approach in offline RL is to learn a state-action value function (Q-function):  $Q^\pi(s, a) =$   
 92  $\mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a]$ , which estimates the expected return obtained by taking  
 93 action  $a$  in state  $s$  and subsequently following policy  $\pi$ . The learned Q-function can then be used for  
 94 either policy improvement or action evaluation.

## 95 2.2 Offline model-based planning

96 Offline model-based methods decouple policy learning from decision-making by first learning a  
 97 dynamic model  $\hat{T}(s' \mid s, a)$ , a reward model  $\hat{r}(s, a)$ , and often a value function  $\hat{Q}(s, a)$  from the  
 98 offline dataset. Then, these methods leverage these components to perform decision-making at test  
 99 time, rather than directly executing the policy, by simulating candidate trajectories, evaluating them  
 100 using the learned models, and selecting the most promising actions according to a planning objective.

101 Offline model-based planning can be formulated as a receding-horizon optimization problem. Given  
 102 a state  $s_0$ , planning horizon  $H$ , and discount factor  $\gamma$ , the planner first generates a set of  $K$  candidate  
 103 actions  $\{\mathbf{a}_0^k\}_{k=1}^K$ , where each action is sampled from a stochastic single-step action proposer  $u(a \mid s)$ .

104 The expected return of each candidate action is computed as  $G^k = \mathcal{G}(s_0, \mathbf{a}_0^k; \hat{T}, \hat{r}, \hat{Q}, H, \gamma)$ , where  $\mathcal{G}$   
 105 defines the rollout strategy and return computation under the learned dynamics model  $\hat{T}$ , combining  
 106 predicted rewards with terminal value estimates along the simulated trajectory.

107 Finally, instead of strictly restricting the choice to the single highest-scoring candidate, the executed  
 108 action  $a^*$  is determined by an action selection rule  $\mathcal{A}$ :

$$a^* = \mathcal{A}(\{\mathbf{a}_0^k, G^k\}_{k=1}^K). \quad (1)$$

109 We summarize the overall planning procedure in Algorithm 1, where the specific instantiations of the  
 110 proposer  $u$ , return estimator  $\mathcal{G}$ , and selection rule  $\mathcal{A}$  determine the resulting planning algorithm.

---

### Algorithm 1 General offline model-based planning

---

- 1: **Input:** Action proposer  $u(a \mid s)$ , dynamic model  $\hat{T}$ , value function  $\hat{Q}$ , reward model  $\hat{r}$ , return  
function  $\mathcal{G}$ , action selection function  $\mathcal{A}$ , candidate number  $K$ , horizon  $H$
  - 2: Observe initial state  $s_0$
  - 3: **while** episode not terminated **do**
  - 4:   **for**  $k = 1$  **to**  $K$  **do**
  - 5:      $a_0^k \sim u(\cdot \mid \hat{s}_0^k)$  ▷ Action candidate generation
  - 6:      $G^k = \mathcal{G}(s_0, a_0^k; \hat{T}, \hat{r}, \hat{Q}, H, \gamma)$ , ▷ Rollout the trajectory and compute the return:
  - 7:   **end for**
  - 8:   Select action  $a^* = \mathcal{A}(\{\mathbf{a}_0^k, G^k\}_{k=1}^K)$  ▷ Select an action to be executed
  - 9:   Execute  $a^*$ , observe real next state  $s_0 \leftarrow s'$
  - 10: **end while**
- 

## 111 2.3 Plug-and-play planning methods

112 To leverage the knowledge encoded in offline RL policies  $\pi$  and value functions, plug-and-play  
 113 planning methods [13, 16, 17] incorporate action proposers  $u$  and value functions  $\hat{Q}$  learned via  
 114 offline RL into the planning process. Model-Based Offline Planning (MBOP) [13] is among the first  
 115 works to introduce model-based planning to offline scenarios and is among the simplest planning  
 116 methods. It rolls out  $K$  candidate actions using the dynamic models and evaluates their expected  
 117 return at horizon  $H$ . The best action is chosen in a softmax-like manner. We consider the plug-and-  
 118 play version in this work, where the action proposer and value function are from offline RL algorithms.  
 119 Learning off-policy with online planning (LOOP) [16] constrains planned action sequences with  
 120 a KL-divergence regularizer to limit deviation from the action proposer. We consider its offline  
 121 variant, LOOP-offline, as a plug-and-play offline model-based planning method. Unlike MBOP,  
 122 LOOP performs ensemble rollouts over future states to capture model uncertainty via the standard  
 123 deviation of predicted returns. More details on MBOP and LOOP are provided in Appendix C.1.

## 124 3 Making planning horizon insensitive

125 Existing plug-and-play planning algorithms have been shown to improve upon the base policies [16,  
 126 17], while appealing, we find that their performance is highly sensitive to the hyperparameter:

127 planning horizon  $H$ . We hypothesize that this sensitivity arises because prior methods employ a  
 128 **fixed** planning horizon across all states during online evaluation. To investigate this issue, we first  
 129 present an empirical study in Section 3.1 that illustrates why a fixed horizon can be problematic. We  
 130 then introduce **Soft Horizon AggRegation for Planning** in Section 3.2, a novel planning algorithm  
 131 designed to mitigate horizon sensitivity. Finally, in Section 3.3, we demonstrate that our method is  
 132 substantially more robust to the choice of planning horizon  $H$ , leading to improved performance in  
 133 settings where hyperparameter tuning is impractical.

### 134 3.1 Why fixed horizons fail?

135 The planning horizon determines how far into the future the planner should roll out trajectories,  
 136 balancing long-term information gain with the accumulation of model errors. However, this  
 137 trade-off may be state-dependent due to uneven state coverage in the offline dataset.  
 140

141 To investigate this phenomenon, we conduct an experiment comparing the learned dynamics  
 142 model with the ground-truth dynamics of the environment. On the D4RL [18] *HalfCheetah*-  
 143 *v2* environment, we train the dynamics model and the IQL policy using the full-replay dataset.  
 144 Training details are provided in Appendix D and E. Starting from the same states and using  
 145 identical action sequences sampled from the policy, we roll out both the learned dynamics  
 146 model and the real environment, and measure the resulting L2 discrepancies between observations  
 147 over horizons. As shown in Figure 1, the discrepancy grows as the planning horizon in-  
 148 creases, while also varying substantially across states. These results suggest that the optimal planning  
 149 horizon is state-dependent: some states benefit from longer rollouts that capture more long-term  
 150 information, whereas others require shorter horizons to mitigate the accumulation of model error.  
 151 Therefore, previous planning methods that use a fixed horizon across all states require careful online  
 152 hyperparameter tuning, which limits their practical use.

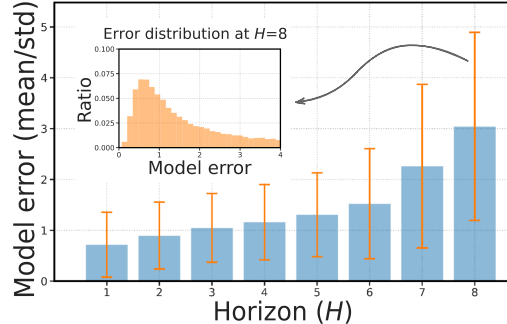


Figure 1: **Optimal planning horizon is state-dependent.** To study why a single fixed horizon can be suboptimal, we analyze model error across different planning horizons on *HalfCheetah-full-replay*. As the horizon increases, the mean error grows, but more importantly, the error variance also increases substantially. This indicates that long-horizon planning is not uniformly harmful across all states.

### 160 3.2 SHARP: Soft Horizon AggRegation for Planning

161 To address the state-dependent model error issue, we propose **Soft Horizon AggRegation for Planning** (SHARP), a plug-and-play, offline model-based planning method that is robust to planning  
 162 horizon  $H$ . Instead of relying on a single, rigid rollout length, the key idea is to dynamically adjust  
 163 the effective planning horizon for each state by normalizing predicted returns based on ensemble  
 164 variance. SHARP achieves this through the following steps: (1) simulating diverse future trajectories  
 165 via ensemble rollouts, (2) estimating horizon-specific returns and precision, (3) soft aggregating these  
 166 horizon-wise returns using precision as weights, and (4) selecting the final action for execution.

168 **Simulating diverse future trajectories.** Similar to LOOP [16], we employ an ensemble of dynamics  
 169 models. For each model, we maintain  $P$  particles, each corresponding to an independent rollout  
 170 branch. For a candidate action sequence  $k$ , future states and actions are sampled according to

$$\hat{s}_{h+1,e,p}^k \sim \hat{T}_e(\cdot | \hat{s}_{h,e,p}^k, a_{h,e,p}^k), \quad a_{h+1,e,p}^k \sim u(\cdot | \hat{s}_{h+1,e,p}^k), \quad (2)$$

171 where the rollout is initialized with  $\hat{s}_{0,e,p}^k \leftarrow s_0$  and  $a_{0,e,p}^k \leftarrow a_0^k$  for all ensemble members  $e$  and  
 172 particles  $p$ .

173 **Horizon-specific return and confidence estimation.** To safely estimate performance, we calculated  
 174 the expected return  $G_{e,p,h}^k$  at every horizon  $h$  for each candidate  $k$ , across every branch (ensemble  $e$   
 175 and particle  $p$ ):

$$G_{e,p,h}^k = \sum_{t=0}^h \gamma^t \hat{r}(\hat{s}_{t,e,p}^k, a_{t,e,p}^k) + \gamma^{h+1} \hat{Q}(\hat{s}_{h+1,e,p}^k, a_{h+1,e,p}^k), \quad (3)$$

176 Then we calculate the mean  $\mu_h^k$  and precision  $\rho_h^k$  across all ensemble members and particles:

$$\mu_h^k = \text{mean}_{e,p}(G_{e,p,h}^k), \quad \rho_h^k = 1/\text{var}_{e,p}(G_{e,p,h}^k) \quad (4)$$

177 **Soft horizon aggregation.** We compute a final, normalized return  $G^k$  for each candidate action by  
 178 taking a precision-weighted sum of the returns across all horizons:

$$G^k = \sum_{h=0}^{H-1} \left( \frac{\rho_h^k}{\sum_{h'=0}^{H-1} \rho_{h'}^k} \right) \mu_h^k \quad (5)$$

179 Intuitively, a large variance indicates high uncertainty about a specific rollout depth, so that the  
 180 horizon is downweighted to prevent error accumulation. Conversely, a smaller variance indicates  
 181 greater certainty, leading to a larger weight being assigned to that horizon’s mean return. We call this  
 182 **soft horizon aggregation** because it “softly” combines return estimates across multiple horizons.

183 **Selecting the final action.** Finally, the action is then selected as:

$$a^* = \frac{\sum_{k=1}^K \exp(\kappa G^k) a_0^k}{\sum_{k=1}^K \exp(\kappa G^k)} \quad (6)$$

184 where  $\kappa$  is a hyperparameter controlling how strongly to emphasize actions with larger estimated  
 185 return  $G^k$ . The full algorithm is presented in Algorithm 2.

---

**Algorithm 2** Soft Horizon Aggregation for Planning

---

```

1: Input: Action proposer  $u(a | s)$ , ensemble of dynamic models  $\{\hat{T}_e\}_{e=1}^E$ , value function  $\hat{Q}$ ,
   reward model  $\hat{r}$ , candidate number  $K$ , horizon  $H$ , Particle numer  $P$ , discount factor  $\gamma$ 
2: Observe initial state  $s_0$ 
3: while episode not terminated do
4:   for  $k = 1$  to  $K$  do
5:      $a_0^k \leftarrow u(\cdot | s_0)$ 
6:     Set  $\hat{s}_{0,e,p}^k \leftarrow s_0, a_{0,e,p}^k \leftarrow a_0^k \quad \forall e, p$ 
7:     for  $e = 1$  to  $E$  do ▷ Simulate diverse future trajectories
8:       for  $p = 1$  to  $P$  do
9:         for  $h = 0$  to  $H - 1$  do
10:          Sample next state  $\hat{s}_{h+1,e,p}^k \sim \hat{T}_e(\cdot | \hat{s}_{h,e,p}^k, a_{h,e,p}^k)$ 
11:          Sample next action  $a_{h+1,e,p}^k \sim u(\cdot | \hat{s}_{h+1,e,p}^k)$ 
12:           $G_{e,p,h}^k = \sum_{t=0}^h \gamma^t \hat{r}(\hat{s}_{t,e,p}^k, a_{t,e,p}^k) + \gamma^{h+1} \hat{Q}(\hat{s}_{h+1,e,p}^k, a_{h+1,e,p}^k)$ 
13:        end for
14:      end for
15:       $\mu_{e,p}^k = \text{mean}_{e,p}(G_{e,p,h}^k)$  ▷ Horizon-Specific Return
16:       $\rho_{e,p}^k = 1/\text{var}_{e,p}(G_{e,p,h}^k)$  ▷ Horizon-Specific Precision
17:    end for
18:     $G^k = \sum_{h=0}^{H-1} \left( \frac{\rho_{e,p,h}^k}{\sum_{h'=0}^{H-1} \rho_{e,p,h'}^k} \right) \mu_{e,p}^k$  ▷ Soft Horizon Aggregation
19:  end for
20:  Select  $a^* = \frac{\sum_{k=1}^K \exp(\kappa G^k) a_0^k}{\sum_{k=1}^K \exp(\kappa G^k)}$  ▷ Weighted action selection
21:  Execute  $a^*$ , observe real next state  $s_0 \leftarrow s'$ 
22: end while

```

---

### 186 3.3 Experiment

187 We compare our proposed SHARP with prior plug-and-play planning methods, MBOP [13] and  
 188 LOOP [16]. We train Q-functions using two classic offline RL algorithms, IQL and CQL, and use  
 189 their trained policies as the action proposer  $q$ . We evaluate on the D4RL locomotion benchmarks [18]  
 190 across three environments (*HalfCheetah*, *Hopper*, and *Walker2d*), each with five dataset configurations  
 191 (D): random (R), medium (M), medium-replay (MR), medium-expert (ME), and full-replay (FR).

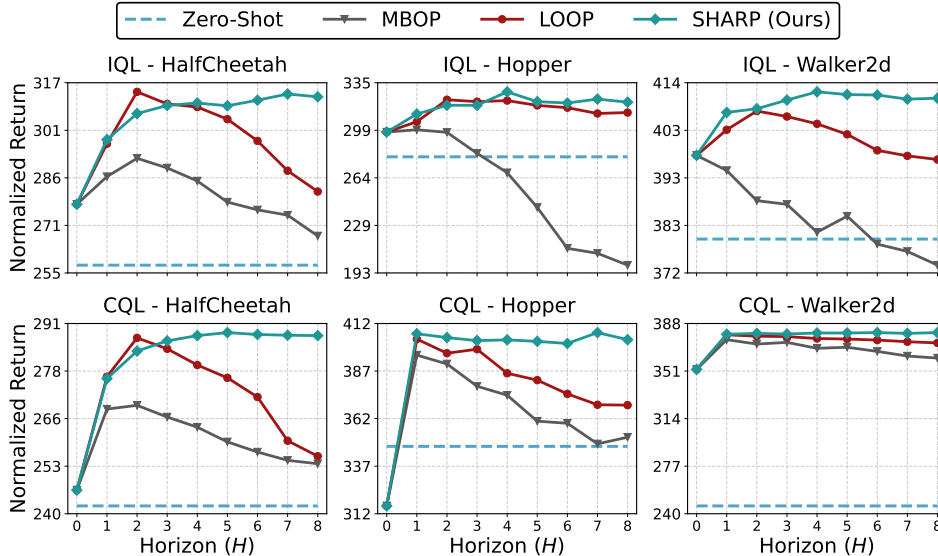


Figure 2: **SHARP is robust to different horizons.** We vary the test-time planning horizon  $H$  for IQL- and CQL-based planners on D4RL locomotion tasks. MBOP and LOOP are sensitive to horizon choice, often improving at short horizons but degrading as  $H$  increases due to accumulated model errors. In contrast, SHARP maintains stable performance across horizons.

192 Details of the environments, datasets, and offline RL training setup, including the IQL and CQL  
 193 training details, are provided in Appendix C.3 and Appendix D.

194 To examine the stability of all methods across different planning horizons, we vary the planning  
 195 horizon  $H$  from 1 to 8 and evaluate their effects on planning performance. For each horizon, we  
 196 report the best performance obtained by searching over  $\kappa \in \{0.1, 0.5, 1.0, 5.0, 10.0\}$  averaged over 3  
 197 seeds and 20 evaluation episodes.

198 The results in Figure 2 show that the planning horizon  $H$  plays a critical yet double-edged role in the  
 199 performance of existing plug-and-play methods such as MBOP and LOOP. As the planning horizon  
 200 increases, performance initially improves, highlighting the benefits of longer-horizon planning.  
 201 However, beyond a certain point, the performance of MBOP and LOOP declines significantly. This  
 202 degradation arises because the compounding of model errors eventually outweighs the advantages of  
 203 deeper planning. As a result, these methods are highly sensitive to the choice of  $H$ , requiring careful  
 204 environment-specific hyperparameter tuning.

205 In contrast, our proposed SHARP addresses this limitation through soft-horizon aggregation. This  
 206 mechanism enables the use of a universally large maximum horizon  $H$  by automatically adapting the  
 207 contribution of deeper planning steps. Specifically, rollouts are weighted according to uncertainty  
 208 estimates derived from the ensemble branches, allowing the method to progressively down-weight  
 209 distant and error-prone predictions. As a result, SHARP maintains stable and strong performance  
 210 even as the planning horizon increases, effectively eliminating the need for manual horizon tuning.

211 To rigorously evaluate planning performance without online hyperparameter tuning, we conduct  
 212 a grid search over the planning hyperparameters, varying the horizon  $H$  from 1 to 8 and  $\kappa \in$   
 213  $\{0.1, 0.5, 1.0, 5.0, 10.0\}$  and report the average performance across all configurations in Table 1.  
 214 Under this challenging evaluation protocol, our method consistently outperforms both the original  
 215 offline RL policies (Original) and the MBOP and LOOP baselines. The consistent gains across  
 216 both IQL- and CQL-based policies demonstrate that our approach not only improves performance  
 217 substantially but also significantly reduces sensitivity to online hyperparameter tuning. We report the  
 218 best performance after tuning for all planning methods in Appendix B.1.

#### 219 4 Does a stronger policy make a better action proposer?

220 The action proposer is a fundamental yet underexplored component in offline model-based planning.  
 221 Existing works [16, 17] typically adopt policies trained via offline RL as action proposers, motivated

Table 1: **SHARP achieves the best test-time planning performance without tuning  $H$  or  $\kappa$ .** We report the normalized return averaged over all tested planning horizons  $H$  and action-selection temperatures  $\kappa$ , simulating a setting where online selection of the best planning configuration is unavailable. SHARP achieves stronger average performance than both the zero-shot offline RL policies and plug-and-play planning baselines across IQL and CQL.

Environment	D	IQL				CQL			
		Zero-shot	MBOP	LOOP	SHARP (Ours)	Zero-shot	MBOP	LOOP	SHARP (Ours)
HalfCheetah	R	6.8	17.7	<u>19.0</u>	<b>19.3</b>	4.4	4.2	<u>5.2</u>	<b>5.4</b>
	M	46.9	51.6	<u>53.6</u>	<b>54.5</b>	44.9	50.3	<u>52.4</u>	<b>53.4</b>
	MR	41.7	43.1	<u>46.4</u>	<b>47.9</b>	29.3	31.0	<u>32.3</u>	<b>33.0</b>
	ME	<b>90.7</b>	72.5	<u>92.7</u>	77.4	91.5	92.3	<u>99.3</u>	<b>102.3</b>
	FR	71.7	71.8	<u>74.8</u>	<b>80.6</b>	72.4	76.5	<u>78.2</u>	<b>81.0</b>
Hopper	R	7.6	7.8	<u>7.9</u>	7.8	4.3	<b>5.1</b>	<u>5.0</u>	4.5
	M	51.8	41.2	<u>56.4</u>	<b>58.9</b>	58.5	63.1	<u>69.1</u>	<b>71.9</b>
	MR	55.9	70.7	<u>87.9</u>	<b>88.2</b>	20.0	79.2	<u>84.2</u>	<b>93.3</b>
	ME	<b>70.9</b>	30.9	34.7	31.5	62.6	105.9	<b>107.9</b>	<u>106.0</u>
	FR	<b>93.4</b>	69.1	<u>91.6</u>	<u>91.4</u>	101.0	103.5	<u>104.3</u>	<b>104.6</b>
Walker2d	R	4.8	4.8	<u>5.1</u>	<b>5.5</b>	<b>0.1</b>	<b>0.1</b>	<b>0.1</b>	<b>0.1</b>
	M	83.5	80.7	<u>87.6</u>	<b>89.3</b>	73.5	81.7	<u>83.5</u>	<b>84.0</b>
	MR	83.9	82.5	<u>92.8</u>	<b>95.1</b>	74.2	79.7	<u>84.4</u>	<b>86.2</b>
	ME	<b>109.5</b>	108.2	<u>109.1</u>	108.6	108.6	<u>109.2</u>	<b>110.0</b>	<b>110.0</b>
	FR	97.9	96.5	<u>99.9</u>	<b>100.7</b>	91.1	93.5	<u>94.2</u>	<b>95.2</b>

222 by the intuition that higher-quality policies should naturally yield better planning performance. Under  
 223 this view, one would expect that improving the policy, through better objectives or hyperparameter  
 224 tuning, should directly translate into improved planning outcomes.

225 In this section, we re-examine this assumption by systematically varying both the policy extraction  
 226 method and the degree of optimization, and evaluating the resulting policies as action proposers  
 227 during planning. Specifically, we consider two policy-extraction objectives (described in Section 4.1)  
 228 and analyze how planning performance varies with optimization strength (Section 4.2).

#### 229 4.1 Policy extraction in Offline-RL

230 Policy extraction is a crucial design choice in Offline-RL, aiming to maximize policy value while  
 231 staying close to the behavior policy represented in the dataset. We consider the following two widely  
 232 used objectives for offline policy extraction:

233 **(1) Weighted Behavioral Cloning (AWR) [19].** We adopt advantage-weighted regression (AWR),  
 234 a simple and widely used offline policy extraction method:

$$\max_{\pi} \mathcal{J}_{\text{AWR}}(\pi) = \mathbb{E}_{s, a \sim \mathcal{D}} \left[ e^{\alpha(Q(s, a) - V(s))} \log \pi(a | s) \right], \quad (7)$$

235 where  $\alpha$  is a temperature parameter. As  $\alpha \rightarrow 0$ , the objective reduces to standard behavioral cloning,  
 236 while larger  $\alpha$  emphasizes high-advantage actions.

237 **(2) Behavior-Constrained Policy Gradient (DDPG+BC) [9].** We also consider a behavior-  
 238 constrained objective that combines policy gradient with behavior cloning:

$$\max_{\pi} \mathcal{J}_{\text{DDPG+BC}}(\pi) = \mathbb{E}_{s, a \sim \mathcal{D}} [Q(s, \mu^{\pi}(s)) + \alpha \log \pi(a | s)], \quad (8)$$

239 where  $\mu^{\pi}(s) = \mathbb{E}_{a \sim \pi(\cdot | s)}[a]$ . When  $\alpha \rightarrow \infty$ , the objective reduces to behavioral cloning, where a  
 240 smaller  $\alpha$  places more emphasis on maximizing  $Q$ .

#### 241 4.2 Proposal policy experiment

242 To investigate whether the quality of the action proposer truly influences planning performance, we  
 243 apply both AWR and DDPG+BC with varying values of  $\alpha$ , using value functions trained by IQL  
 244 to extract the policy  $\pi$ . We then use the resulting policy as the action proposer for all planning  
 245 methods, MBOP, LOOP, SHARP (Ours). For each planning method, following prior works [16, 17],  
 246 we tune the planning hyperparameter over  $H \in \{2, 4\}$  and  $\kappa \in \{0.1, 0.5, 1.0, 5.0, 10.0\}$ , and report  
 247 the best-performing setting averaged over 3 seeds and 20 evaluation episodes.

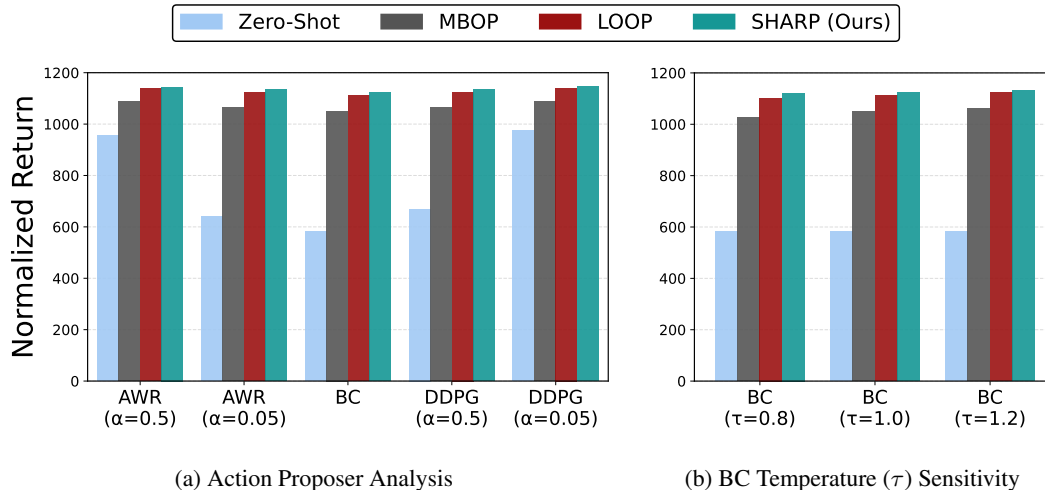


Figure 3: **Planning is robust to the choice of action proposer.** We compare action proposers obtained from AWR, DDPG+BC, and BC using IQL-trained value functions. Despite large differences in the zero-shot performance of the RL policies, planning performance remains similar across different action proposers and BC temperatures.

248 The results in Figure 3a show that different  $\alpha$  leads to substantial variation in zero-shot policy perfor-  
 249 mance, as expected. Surprisingly, however, the planning performance remains largely unchanged  
 250 across different action proposers and values of  $\alpha$ , despite the wide gap in their zero-shot performance.  
 251 Moreover, we find that even a simple behavior cloning (BC) policy achieves planning performance  
 252 comparable to that of more sophisticated policy extraction methods. These results suggest that,  
 253 although policy extraction methods require careful objective and hyperparameter tuning to achieve  
 254 strong zero-shot performance, such optimization plays a much less important role when the policy is  
 255 used as an action proposer for planning. We include the analyses across dataset types in Appendix B.3.

256 We further analyze the effect of varying the BC temperature on planning performance, where higher  
 257 temperatures correspond to larger action distributions. The results in Figure 3b show that planning  
 258 performance is likewise largely insensitive to the BC temperature. Taken together, these findings  
 259 reveal a surprising decoupling between policy quality and planning performance. In practice, the  
 260 role of the action proposer appears to be primarily to provide a sufficiently diverse set of candidate  
 261 actions, while the policy’s optimality becomes largely irrelevant once sufficient coverage is achieved.

262 Together, these findings substantially simplify action proposer design: simple BC is sufficient—  
 263 sophisticated policy extraction methods and additional hyperparameter tuning are not needed.

## 264 5 Comparison with end-to-end planning methods

265 We combine the findings from the previous sections to derive SHARP-BC, which combines soft  
 266 horizon aggregation with a simple BC policy as an action proposer. We compare SHARP-BC against  
 267 classical end-to-end planning baselines, *e.g.*, sequence-based methods (Diffuser [20], DT [21]) and  
 268 value-based approaches (MBOP [13], IQL-TD-MPC [22]). See Appendix C.2 for more details.

269 We report the performance of SHARP-BC with and without hyperparameter search (by averaging  
 270 over all  $H$  from 1 to 8 and  $\kappa \in \{0.1, 0.5, 1.0, 5.0, 10.0\}$ ). The results are presented in Table 2, where  
 271 \* denotes the default setting without hyperparameter tuning and unmarked SHARP-BC columns  
 272 denote tuned settings. Importantly, the baseline methods report results using their best-performing  
 273 hyperparameters obtained through online hyperparameter search. In contrast, even without any  
 274 tuning, our method consistently achieves competitive or superior performance relative to these fully  
 275 tuned baselines. Moreover, when a tuning budget is available, SHARP-BC can further leverage  
 276 hyperparameter search to attain even stronger performance.

277 Overall, these findings demonstrate that our proposed SHARP-BC requires less tuning and achieves  
 278 better planning, and can scale effectively with additional tuning resources.

Table 2: **Performance comparison with end-to-end planning baselines on D4RL locomotion tasks.** We compare SHARP-BC against sequence-based and value-based offline planning baselines. The baseline results are reported using their best-performing settings. For SHARP-BC, columns marked with \* denote the default setting without hyperparameter tuning, while unmarked columns use a full tuning budget. Results show that SHARP-BC remains competitive without tuning and further improves when tuning is available.

Environment	Dataset	Sequence-based		Value-based		SHARP-BC (Ours)			
		DT	Diffuser	IQL-TD-MPC	MBOP	IQL*	IQL	CQL*	CQL
HalfCheetah	M	42.6	42.8	<u>57.4</u>	44.6	53.9	<b>59.6</b>	52.9	57.2
	MR	36.6	37.7	<u>49.2</u>	42.3	48.5	<b>52.0</b>	31.7	32.8
	ME	86.8	88.9	44.8	<u>105.9</u>	95.8	103.4	102.8	<b>106.9</b>
Hopper	M	67.6	74.3	–	12.4	68.2	<u>83.1</u>	73.6	<b>88.4</b>
	MR	82.7	93.6	–	48.8	92.2	<u>99.2</u>	96.0	<b>100.1</b>
	ME	107.6	103.3	–	55.1	101.4	<u>109.7</u>	106.7	<b>110.7</b>
Walker2d	M	74.0	79.6	–	9.7	<u>85.4</u>	<b>91.0</b>	83.7	84.8
	MR	66.6	70.6	–	41.0	87.6	<b>95.2</b>	85.7	<u>89.1</u>
	ME	108.1	106.9	–	70.2	109.8	<b>110.4</b>	109.9	<u>110.2</u>
<b>Average</b>		74.7	77.5	50.5	47.8	82.5	<b>89.3</b>	82.5	<u>86.7</u>

## 279 6 Related work

280 **Offline RL** [1] learns policies from static datasets without online interaction. Its core challenge is  
 281 distribution shift: the learned policy may choose out-of-distribution (OOD) actions with unreliable  
 282 value estimates. Prior work addresses this through conservative value estimation [8, 23], behavior  
 283 regularization [24, 9, 10], in-sample learning [7, 25, 26], and uncertainty-aware regularization [27, 28].  
 284 Several studies further show that offline RL performance depends strongly on implementation and  
 285 policy extraction design choices [29, 10, 30]. In contrast, we focus on test-time planning.

286 **Offline model-based planning** learns a dynamics model from static data and performs test-time  
 287 planning over candidate actions. Unlike standard offline RL, these methods adapt decisions online  
 288 through search. MBOP [13] is an early example, while later approaches improve robustness to  
 289 model error through uncertainty-aware planning. MOPP [14] removes uncertain branches, IQL-  
 290 TD-MPC [22] combines TD-MPC with offline RL objectives, and plug-and-play planners such as  
 291 LOOP [16] and RefPlan [17] integrate offline RL policies and value functions into planning.

292 **Offline model-based RL** uses learned dynamics models to improve policy optimization during  
 293 training rather than for test-time planning. Since policies can exploit model errors on OOD state-  
 294 action pairs, prior methods employ pessimistic rollouts [31–33], conservative value regularization [34],  
 295 and adversarial dynamics models [35]. Most related to our work, STEVE [36] adapts rollout horizons  
 296 using ensemble uncertainty estimates; we instead apply this principle to test-time planning.

## 297 7 Conclusions

298 We introduce Soft Horizon AggRegation for Planning (SHARP), which addresses horizon sensitivity  
 299 by dynamically weighting returns based on model uncertainty, eliminating the need for manually  
 300 tuned, state-invariant horizons. We further show that sophisticated policy extraction is often unnec-  
 301 essary, as a simple behavior cloning (BC) proposer can achieve comparable planning performance.  
 302 Combining these insights, we propose SHARP-BC, a robust plug-and-play offline planning method  
 303 that consistently outperforms existing baselines across diverse tasks and dataset settings. Overall, our  
 304 approach simplifies offline planning by reducing reliance on extensive hyperparameter tuning while  
 305 maintaining strong performance.

306 **References**

- 307 [1] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning:  
308 Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- 309 [2] Physical Intelligence, Ali Amin, Raichelle Aniceto, Ashwin Balakrishna, Kevin Black, Ken  
310 Conley, Grace Connors, James Darpinian, Karan Dhabalia, Jared DiCarlo, et al.  $\pi_{0.6}^*$ : a vla that  
311 learns from experience. *arXiv preprint arXiv:2511.14759*, 2025.
- 312 [3] Yevgen Chebotar, Quan Vuong, Karol Hausman, Fei Xia, Yao Lu, Alex Irpan, Aviral Kumar,  
313 Tianhe Yu, Alexander Herzog, Karl Pertsch, et al. Q-transformer: Scalable offline reinforcement  
314 learning via autoregressive q-functions. In *Conference on Robot Learning*, 2023.
- 315 [4] Tianyu Shi, Dong Chen, Kaian Chen, and Zhaojian Li. Offline reinforcement learning for  
316 autonomous driving with safety and exploration enhancement. *arXiv preprint arXiv:2110.07067*,  
317 2021.
- 318 [5] Antonio Guillen-Perez. From imitation to optimization: A comparative study of offline learning  
319 for autonomous driving. *arXiv preprint arXiv:2508.07029*, 2025.
- 320 [6] Shengpu Tang and Jenna Wiens. Model selection for offline reinforcement learning: Practical  
321 considerations for healthcare settings. In *Machine Learning for Healthcare Conference*, 2021.
- 322 [7] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit  
323 q-learning. In *International Conference on Learning Representations*, 2022.
- 324 [8] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for  
325 offline reinforcement learning. In *Advances in neural information processing systems*, 2020.
- 326 [9] Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning.  
327 In *Advances in neural information processing systems*, 2021.
- 328 [10] Denis Tarasov, Vladislav Kurenkov, Alexander Nikulin, and Sergey Kolesnikov. Revisiting  
329 the minimalist approach to offline reinforcement learning. In *Advances in Neural Information  
330 Processing Systems*, 2023.
- 331 [11] Dibya Ghosh, Anurag Ajay, Pulkit Agrawal, and Sergey Levine. Offline rl policies should be  
332 trained to be adaptive. In *International Conference on Machine Learning*, 2022.
- 333 [12] Xiong-Hui Chen, Yang Yu, Qingyang Li, Fan-Ming Luo, Zhiwei Qin, Wenjie Shang, and  
334 Jieping Ye. Offline model-based adaptable policy learning. In *Advances in Neural Information  
335 Processing Systems*, 2021.
- 336 [13] Arthur Argenson and Gabriel Dulac-Arnold. Model-based offline planning. In *International  
337 Conference on Learning Representations*, 2021.
- 338 [14] Xianyu Zhan, Xiangyu Zhu, and Haoran Xu. Model-based offline planning with trajec-  
339 tory pruning. In *Proceedings of the Thirty-First International Joint Conference on Artificial  
340 Intelligence, IJCAI-22*, 2022.
- 341 [15] Jacques Richalet, André Rault, Jean-Louis Testud, and Jean Papon. Model predictive heuristic  
342 control. *Automatica (journal of IFAC)*, 14, 1978.
- 343 [16] Harshit Sikchi, Wenxuan Zhou, and David Held. Learning off-policy with online planning. In  
344 *5th Annual Conference on Robot Learning*, 2021.
- 345 [17] Jihwan Jeong, Xiaoyu Wang, Jingmin Wang, Scott Sanner, and Pascal Poupart. Reflect-  
346 then-plan: Offline model-based planning through a doubly bayesian lens. In *Forty-second  
347 International Conference on Machine Learning*, 2025.
- 348 [18] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for  
349 deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- 350 [19] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression:  
351 Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.

- 352 [20] Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for  
353 flexible behavior synthesis. In *International Conference on Machine Learning*, 2022.
- 354 [21] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter  
355 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning  
356 via sequence modeling. In *Advances in Neural Information Processing Systems*, 2021.
- 357 [22] Rohan Chitnis, Yingchen Xu, Bobak Hashemi, Lucas Lehnert, Urun Dogan, Zheqing Zhu, and  
358 Olivier Delalleau. Iql-td-mpc: Implicit q-learning for hierarchical model predictive control. In  
359 *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024.
- 360 [23] Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor  
361 critic for offline reinforcement learning. In *International Conference on Machine Learning*,  
362 2022.
- 363 [24] Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement  
364 learning. *arXiv preprint arXiv:1911.11361*, 2019.
- 365 [25] Haoran Xu, Li Jiang, Jianxiong Li, Zhuoran Yang, Zhaoran Wang, Victor Wai Kin Chan,  
366 and Xianyuan Zhan. Offline rl with no ood actions: In-sample learning via implicit value  
367 regularization. In *International Conference on Learning Representations*, 2023.
- 368 [26] Divyansh Garg, Joey Hejna, Matthieu Geist, and Stefano Ermon. Extreme q-learning: Maxent  
369 RL without entropy. In *The Eleventh International Conference on Learning Representations*,  
370 2023.
- 371 [27] Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline re-  
372 inforcement learning with diversified q-ensemble. In *Advances in neural information processing*  
373 *systems*, 2021.
- 374 [28] Seyed Kamyar Seyed Ghasemipour, Shixiang Shane Gu, and Ofir Nachum. Why so pessimistic?  
375 estimating uncertainties for offline RL through ensembles, and why their independence matters.  
376 In *Advances in Neural Information Processing Systems*, 2022.
- 377 [29] Yuwei Fu, Di Wu, and Benoit Boulet. A closer look at offline RL agents. In *Advances in Neural*  
378 *Information Processing Systems*, 2022.
- 379 [30] Seohong Park, Kevin Frans, Sergey Levine, and Aviral Kumar. Is value learning really the  
380 main bottleneck in offline RL? In *The Thirty-eighth Annual Conference on Neural Information*  
381 *Processing Systems*, 2024.
- 382 [31] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea  
383 Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. In *Advances in neural*  
384 *information processing systems*, 2020.
- 385 [32] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel:  
386 Model-based offline reinforcement learning. In *Advances in neural information processing*  
387 *systems*, 2020.
- 388 [33] Cong Lu, Philip Ball, Jack Parker-Holder, Michael Osborne, and Stephen J. Roberts. Revisiting  
389 design choices in offline model based reinforcement learning. In *International Conference on*  
390 *Learning Representations*, 2022.
- 391 [34] Tianhe Yu, Aviral Kumar, Rafael Rafailov, Aravind Rajeswaran, Sergey Levine, and Chelsea  
392 Finn. Combo: Conservative offline model-based policy optimization. In *Advances in neural*  
393 *information processing systems*, 2021.
- 394 [35] Marc Rigter, Bruno Lacerda, and Nick Hawes. Rambo-rl: Robust adversarial model-based  
395 offline reinforcement learning. In *Advances in neural information processing systems*, 2022.
- 396 [36] Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample-  
397 efficient reinforcement learning with stochastic ensemble value expansion. In *Advances in*  
398 *Neural Information Processing Systems*, 2018.

- 399 [37] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based  
400 control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012.
- 401 [38] Soichiro Nishimori. Jax-corl: Clean single-file implementations of offline rl algorithms in jax.  
402 2024. URL <https://github.com/nissymori/JAX-CORL>.

404 **Appendix**

405 **Table of Contents**

---

407	<b>A Limitations</b>	<b>13</b>
408	<b>B Additional experimental results</b>	<b>13</b>
409	B.1 Best performance after online hyperparameter tuning . . . . .	13
410	B.2 Mean performance across all hyperparameter settings . . . . .	14
411	B.3 Action proposer analysis across dataset types . . . . .	14
412	<b>C Baselines and dataset descriptions</b>	<b>14</b>
413	C.1 Plug-and-play planning methods . . . . .	14
414	C.2 End-to-dnd planning baselines . . . . .	15
415	C.3 Tasks and dataset . . . . .	16
416	<b>D Offline RL implementation details</b>	<b>16</b>
417	<b>E Planning implementation details</b>	<b>16</b>
418	<b>F Experiments compute resources</b>	<b>17</b>

---

422 **A Limitations**

423 The primary goal of this work is to design a tuning-free, plug-and-play offline reinforcement learning  
424 algorithm. While our proposed method, SHARP-BC, effectively eliminates the need for online tuning  
425 of the planning horizon and simplifies the selection of an action proposer, certain limitations remain  
426 that provide avenues for future investigation:

427 **Sensitivity to action-selection temperature ( $\kappa$ ):** Although we have successfully addressed the  
428 challenges of horizon selection and action proposer design, our current method remains slightly  
429 sensitive to the hyperparameter  $\kappa$ . This sensitivity led to lower performance in some experimental  
430 results when evaluated under a strictly tuning-free protocol. Developing a mechanism to automatically  
431 adapt or eliminate the need for tuning  $\kappa$  is a critical area for future work.

432 **Hyperparameters in value function learning:** To move a significant step closer toward a truly  
433 end-to-end tuning-free model-based planning framework, it will be necessary to investigate the  
434 hyperparameters involved in training the underlying value functions (e.g.,  $Q$ -function learning).  
435 Currently, these still rely on standard offline RL training configurations, which may not be optimal  
436 for all planning scenarios. By addressing these remaining sensitivities, we aim to further simplify  
437 the deployment of offline model-based agents in practical, real-world applications where online  
438 evaluation is impossible.

439 **Focus on non-generative policies:** To simplify the scope of this work, we focus exclusively on  
440 non-generative policies in this paper. Nevertheless, generative models are capable of producing  
441 diverse action distributions and may therefore serve as effective action proposers. We leave the  
442 exploration of generative models in this context to future work.

443 **B Additional experimental results**

444 **B.1 Best performance after online hyperparameter tuning**

445 In the main paper, Table 1 evaluates robustness by averaging performance over all tested planning  
446 horizons  $H$  and action-selection temperatures  $\kappa$ . This reflects a setting where online hyperparameter  
447 selection is unavailable. Here, we also report the max performance after hyper search.

448 As shown in Table 3, SHARP still performs better than baselines by doing soft horizon aggregation.

Table 3: **Max planning performance.** We report the max normalized return across all evaluated settings for each planning method.

Environment	Dataset	IQL				CQL			
		Zero-shot	MBOP	LOOP	SHARP (Ours)	Zero-shot	MBOP	LOOP	SHARP (Ours)
HalfCheetah	R	6.8	21.4±0.6	23.2±0.8	<b>23.8±1.0</b>	4.4	5.7±1.4	6.7±1.2	<b>7.2±1.5</b>
	M	46.9	55.3±0.2	<u>58.4±0.4</u>	<b>60.1±0.3</b>	44.9	53.1±0.3	<u>55.8±0.4</u>	<b>57.8±0.5</b>
	MR	41.7	45.5±0.1	<u>48.9±0.2</u>	<b>52.5±0.6</b>	29.3	32.4±22.9	<u>33.5±23.6</u>	<b>34.2±24.3</b>
	ME	90.7	91.1±1.7	<b>101.5±1.0</b>	<u>93.4±0.4</u>	91.5	99.3±1.0	<b>107.7±0.3</b>	<u>107.0±0.3</u>
	FR	71.7	79.3±1.3	<u>83.6±1.1</u>	<b>83.7±1.2</b>	72.4	79.6±0.5	<b>83.4±0.6</b>	<u>82.5±0.2</u>
Hopper	R	7.6	<u>7.9±0.2</u>	<b>8.2±0.4</b>	<b>8.2±0.3</b>	4.3	<b>6.1±0.7</b>	<u>5.8±0.4</u>	5.0±0.6
	M	51.8	53.4±3.3	<u>73.8±3.7</u>	<b>80.4±6.7</b>	58.5	74.4±3.7	<u>82.7±3.6</u>	<b>84.9±2.3</b>
	MR	55.9	98.7±0.7	<b>99.1±0.2</b>	<u>98.8±0.3</u>	20.0	99.5±0.6	<u>100.1±0.7</u>	<b>100.2±0.6</b>
	ME	<b>70.9</b>	39.3±3.1	<u>41.7±8.9</u>	<u>36.7±2.4</u>	62.6	<u>111.3±0.1</u>	110.9±0.3	<b>111.4±0.3</b>
	FR	93.4	102.7±1.2	<b>106.0±0.1</b>	<u>104.8±0.3</u>	101.0	<u>105.4±0.7</u>	<b>106.2±0.2</b>	<b>106.2±0.5</b>
Walker2d	R	4.8	<u>5.7±0.7</u>	<b>6.5±0.6</b>	<b>6.5±0.6</b>	<b>0.1</b>	<b>0.1±0.0</b>	<b>0.1±0.0</b>	<b>0.1±0.0</b>
	M	83.5	87.0±2.0	<u>92.4±0.3</u>	<b>93.7±0.5</b>	73.5	84.2±0.3	<u>84.6±0.2</u>	<b>84.7±0.1</b>
	MR	83.9	93.1±1.3	<u>97.2±1.1</u>	<b>98.7±1.1</b>	74.2	86.4±0.4	<u>87.8±1.7</u>	<b>89.3±0.6</b>
	ME	<u>109.5</u>	108.9±0.1	<b>110.1±0.4</b>	109.4±0.2	108.6	<u>110.0±0.0</u>	<b>110.4±0.3</b>	<b>110.4±0.3</b>
	FR	97.9	99.9±1.8	<u>103.5±0.3</u>	<b>104.2±0.2</b>	91.1	95.1±0.9	<u>96.7±0.7</u>	<b>97.1±1.6</b>

Table 4: **Mean planning performance.** We report the mean normalized return across all evaluated settings for each planning method.

Environment	Dataset	IQL				CQL			
		Zero-shot	MBOP	LOOP	SHARP (Ours)	Zero-shot	MBOP	LOOP	SHARP (Ours)
HalfCheetah	R	6.8	17.7±0.7	<u>19.0±0.9</u>	<b>19.3±0.9</b>	4.4	4.2±1.1	5.2±1.2	<b>5.4±1.1</b>
	M	46.9	51.6±0.2	<u>53.6±0.1</u>	<b>54.5±0.2</b>	44.9	50.3±0.2	<u>52.4±0.3</u>	<b>53.4±0.3</b>
	MR	41.7	43.1±0.5	<u>46.4±0.2</u>	<b>47.9±0.2</b>	29.3	31.0±21.9	<u>32.3±22.8</u>	<b>33.0±23.4</b>
	ME	<u>90.7</u>	72.5±3.2	<b>92.7±1.2</b>	77.4±6.0	91.5	92.3±2.9	<u>99.3±1.3</u>	<b>102.3±1.3</b>
	FR	71.7	71.8±1.3	<u>74.8±1.0</u>	<b>80.6±1.1</b>	72.4	76.5±1.2	<u>78.2±1.5</u>	<b>81.0±0.7</b>
Hopper	R	7.6	<u>7.8±0.2</u>	<b>7.9±0.3</b>	<u>7.8±0.3</u>	4.3	<b>5.1±0.6</b>	<u>5.0±0.6</u>	4.5±0.9
	M	51.8	41.2±2.3	<u>56.4±4.0</u>	<b>58.9±6.2</b>	58.5	63.1±2.8	<u>69.1±3.4</u>	<b>71.9±3.1</b>
	MR	55.9	70.7±3.0	<u>87.9±1.2</u>	<b>88.2±1.0</b>	20.0	79.2±6.8	<u>84.2±6.4</u>	<b>93.3±2.5</b>
	ME	70.9	30.9±2.8	<b>34.7±3.5</b>	<u>31.5±3.6</u>	62.6	105.9±4.0	<b>107.9±2.8</b>	<u>106.0±4.4</u>
	FR	<b>93.4</b>	69.1±6.1	<u>91.6±2.9</u>	<u>91.4±3.5</u>	101.0	103.5±0.3	<u>104.3±0.5</u>	<b>104.6±0.5</b>
Walker2d	R	4.8	4.8±0.6	<u>5.1±0.5</u>	<b>5.5±0.6</b>	<b>0.1</b>	<b>0.1±0.0</b>	<b>0.1±0.0</b>	<b>0.1±0.0</b>
	M	83.5	80.7±3.0	<u>87.6±1.7</u>	<b>89.3±1.6</b>	73.5	81.7±1.2	<u>83.5±0.5</u>	<b>84.0±0.3</b>
	MR	83.9	82.5±3.5	<u>92.8±2.0</u>	<b>95.1±1.9</b>	74.2	79.7±3.3	<u>84.4±2.1</u>	<b>86.2±1.6</b>
	ME	<u>109.5</u>	108.2±0.2	<b>109.1±0.2</b>	108.6±0.2	108.6	<u>109.2±0.2</u>	<b>110.0±0.2</b>	<b>110.0±0.2</b>
	FR	97.9	96.5±1.5	<u>99.9±0.6</u>	<b>100.7±0.5</b>	91.1	93.5±0.9	<u>94.2±0.9</u>	<b>95.2±1.1</b>

449 **B.2 Mean performance across all hyperparameter settings**

450 In the paper, we report the mean planning performance across all hyperparameter settings. Here we  
 451 report the mean and standard deviation across 3 seeds in Table 4.

452 **B.3 Action proposer analysis across dataset types**

453 In Figure 3, we have shown that using BC as an action proposer yields comparable planning perfor-  
 454 mance. Here, we further compare different action proposers on varying dataset quality. The results  
 455 in Figure 4 indicate that the correlation between planning performance and zero-shot performance  
 456 occurs in low-quality datasets (Random). This highlights that when the dataset quality is sufficient,  
 457 using BC as an action proposer is the simplest option at no cost.

458 **C Baselines and dataset descriptions**

459 **C.1 Plug-and-play planning methods**

460 To leverage the knowledge encoded in offline RL policies  $\pi$  and value functions, plug-and-play  
 461 planning methods [13, 16, 17] incorporate action proposers  $u$  and value functions  $\hat{Q}$  learned via  
 462 offline RL into the planning process. Here, we discuss planning algorithms that can be used in a  
 463 plug-and-play style.

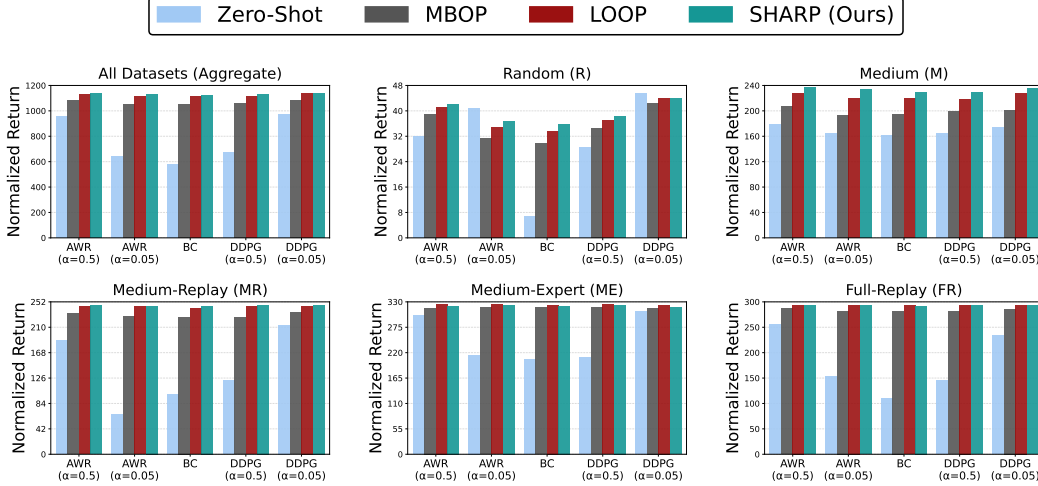


Figure 4: **Planning robustness over action proposers on varying dataset quality.**

464 **MBOP.** Model-Based Offline Planning (MBOP) [13], one of the first works to introduce model-based  
 465 planning to offline scenarios, is among the simplest planning methods. They train an ensemble of  
 466 dynamic models  $\{\hat{T}_e\}_{e=1}^E$ , a behavior cloning policy as an action proposer  $u$ , and use return-to-go as  
 467 the objective to optimize the value functions. The trajectory is rolled out as:

$$a_{h+1}^k \sim u(\cdot | \hat{s}_{h+1}^k), \quad \hat{s}_{h+1}^k \sim \hat{T}_e(\cdot | \hat{s}_h^k, a_h^k) \quad e = k \bmod E \quad (9)$$

468 with  $\hat{s}_0^k = s_0, a_0^k = a^k$ . The return is estimated as:

$$G^k = \sum_{h=0}^{H-1} \gamma^h \hat{r}(\hat{s}_h^k, a_h^k) + \gamma^H \hat{Q}(\hat{s}_H^k, a_H^k), \quad (10)$$

469 The action selection rule  $\mathcal{A}$  is defined as:

$$\mathcal{A}(\{\mathbf{a}_{0:H-1}^k, G^k\}_{k=1}^K) = \frac{\sum_{k=1}^K \exp(\kappa G^k) a_0^k}{\sum_{k=1}^K \exp(\kappa G^k)}, \quad (11)$$

470 where  $\kappa$  is a hyperparameter controlling how strongly to emphasize actions with larger estimated  
 471 return  $G^k$ .

472 We consider the plug-and-play version in this work.

473 **LOOP.** Learning Off-Policy with Online Planning (LOOP) [16] constrains the planned action  
 474 sequence using a KL-divergence regularization term to prevent the executed actions from deviating  
 475 excessively from the action proposer. In this work, we consider their proposed offline variant,  
 476 LOOP-offline, as a plug-and-play method for offline model-based planning. Unlike MBOP, LOOP  
 477 rolls out an ensemble of future states  $\hat{s}_{h,e,p}^k$ , where  $e$  indexes the dynamics model and  $p$  indexes the  
 478 particle, *i.e.*, different rollout samples generated from the same dynamics model. By performing  
 479 multiple rollouts across the dynamics ensembles, the method accounts for model uncertainty through  
 480 the standard deviation of predicted returns. The return function is estimated as:

$$G_{p,e}^k = \sum_{h=0}^{H-1} \gamma^h \hat{r}(\hat{s}_{h,e,p}^k, a_{h,e,p}^k) + \gamma^H \hat{Q}(\hat{s}_{H,e,p}^k, a_{H,e,p}^k), \quad (12)$$

$$G_e^k = \text{mean}_p(G_{e,p}^k) - \beta_{pess} \cdot \text{std}_p(G_{e,p}^k), \quad G^k = \text{mean}_e(G_e^k),$$

481 where  $\beta_{pess}$  controls the weight of the penalty. The action selection rule  $\mathcal{A}$  is the same one as in  
 482 MBOP.

## 483 C.2 End-to-dnd planning baselines

484 To provide a comprehensive comparison, we select methodologies that represent two paradigms in  
 485 offline reinforcement learning and planning: sequence-based and value-based approaches. Below, we  
 486 provide a brief overview of the core mechanics of each baseline.

487 Sequence-based planning methods reformulate the reinforcement learning and planning problem  
488 as a sequence modeling or generative task, typically leveraging generative models to co-learn the  
489 planning process and the dynamic model.

490 **Decision Transformer (DT) [21]:** Decision Transformer abstracts the offline RL problem into a  
491 conditional sequence modeling task. Rather than explicitly fitting value functions or computing  
492 temporal difference errors, DT feeds sequences of states, actions, and desired returns (returns-to-go)  
493 into a causal transformer architecture. By conditioning on a high target return, the model generates  
494 action sequences intended to achieve that outcome, effectively turning planning into autoregressive  
495 generation.

496 **Diffuser [20]:** Diffuser takes a fundamentally different generative approach by conceptualizing  
497 planning as a denoising diffusion process. It models the joint trajectory of states and actions,  
498 iteratively refining a sequence of purely Gaussian noise into a physically plausible and high-reward  
499 trajectory. Because the diffusion process operates simultaneously over the entire trajectory, it provides  
500 a sequence-level perspective on long-horizon planning.

501 Value-based planning methods rely on learning a separate dynamic model and leveraging it for  
502 test-time planning; SHARP falls into this category.

503 **MBOP [13]** is discussed in Appendix C.1. We use its original version as the end-to-end baseline.

504 **IQL-TD-MPC [22]:** This approach integrates the conservative value estimation of Implicit Q-  
505 Learning (IQL) with the latent-space planning capabilities of TD-MPC (Temporal Difference Learning  
506 for Model Predictive Control). By leveraging IQL, the method obtains robust value estimates from  
507 offline data without querying out-of-distribution actions. These learned representations and value  
508 functions are then used to perform efficient, derivative-free trajectory optimization in a learned latent  
509 space during inference.

### 510 C.3 Tasks and dataset

511 We evaluate our method on three standard MuJoCo [37] locomotion environments from D4RL [18]:  
512 *halfcheetah*, *hopper*, and *walker2d*. These continuous-control tasks are commonly used in offline RL  
513 and provide long-horizon decision-making problems for evaluating model-based planning.

514 For each environment, we use five dataset types: **random** (R), collected by an untrained policy;  
515 **medium** (M), collected by a partially trained policy; **medium-replay** (MR), containing the replay  
516 buffer accumulated during training to medium performance; **medium-expert** (ME), mixing medium-  
517 level and expert demonstrations; and **full-replay** (FR), containing the full replay buffer from scratch  
518 to expert performance. These 15 task–dataset combinations span different levels of data quality,  
519 behavioral coverage, and state-action distribution shift, allowing us to assess whether planning  
520 methods remain robust across diverse offline data regimes.

## 521 D Offline RL implementation details

522 We train both IQL and CQL following Jax-CORL’s [38] codebase and hyperparameters. For policy-  
523 extraction experiments, all policies (BC, AWR, and DDPG+BC) follow IQL’s policy implementation.  
524 The hyperparameters are reported in Table 5.

525 For all zero-shot performance, the policy  $\pi$  will be set to deterministic ( $\tau=0$ ). When the policy  $\pi$  is  
526 used as an action proposer  $u$ , it will be set to stochastic ( $\tau=1$ ).

## 527 E Planning implementation details

528 We train dynamic models following the setting of previous work [16]. The hyperparameters are  
529 provided in Table 6.

530 For all the plug-and-play planning methods, we set  $K=32$ ,  $\sigma=0.01$ , and action mixture  $\beta=0$ . For  
531 LOOP and SHARP, we set particle  $P=4$ . We set  $\beta_{\text{pess}}=0.1$  for LOOP.

Table 5: **Offline RL training hyperparameters.** We list the IQL, CQL, and policy extraction hyperparameters used in all experiments, following the Jax-CORL implementation.

Method	Hyperparameter	Value
IQL	Steps	1M
	Batch size	256
	Target update rate	5e-3
	Actor learning rate	3e-4
	Value learning rate	3e-4
	Critic learning rate	3e-4
	Learning rate decay	Cosine
	Expectile ( $\tau$ )	0.7 (0.5 for hopper-me, walker2d-mr)
	AWR alpha ( $\beta$ )	3.0 (6.0 for hopper-me)
	Actor hidden dim	[256, 256]
	Critic/Value hidden dim	[256, 256]
CQL	Steps	1M
	Batch size	256
	Target update rate	5e-3
	Actor lr	1e-4
	Critic lr	3e-4
	Target entropy	$-1 \times \text{action\_dim}$
	Actor hidden dim	[256, 256, 256]
	Critic hidden dim	[256, 256, 256]
	Lagrange	False
	CQL alpha	10
	CQL num action	10
BC/AWR/DDPG+BC	Steps	1M
	Batch size	256
	Learning rate	3e-4
	Actor hidden dim	[256, 256]
	Policy type	Stochastic

532 **F Experiments compute resources**

533 We performed all experiments on the workstations listed in Table 7.

534

Table 6: **Hyperparameters for learned dynamics and reward models.** We use the same network architecture and training configuration across all MuJoCo tasks. The dynamics model is trained as a stochastic ensemble, while the reward model is trained as a deterministic predictor.

Model	Hyperparameter	Value
Dynamic Model	Model type	stochastic ensemble 7
	Hidden dimensions	[200, 200, 200, 200]
	Learning rate	1e-3
	Steps	1M
	Batch size	1024
	Validation ratio	0.01
Reward Model	Model type	deterministic 1
	Hidden dimensions	[200, 200, 200, 200]
	Learning rate	1e-3
	Steps	1M
	Batch size	1024
	Validation ratio	0.01

Table 7: **Computational resources used.** We report the hardware used for training offline RL agents, learning dynamics and reward models, and running planning evaluations.

Workstation	CPU	GPU	RAM
Workstation 1	Intel Xeon W-2255	NVIDIA GeForce RTX 3080 Ti	125 GiB
Workstation 2	Intel Xeon W-2255	NVIDIA GeForce RTX 4070 Ti×2	125 GiB
Workstation 3	Intel Xeon Platinum 8480+	NVIDIA H100x2	503 GiB