

# COPlanner: Plan to Roll Out Conservatively but to Explore Optimistically for Model-Based RL

Anonymous Author(s)

Affiliation

Address

email

**Abstract:** Dyna-style model-based reinforcement learning contains two phases: model rollouts to generate sample for policy learning and real environment exploration using current policy for dynamics model learning. However, due to the complex real-world environment, it is inevitable to learn an imperfect dynamics model with model prediction error, which can further mislead policy learning and result in sub-optimal solutions. In this paper, we propose COPlanner, a planning-driven framework for model-based methods to address the inaccurately learned dynamics model problem with conservative model rollouts and optimistic environment exploration. COPlanner leverages an uncertainty-aware policy-guided model predictive control (UP-MPC) component to plan for multi-step uncertainty estimation. This estimated uncertainty then serves as a penalty during model rollouts and as a bonus during real environment exploration respectively, to choose actions. Consequently, COPlanner can avoid model uncertain regions through conservative model rollouts, thereby alleviating the influence of model error. Simultaneously, it explores high-reward model uncertain regions to reduce model error actively through optimistic real environment exploration. COPlanner is a plug-and-play framework that can be applied to any dyna-style model-based methods. Experimental results on a series of proprioceptive and visual continuous control tasks demonstrate that both sample efficiency and asymptotic performance of strong model-based methods are significantly improved combined with COPlanner.

**Keywords:** Model-based RL, Model prediction error, Uncertainty-based Planning

## 1 Introduction

Model-Based Reinforcement Learning (MBRL) has emerged as a promising approach to improve the sample efficiency of model-free RL methods. Most MBRL methods contain two phases that are alternated during training: 1) the first phase where the agent interacts with the real environment using a policy to obtain samples for dynamics model learning; 2) the second phase where the learned dynamics model rolls out to generate massive samples for updating the policy. Consequently, learning an accurate dynamics model is critical as the model-generated samples with high bias can mislead the policy learning [3, 33].

However, dynamics model errors are inevitable due to the complex real-world environment. Existing methods try to avoid model errors in two main ways. 1) Design different mechanisms such as filtering out error-prone samples to mitigate the influence of model errors after model rollouts [1, 35, 20, 34]. 2) Actively reduce model errors during real environment interaction through uncertainty-guided exploration [28, 24, 25, 18]. While both categories of methods have achieved advancements, each

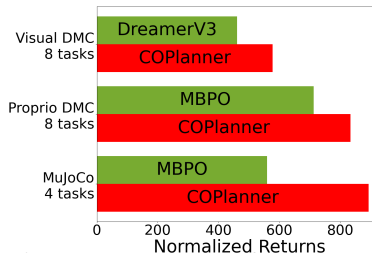


Figure 1: Mean performance of COPlanner compared with base-lines across 3 diverse benchmarks.

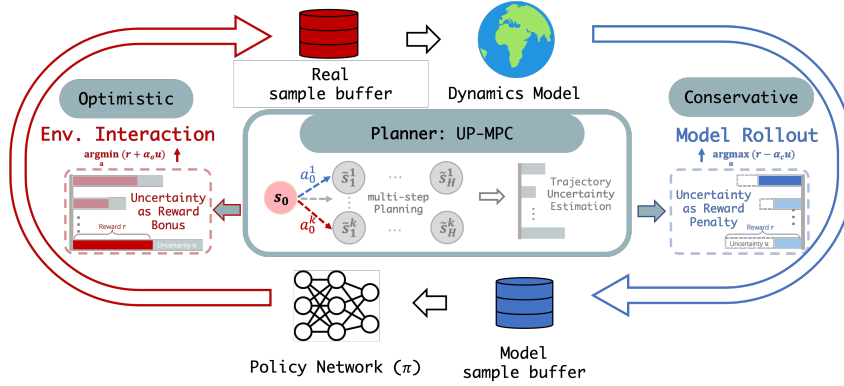


Figure 2: COPlanner Framework. The most essential part of COPlanner is the *Uncertainty-aware Policy-Guided MPC (UP-MPC)* phase in which we plan trajectories of length  $H$ , according to the learned dynamics model and learned policy network  $\pi$ , to select the action with highest trajectory reward. This UP-MPC phase is implemented differently for the two different purposes: *environment exploration* v.s. *dynamics model rollouts*. In *environment exploration*, trajectory reward has an uncertainty bonus term to encourage exploring uncertain regions in the environment. In *dynamics model rollouts*, trajectory reward, on the contrary, has an uncertainty penalty term to encourage policy learning on confident regions of the learned dynamics model.

39 comes with its own set of limitations. For the first category, although these approaches are shown to  
 40 be empirically effective, they primarily concentrate on estimating uncertainty at the current step, often  
 41 neglecting the long-term implications that present samples might have on model rollouts. Moreover,  
 42 post-processing samples after model rollouts can compromise rollout efficiency as many model-  
 43 generated samples are discarded or down-weighted. As for the second category, it is intrinsically  
 44 challenging to achieve low model error and high long-term reward without sacrificing the sample  
 45 efficiency by learning exploration policies.

46 To tackle the aforementioned limitations, we introduce a novel framework, COPlanner, which  
 47 mitigates the model errors from two aspects: 1) avoid being misled by the existing model errors  
 48 via conservative model rollouts, and 2) keep reducing the model error via optimistic environment  
 49 exploration. The two aspects are achieved simultaneously by a novel uncertainty-aware multi-step  
 50 planning method, which requires no extra exploration policy training nor additional samples, resulting  
 51 in stable policy updates and high sample efficiency. COPlanner is structured around three core  
 52 components: *the Planner*, *conservative model rollouts*, and *optimistic environment exploration*. In  
 53 the Planner, we employ an *Uncertainty-aware Policy-guided Model Predictive Control (UP-MPC)*  
 54 to forecast future trajectories in terms of selecting actions and to estimate the long-term uncertainty  
 55 associated with each action. As shown in Figure 2, this long-term uncertainty serves as dual roles. In  
 56 the model rollouts phase, the uncertainty acts as a penalty on the total planning trajectory, guiding the  
 57 selection of conservative actions. Conversely, during the model learning phase, it serves as a bonus  
 58 on the total planning trajectory, steering towards optimistic actions for environment exploration.

59 Compared to previous methods, COPlanner has the following advantages: **(a)** COPlanner has **higher**  
 60 **exploration efficiency**, as it focuses on investigating high-reward uncertain regions to broaden the  
 61 dynamics model, thereby preventing unnecessary excessive exploration of areas with low rewards. **(b)**  
 62 COPlanner has **higher model-generated sample utilization rate**. Through planning for multi-step  
 63 model uncertainty estimation, COPlanner can prevent model rolled out trajectories from falling into  
 64 uncertain areas, thereby avoiding model errors before model rollouts and improving the utility of  
 65 model generated samples. **(c)** COPlanner enjoys an **unified policy framework**. Unlike previous  
 66 methods [25, 18] that require training two separate policies for different usage, COPlanner only  
 67 requires training a single policy and we only change the way model-based planning is utilized,  
 68 thus improving training efficiency and resolving potential policy distribution mismatches. **(d)**  
 69 COPlanner ensures **undistracted policy optimization**. Notably, COPlanner diverges from existing  
 70 approaches by not using long-term uncertainty as an intrinsic reward. Instead, the policy’s objective  
 71 remains focused on maximizing environmental rewards, thereby avoiding the introduction of spurious  
 72 behaviors due to model uncertainty.

73 **Summary of Contributions:** **(1)** We introduce COPlanner framework which can mitigate the  
 74 influence of model errors during model rollouts and explore the environment to actively reduce

75 model errors simultaneously by leveraging our proposed uncertainty-aware policy-guided MPC. **(2)**  
 76 COPlanner is a plug-and-play framework that can be applicable to any dyna-style MBRL method.  
 77 **(3)** After being integrated with other MBRL baseline methods, COPlanner improves the sample  
 78 efficiency of these baselines by nearly double. **(4)** Besides, COPlanner also significantly improves  
 79 the performance on a suite of proprioceptive and visual control tasks compared with other MBRL  
 80 baseline methods (16.9% on proprioceptive DMC, 59.7% on MuJoCo, and 23.9% on visual DMC).

## 81 2 Preliminaries

82 **Model-based reinforcement learning.** We consider a Markov Decision Process (MDP) defined  
 83 by the tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \rho_0, r, \gamma)$ , where  $\mathcal{S}$  and  $\mathcal{A}$  are the state space and action space respectively,  
 84  $\mathcal{T}(s'|s, a)$  is the transition dynamics,  $\rho_0$  is the initial state distribution,  $r(s, a)$  is the reward function  
 85 and  $\gamma$  is the discount factor. In model-based RL, the transition dynamics  $T$  in the real world is  
 86 unknown, and we aim to construct a model  $\hat{T}(s'|s, a)$  of transition dynamics and use it to find an  
 87 optimal policy  $\pi$  which can maximize the expected sum of discounted rewards,

$$\pi = \operatorname{argmax}_{\pi} \mathbb{E}_{s_t \sim \hat{T}(\cdot | s_{t-1}, a_{t-1})} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]. \quad (1)$$

88 **Model predictive control.** Model predictive control (MPC) has a long history in robotics and control  
 89 systems [5, 22]. MPC find the optimal action through trajectory optimization. Specifically, given the  
 90 transition dynamics  $T$  in the real world, the agent obtains a local solution at each step  $t$  by estimating  
 91 optimal actions over a finite horizon  $H$  (i.e., from  $t$  to  $t + H$ ) and executing the first action  $a_t$  from  
 92 the computed optimal sequence at time step  $t$ :

$$a_t = \operatorname{argmax}_{a_{t:t+H}} \mathbb{E} \left[ \sum_{i=t}^H \gamma^i r(s_i, a_i) \right], s_i \sim T(\cdot | s_{i-1}, a_{i-1}), \quad (2)$$

93 where  $\gamma$  is typically set to 1. In model-based control methods, the transition dynamics  $T$  is simulated  
 94 by the learned dynamics model  $\hat{T}$  [2, 30, 11].

## 95 3 The COPlanner Framework

96 In this section, we will introduce COPlanner framework. COPlanner consists of three components:  
 97 the Planner, conservative model rollouts, and optimistic environment exploration. Within the Planner,  
 98 we propose using an Uncertainty-aware Policy-guided MPC to predict potential future trajectories  
 99 when selecting different actions under the current state and estimate the long-term uncertainty  
 100 associated with each action, which will be introduced in Sec 3.1. Depending on the phase, this  
 101 long-term uncertainty is used to further guide the selection of conservative actions for policy learning  
 102 or optimistic actions for environment exploration which will be introduced in Sec 3.2 and Sec 3.3.

### 103 3.1 “The Planner”: Uncertainty-Aware Policy-Guided MPC

104 In this section, we present the core part of our pro-  
 105 posed framework which is called Uncertainty-aware  
 106 Policy-guided MPC (UP-MPC). Inspired by MPC,  
 107 we apply the random shooting method [23] to intro-  
 108 duce a long-term vision. Specifically, given the cur-  
 109 rent state  $s_t$ , before each interaction with the model  
 110 or real environment, we first generate an action candidate set containing  $K$  actions using the policy:  
 111  $\mathbf{a}_t = \{a_t^{(1)}, a_t^{(2)}, \dots, a_t^{(k)}\}$ . Then, for each action candidate, we perform  $H_p$ -step planning and  
 112 calculate the reward  $r$ , and model uncertainty  $u$  for each step. Finally, we select the action according  
 113 to accumulated reward and model uncertainty, (to interact with the learned dynamics for the model  
 114 rollouts or to interact with the environment for the model learning), as will be discussed in details in  
 115 Sec 3.2 and Sec 3.3.

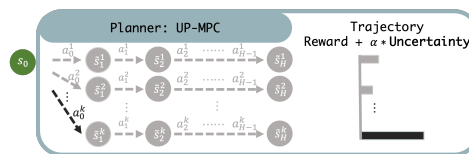


Figure 3: The Planner.

116 Incorporating model uncertainty is crucial for action selection to compensate for model error. As  
 117 illustrated in Algorithm 1, we calculate the model uncertainty  $u$  through the model disagreement [21]  
 118 method. Model disagreement is closely related to model learning and is currently the most common  
 119 way to estimate model uncertainty in MBRL [35, 14, 20, 25, 34, 18]. We train a dynamics model  
 120 ensemble  $\hat{T}_\theta = \{\hat{T}_\theta^{(1)}, \hat{T}_\theta^{(2)}, \dots, \hat{T}_\theta^{(n)}\}$  to predict the next state given the current state-action pair  
 121  $(s_t, a_t)$  as input. Utilizing the ensemble, we approximate the model uncertainty by calculating the  
 122 variance over predicted states of the different ensemble members. This estimation closely represents  
 123 the expected information gain [21]:

$$u(s_t, a_t) = \frac{1}{N-1} \sum_n (\hat{T}_\theta^{(n)}(s_t, a_t) - \mu')^2, \quad \mu' = \frac{1}{N} \sum_n \hat{T}_\theta^{(n)}(s_t, a_t). \quad (3)$$

124 See Figure 3 for the illustration of the process. The pseudocode for the Planner, i.e., the UP-MPC  
 125 process, is summarized in Algorithm 1.  
 126

---

**Algorithm 1** The Planner: **UP-MPC**  $(\pi_\phi, s, \hat{T}_\theta, K, H, \alpha)$

---

**Require:** Policy  $\pi_\phi$ , State  $s$ , learned dynamics model  $\hat{T}_\theta$ , number of candidates actions  $K$ , planning horizon  $H_p$ , optimistic/conservative parameter  $\alpha$

- 1: Initialize  $R^{(k)} = 0$  for  $k = 1, \dots, K$ ,  $s_0^{(k)} = s$  for  $k = 1, \dots, K$
  - 2: **for**  $k = 1$  to  $K$  **do**
  - 3:   **for**  $t = 0$  to  $H_p - 1$  **do**
  - 4:     Sample  $a_t^{(k)} \sim \pi_\phi(\cdot | s_t^{(k)})$
  - 5:     Rollout dynamics model  $r_t^{(k)} = \hat{R}(\cdot | s_t^{(k)}, a_t^{(k)})$ ,  $s_{t+1}^{(k)} \sim \hat{T}_\theta(\cdot | s_t^{(k)}, a_t^{(k)})$
  - 6:     Compute model uncertainty  $u_t^{(k)}$  according to Eq. 3
  - 7:      $R^{(k)} = R^{(k)} + r_t^{(k)} + \alpha u_t^{(k)}$
  - 8:   Select  $k^* = \arg \max_{k=1, \dots, K} R^{(k)}$
  - 9: **return**  $a_0^{(k^*)}$
- 

127 Although in Algorithm 1 model uncertainty  $u$  is implemented through model disagreement, our  
 128 proposed UP-MPC is a generic framework, any method for calculating intrinsic rewards to encourage  
 129 exploration can be embedded into our framework for computing  $u$ . In Appendix D.6 we provide an  
 130 ablation study of uncertainty estimation methods to further illustrate this point.

### 131 3.2 Conservative model rollouts

132 In model-based RL, due to the limited samples available for model learning, model prediction errors  
 133 are inevitable. If a policy is trained using model-generated samples with a large error, these samples  
 134 will not provide correct gradient and may mislead the policy update. Previous methods estimate the  
 135 model uncertainty of each sample after generation and re-weight or discarded samples with high  
 136 uncertainty. However, re-weighting samples based on uncertainty still leads to samples with high  
 137 uncertainty participating in the policy learning process, while filtering requires manually setting  
 138 an uncertainty threshold, and determining the optimal threshold is difficult. Discarding too many  
 139 samples can result in inefficient rollouts.

140 We apply our Planner to plan for maximizing the future reward while minimizing the model uncer-  
 141 tainty during model rollouts before executing the action. After calculating the reward and model  
 142 uncertainty for the  $H_p$ -step trajectories of  $K$  action candidates (line 5 and 6 in Algorithm 1), we  
 143 replace  $\alpha = -\alpha_c$ , for a positive  $\alpha_c > 0$  at line 7 in Algorithm 1. Mathematically, we select the action  
 144 according to Eq. 4 to interact with the model for model rollouts:

$$a = \operatorname{argmax}_{a_t \in \mathcal{A}_t} \left[ r(s_t, a_t) + \sum_{i=1}^{H_p} r(\hat{s}_{t+i}, \pi(\hat{s}_{t+i})) - \alpha_c \sum_{i=1}^{H_p} u(\hat{s}_{t+i}, \pi(\hat{s}_{t+i})) \right], \hat{s}_{t+i} \sim \hat{T}(\cdot | \hat{s}_{t+i-1}, a_{t+i-1}). \quad (4)$$

145 The negative  $-\alpha_c$  is a coefficient that adds the model uncertainty as a penalty term to the trajectory  
 146 total reward. By employing this approach, we can prevent model rollout trajectories from falling into  
 147 model-uncertain regions while obtaining samples with higher rewards.



148 **3.3 Optimistic environment exploration**

149 In addition to model rollouts, another crucial part of MBRL is interacting with the real environment  
 150 to obtain samples to improve the dynamics model. Since the main purpose of MBRL is to improve  
 151 sample efficiency, we should acquire more meaningful samples for improving the dynamics model  
 152 within a limited number of interactions. Therefore, unlike previous methods that merely aimed to  
 153 thoroughly explore the environment to obtain a comprehensive model [28, 24, 25], we do not expect  
 154 the dynamics model to learn all samples in the state space. This is because many low-reward samples  
 155 do not contribute to policy improvement. Instead, we hope to obtain samples with both high rewards  
 156 and high model uncertainty to sufficiently expand the model and reduce model uncertainty.

157 Similar to model rollouts, we also employ our Planner in the process of selecting actions when  
 158 interacting with the environment. However, the difference lies in that we replace  $\alpha = \alpha_o$ , for a  
 159 positive  $\alpha_o > 0$  at line 7 in Algorithm 1. Mathematically, we choose the action with both high  
 160 cumulative rewards and model uncertainty according to Eq. 5, which is a symmetric form of Eq. 4.  
 161  $\alpha_o$  is a hyperparameter to balance the reward and exploration. Such an action can guide the trajectory  
 162 towards regions with high rewards and model uncertainty in the real environment, thereby effectively  
 163 expanding the learned dynamics model.

$$a = \operatorname{argmax}_{a_t \in \mathcal{A}_t} \left[ r(s_t, a_t) + \sum_{i=1}^{H_p} r(\hat{s}_{t+i}, \pi(\hat{s}_{t+i})) + \alpha_o \sum_{i=1}^{H_p} u(\hat{s}_{t+i}, \pi(\hat{s}_{t+i})) \right], \hat{s}_{t+i} \sim \hat{T}(\cdot | \hat{s}_{t+i-1}, a_{t+i-1}) \quad (5)$$

164 In summary, by simultaneously using conservative model rollouts and optimistic environment ex-  
 165 ploration, COPlanner effectively alleviates the model error problem in MBRL. As we will show in  
 166 Section 5, this is of great help in improving the sample efficiency and performance. The pseudocode  
 167 of COPlanner is shown in Algorithm 2, and a more detailed figure is shown in Appendix A. Very  
 168 importantly, COPlanner achieves both conservative model rollouts and optimistic environment ex-  
 169 ploration using a single policy. Different from prior exploration methods, the policy that COPlanner  
 170 learns does not have to be an “exploration” policy which is inevitably suboptimal.

---

**Algorithm 2** Main Algorithm: COPlanner

---

**Require:** Interaction epochs  $I$ , rollout horizon  $H_r$ , planning horizon  $H_p$ , number of candidates  
 actions  $K$ , conservative rate  $\alpha_c$ , optimistic rate  $\alpha_o$

- 1: Initialize policy  $\pi_\phi$ , dynamics model  $\hat{T}$ , real sample buffer  $\mathcal{D}_e$ , model sample buffer  $\mathcal{D}_m$
- 2: **for**  $I$  epochs **do**
- 3:   **while** not Done **do**
- 4:     Select action  $a_t = \text{UP-MPC}(\pi_\phi, s_t, \hat{T}_\theta, K, H_p, \alpha_o)$
- 5:     Execute in real environment, add  $(s_t, a_t, r_t, s_{t+1})$  to  $\mathcal{D}_e$
- 6:     Train dynamics model  $\hat{T}_\theta$  with  $\mathcal{D}_e$
- 7:     **for**  $M$  model rollouts **do**
- 8:       Sample initial states from real sample buffer  $\mathcal{D}_e$
- 9:       **for**  $h = 0$  to  $H_r$  **do**
- 10:          Select action  $\hat{a}_{t+h} = \text{UP-MPC}(\pi_\phi, \hat{s}_h, \hat{T}_\theta, K, H_p, -\alpha_c)$
- 11:          Rollout learned dynamics model and add to  $\mathcal{D}_m$
- 12:       Update current policy  $\pi_\phi$  with  $\mathcal{D}_m$

---

171 **4 Related work**

172 **Mitigating model error by improving rollout strategies.** Prior methods primarily focus on using  
 173 dynamics model ensembles [15, 2] to assess model uncertainty of samples after they were generated  
 174 by the model, and then apply weighting techniques [1, 34], penalties [14, 35] or filtering [20, 32] to  
 175 those high uncertainty samples to mitigate the influence of model error. These methods only quantify  
 176 uncertainty after generating the samples and since their uncertainty metrics are based on the current  
 177 step and are myopic, these metrics can not evaluate the potential influence of the current sample  
 178 on future trajectories. Therefore, they fail to prevent the trajectories, which is generated through

179 model rollout on the current policy, from entering high uncertainty regions, eventually leading to  
180 a failed policy update. Wu et al. [33] proposed Plan to Predict (P2P), which reverses the roles of  
181 the model and policy during model learning to learn an uncertainty-foreseeing model, aiming to  
182 avoid model uncertain regions during model rollouts. Combined with MPC, their method achieved  
183 promising results. However, their approach lacks effective exploration of the environment. Branched  
184 rollout [13] and bidirectional rollout [16] take advantage of small model errors in the early stages  
185 of rollouts and uses shorter rollout horizons to avoid model errors, but these approaches limit the  
186 planning capabilities of the learned dynamics model. Besides, different model learning objectives  
187 [27, 4, 31, 37] are designed to solve objective mismatch [17] in model-based RL and further mitigate  
188 model error during model rollouts.

189 **Reducing model error by improving environment exploration.** Another approach to mitigate  
190 model error is to expand the dynamics model by obtaining more diverse samples through exploration  
191 during interactions with the environment. However, previous methods mostly focused on pure  
192 exploration, i.e., how to make the dynamics model learn more comprehensively [28, 24, 25, 18, 12].  
193 In complex environments, thoroughly exploring the entire environment is very sample-inefficient  
194 and not practical in real-world applications. Moreover, using pure exploration to expand the model  
195 may lead to the discovery of many low-reward samples (e.g., different ways an agent may fall in  
196 MuJoCo environment [29]), which are not very useful for policy learning. Mendonca et al. [18]  
197 proposed Latent Explorer Achiever (LEXA) which involves a explorer for exploring the environment  
198 and one achiever for solving diverse tasks based on collected samples, but the explorer and achiever  
199 may experience policy distribution shift under specific single-task settings, causing the achiever to  
200 potentially not converge to the optimal solution.

201 **Mitigating model error from both sides.** One most relevant work is Model-Ensemble Exploration  
202 and Exploitation (MEEE) [34] which simultaneously expands the dynamics model and reduces the  
203 impact of model error during model rollouts. During the rollout process, it uses uncertainty to weight  
204 the loss calculated for each sample to update the policy and the critic. Before interacting with the  
205 environment, they first generate  $k$  action candidates and then select the action with the highest sum of  
206 Q-value and one-step model uncertainty to execute. However, as we mentioned earlier, weighting  
207 samples cannot fundamentally prevent the impact of model errors on policy learning, and it may  
208 still mislead policy updates. Moreover, since the one-step prediction error of dynamics models  
209 is often small [20], relying only on the sum of Q-values and one-step model uncertainty may not  
210 effectively differentiate action candidates. As a result, samples collected during interactions with the  
211 environment might not efficiently expand the model.

## 212 5 Experiment

213 In this section, we combine COPlanner with strong MBRL baseline methods and conduct experi-  
214 ments on both proprioceptive control environments and visual control environments to demonstrate  
215 the effectiveness of our method. Due to space constraints, further discussions about the method and  
216 ablation studies can be found in Appendix D.

### 217 5.1 Experiment on proprioceptive control tasks

218 **Baselines:** In this section, we conduct experiments to demonstrate the effectiveness of COPlanner  
219 on proprioceptive control MBRL methods. We combine COPlanner with MBPO [13], the most  
220 classic method in proprioceptive control dyna-style MBRL, and we name the combined method  
221 as **COPlanner-MBPO**. The implementation details can be found in Appendix B. Consequently,  
222 **MBPO** naturally becomes one of our baselines. The other two baselines are **P2P-MPC** [33] and  
223 **MEEE** [34]. These two methods also aim to mitigate the impact of model errors in model-based RL.  
224 More details of P2P-MPC and MEEE can be found in Section 4. We also provide comparison with  
225 more proprioceptive control MBRL methods in Appendix D.1.

226 **Environment and hyperparameter settings:** We conduct experiments on 8 proprioceptive con-  
227 tinuous control tasks of DeepMind Control (DMC) and 4 proprioceptive control tasks of MuJoCo.

228 MBPO trains an ensemble of 7 networks as the dynamics model while using the Soft Actor-Critic  
 229 (SAC) as the policy network. In COPlanner-MBPO, we adopt the setting of MBPO and directly use  
 230 the dynamics model ensemble to calculate model uncertainty for action selection in Policy-Guided  
 231 MPC. For hyperparameter setting, we set optimistic rate  $\alpha_o$  to be 1, conservative rate  $\alpha_c$  to be 2 in  
 232 most tasks. We set action candidate number  $K$  and planning horizon  $H_p$  equal to 5 in all tasks. The  
 233 specific setting are shown in the Appendix C.1.

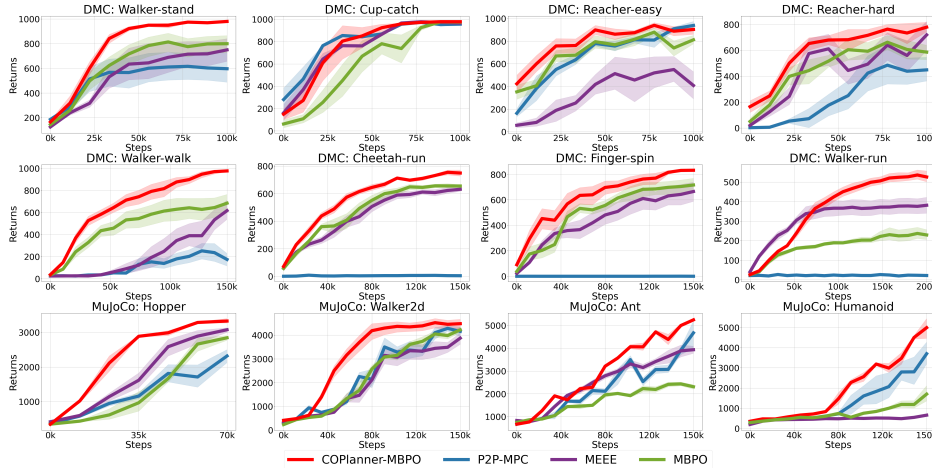


Figure 4: Experiment results of COPlanner-MBPO and other three baselines on proprioceptive control environments. The curves in the first eight figures originate from DM Control tasks, while those in the last four are from MuJoCo tasks. The results are averaged over 8 random seeds, and shaded regions correspond to the 95% confidence interval among seeds. During evaluation, for each seed of each method, we test for up to 1000 steps in the test environment and perform 10 evaluations to obtain an average value. The evaluation interval is every 1000 environment steps.

234 COPlanner **significantly improves the sample efficiency and performance of MBPO**: Through  
 235 the results in Figure 4 we can find that both sample efficiency and performance of MBPO have  
 236 a significant improvement after combining COPlanner. **(a) Sample efficiency**: In proprioceptive  
 237 control DMC, the sample efficiency is improved by 40% on average compared to MBPO. For  
 238 example, in the Walker-walk task, MBPO requires 100k steps for the performance to reach 700,  
 239 while COPlanner-MBPO only needs approximately 60k steps. In more complex MuJoCo tasks, the  
 240 improvement brought by COPlanner is even more significant. Compared to MBPO, the sample  
 241 efficiency of COPlanner-MBPO has almost doubled. **(b) Performance**: From the performance  
 242 perspective, as shown in Figure 1, the performance of MBPO has improved by 16.9% after combining  
 243 COPlanner. Moreover, it is worth noting that our method successfully solves the Walker-run task,  
 244 which MBPO fails to address, further demonstrating the effectiveness of our proposed framework. In  
 245 MuJoCo tasks, the average performance at 150k environment steps has increased by 59.7%. Besides,  
 246 COPlanner-MBPO also outperforms other two baselines.

## 247 5.2 Experiment on visual control tasks

248 **Baselines**: We conduct experiments to demonstrate the effectiveness of our proposed framework on  
 249 visual control environments. We integrate our algorithm with DreamerV3 [10], the state-of-the-art  
 250 Dyna-style model-based RL approach recently introduced for visual control. The implementation  
 251 details can be found in Appendix B. We choose LEXA [18] as our another baseline. LEXA uses  
 252 Plan2Explore [25] as intrinsic reward to explore the environment and learn a world model, then using  
 253 this model to train a policy to solve diverse tasks such as goal achieving. Here we adopt LEXA  
 254 on DreamerV3 to address continuous control tasks and name it as LEXA-DreamerV3. Since pure  
 255 exploration base on Plan2Explore is sample inefficient for model learning when solving specific tasks,  
 256 we use the real reward provided by environment as extrinsic reward and add it to intrinsic reward  
 257 provided by Plan2Explore to train the explorer. We call this baseline LEXA-reward-DreamerV3.

258 We also provide comparison with more visual control MBRL methods including TDMPC [11] and  
 259 PlaNet [8] in Appendix D.2.

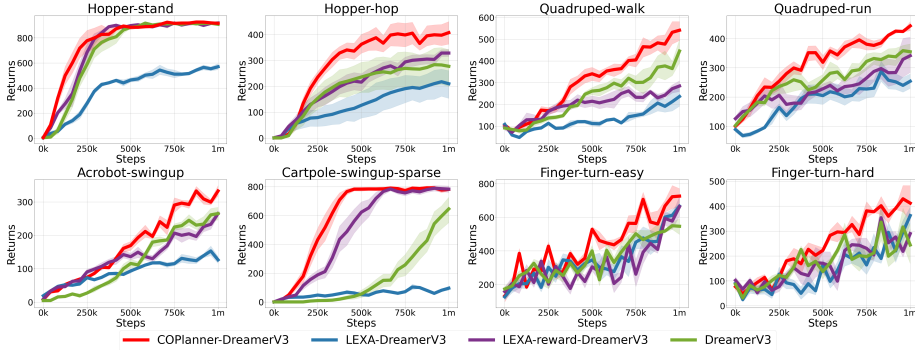


Figure 5: Experiment results of COPlanner-DreamerV3 and other three baselines on pixel-input DMC. The results are averaged over 8 random seeds, and shaded regions correspond to the 95% confidence interval among seeds. During evaluation, for each seed of each method, we test for up to 1000 steps in the test environment and perform 10 evaluations to obtain an average value. The evaluation interval is every 1000 environment steps.

260 **Environment and hyperparameter settings:** We use 8 visual control tasks of DMC as our en-  
 261 vironment. In COPlanner-DreamerV3, we learn a latent one-step prediction dynamics model as  
 262 Plan2Explore [25], the ensemble size is 8. We set action candidate number  $K$  and planning horizon  
 263  $H_p$  equal to 4 in all tasks. For optimistic rate  $\alpha_o$  and conservative rate  $\alpha_c$ , we set them to be 1 and  
 264 0.5, respectively. All other hyperparameters remain consistent with the original DreamerV3 paper.

265 **COPlanner significantly improves the sample efficiency and performance of DreamerV3:** From  
 266 the experiment results in Figure 5, we observe that COPlanner-DreamerV3 improves the sample  
 267 efficiency and performance significantly over DreamerV3. The sample efficiency of COPlanner-  
 268 DreamerV3 is more than twice that of DreamerV3, and the performance is improved by 23.9%.  
 269 Besides, LEXA-DreamerV3 has a low sample efficiency and do not perform well. This demonstrates  
 270 the limitation of pure exploration when the goal is to solve specific tasks instead of learning a  
 271 dynamics model applicable to a variety of tasks. After adding real reward as extrinsic reward for  
 272 explorer learning, LEXA-reward-DreamerV3 delivers performance comparable to DreamerV3 in most  
 273 environments. It outperforms DreamerV3 in Cartpole-swingup-sparse and Hopper-stand. However,  
 274 its performance and sample efficiency are still worse than COPlanner-DreamerV3, further indicates  
 275 the effectiveness of COPlanner.

## 276 6 Conclusion and discussion

277 We investigate how to effectively address the inaccurate learned dynamics model problem in MBRL.  
 278 We propose COPlanner, a general framework that can be applied to any dyna-style MBRL method.  
 279 COPlanner utilizes Uncertainty-aware Policy-Guided MPC phase to predict the cumulative uncer-  
 280 tainty of future steps and symmetrically uses this uncertainty as a penalty or bonus to select actions  
 281 for conservative model rollouts or optimistic environment exploration. In this way, COPlanner can  
 282 avoid model uncertain areas before model rollouts to minimize the impact of model error, while also  
 283 exploring high-reward model-uncertain areas in the environment to expand the model and reduce  
 284 model error. Experiments on a range of continuous control tasks demonstrates the effectiveness of our  
 285 method. One drawback of COPlanner is that MPC can lead to additional computational time and we  
 286 provide a detailed computational time consumption in Appendix D.7. We can improve computational  
 287 efficiency by parallelizing planning, which we leave for future work.

## 288 References

289 [1] Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample-  
 290 efficient reinforcement learning with stochastic ensemble value expansion. *Advances in neural*

- 291 *information processing systems*, 31, 2018.
- 292 [2] Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement  
293 learning in a handful of trials using probabilistic dynamics models. *Advances in Neural*  
294 *Information Processing Systems*, 31, 2018.
- 295 [3] Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach  
296 to policy search. In *Proceedings of the 28th International Conference on machine learning*  
297 (*ICML-11*), pages 465–472, 2011.
- 298 [4] Benjamin Eysenbach, Alexander Khazatsky, Sergey Levine, and Russ R Salakhutdinov. Mis-  
299 matched no more: Joint model-policy optimization for model-based rl. *Advances in Neural*  
300 *Information Processing Systems*, 35:23230–23243, 2022.
- 301 [5] Carlos E Garcia, David M Prett, and Manfred Morari. Model predictive control: Theory and  
302 practice—a survey. *Automatica*, 25(3):335–348, 1989.
- 303 [6] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-  
304 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International*  
305 *conference on machine learning*, pages 1861–1870. PMLR, 2018.
- 306 [7] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control:  
307 Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- 308 [8] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and  
309 James Davidson. Learning latent dynamics for planning from pixels. In *International conference*  
310 *on machine learning*, pages 2555–2565. PMLR, 2019.
- 311 [9] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with  
312 discrete world models. *arXiv preprint arXiv:2010.02193*, 2020.
- 313 [10] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains  
314 through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 315 [11] Nicklas Hansen, Xiaolong Wang, and Hao Su. Temporal difference learning for model predictive  
316 control. *arXiv preprint arXiv:2203.04955*, 2022.
- 317 [12] Edward S Hu, Richard Chang, Oleh Rybkin, and Dinesh Jayaraman. Planning goals for  
318 exploration. *arXiv preprint arXiv:2303.13002*, 2023.
- 319 [13] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model:  
320 Model-based policy optimization. *Advances in neural information processing systems*, 32,  
321 2019.
- 322 [14] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel:  
323 Model-based offline reinforcement learning. *Advances in neural information processing systems*,  
324 33:21810–21823, 2020.
- 325 [15] Thanard Kurutach, Ignasi Clavera, Yan Duan, Aviv Tamar, and Pieter Abbeel. Model-ensemble  
326 trust-region policy optimization. In *International Conference on Learning Representations*.
- 327 [16] Hang Lai, Jian Shen, Weinan Zhang, and Yong Yu. Bidirectional model-based policy op-  
328 timization. In *International Conference on Machine Learning*, pages 5618–5627. PMLR,  
329 2020.
- 330 [17] Nathan Lambert, Brandon Amos, Omry Yadan, and Roberto Calandra. Objective mismatch in  
331 model-based reinforcement learning. *arXiv preprint arXiv:2002.04523*, 2020.
- 332 [18] Russell Mendonca, Oleh Rybkin, Kostas Daniilidis, Danijar Hafner, and Deepak Pathak. Dis-  
333 covering and achieving goals via world models. *Advances in Neural Information Processing*  
334 *Systems*, 34:24379–24391, 2021.

- 335 [19] Andrew S Morgan, Daljeet Nandha, Georgia Chalvatzaki, Carlo D’Eramo, Aaron M Dollar,  
336 and Jan Peters. Model predictive actor-critic: Accelerating robot skill acquisition with deep  
337 reinforcement learning. In *2021 IEEE International Conference on Robotics and Automation*  
338 (*ICRA*), pages 6672–6678. IEEE, 2021.
- 339 [20] Feiyang Pan, Jia He, Dandan Tu, and Qing He. Trust the model when it is confident: Masked  
340 model-based actor-critic. *Advances in neural information processing systems*, 33:10537–10546,  
341 2020.
- 342 [21] Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-supervised exploration via disagree-  
343 ment. In *International conference on machine learning*, pages 5062–5071. PMLR, 2019.
- 344 [22] S Joe Qin and Thomas A Badgwell. A survey of industrial model predictive control technology.  
345 *Control engineering practice*, 11(7):733–764, 2003.
- 346 [23] Anil V Rao. A survey of numerical methods for optimal control. *Advances in the Astronautical*  
347 *Sciences*, 135(1):497–528, 2009.
- 348 [24] Neale Ratzlaff, Qinxun Bai, Li Fuxin, and Wei Xu. Implicit generative modeling for efficient  
349 exploration. In *International Conference on Machine Learning*, pages 7985–7995. PMLR,  
350 2020.
- 351 [25] Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak  
352 Pathak. Planning to explore via self-supervised world models. In *International Conference on*  
353 *Machine Learning*, pages 8583–8592. PMLR, 2020.
- 354 [26] Younggyo Seo, Lili Chen, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. State en-  
355 tropy maximization with random encoders for efficient exploration. In *International Conference*  
356 *on Machine Learning*, pages 9443–9454. PMLR, 2021.
- 357 [27] Jian Shen, Han Zhao, Weinan Zhang, and Yong Yu. Model-based policy optimization with  
358 unsupervised model adaptation. *Advances in Neural Information Processing Systems*, 33:  
359 2823–2834, 2020.
- 360 [28] Pranav Shyam, Wojciech Jaśkowski, and Faustino Gomez. Model-based active exploration. In  
361 *International conference on machine learning*, pages 5779–5788. PMLR, 2019.
- 362 [29] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012*  
363 *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033, 2012.
- 364 [30] Tingwu Wang and Jimmy Ba. Exploring model-based planning with policy networks. *arXiv*  
365 *preprint arXiv:1906.08649*, 2019.
- 366 [31] Xiyao Wang, Wichayaporn Wongkamjan, Ruonan Jia, and Furong Huang. Live in the moment:  
367 Learning dynamics model adapted to evolving policy. In *International Conference on Machine*  
368 *Learning*. PMLR, 2023.
- 369 [32] Zhihai Wang, Jie Wang, Qi Zhou, Bin Li, and Houqiang Li. Sample-efficient reinforcement  
370 learning via conservative model-based actor-critic. In *Proceedings of the AAAI Conference on*  
371 *Artificial Intelligence*, volume 36, pages 8612–8620, 2022.
- 372 [33] Zifan Wu, Chao Yu, Chen Chen, Jianye Hao, and Hankz Hankui Zhuo. Plan to predict: Learning  
373 an uncertainty-foreseeing model for model-based reinforcement learning. *Advances in Neural*  
374 *Information Processing Systems*, 35:15849–15861, 2022.
- 375 [34] Yao Yao, Li Xiao, Zhicheng An, Wanpeng Zhang, and Dijun Luo. Sample efficient reinforce-  
376 ment learning via model-ensemble exploration and exploitation. In *2021 IEEE International*  
377 *Conference on Robotics and Automation (ICRA)*, pages 4202–4208. IEEE, 2021.



- 378 [35] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Y Zou, Sergey Levine, Chelsea  
379 Finn, and Tengyu Ma. Mopo: Model-based offline policy optimization. *Advances in Neural*  
380 *Information Processing Systems*, 33:14129–14142, 2020.
- 381 [36] Tianjun Zhang, Paria Rashidinejad, Jiantao Jiao, Yuandong Tian, Joseph E Gonzalez, and Stuart  
382 Russell. Made: Exploration via maximizing deviation from explored regions. *Advances in*  
383 *Neural Information Processing Systems*, 34:9663–9680, 2021.
- 384 [37] Ruijie Zheng, Xiyao Wang, Huazhe Xu, and Furong Huang. Is model ensemble necessary?  
385 model-based rl via a single model with lipschitz regularized value function. In *International*  
386 *Conference on Learning Representations*, 2023.

# Appendix

## 388 A Detailed figure of COPlanner

389 We present a more detailed figure to illustrate our COPlanner framework. During **environment**  
 390 **exploration**, we first choose an action using UP-MPC with multi-step uncertainty bonus, then interact  
 391 with the real environment to obtain real samples for dynamics model learning. In **dynamics model**  
 392 **rollouts**, at each rollout step, we select the actions using UP-MPC with multi-step uncertainty penalty  
 393 to avoid model uncertain regions and interact with the learned dynamics model to get model-generated  
 394 samples to update the policy.

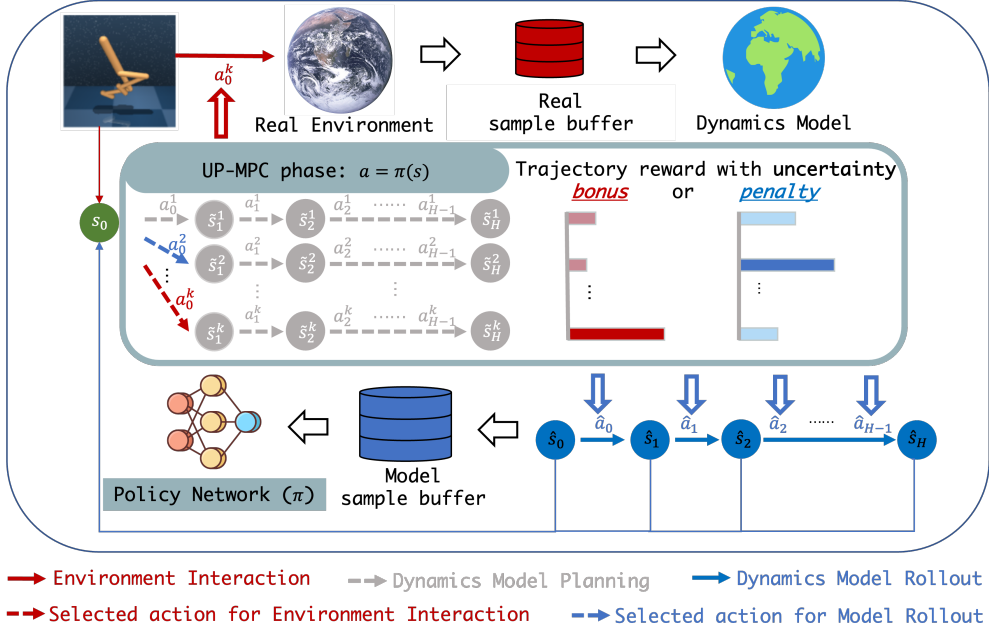


Figure 6: Figure illustration of COPlanner framework with more details.

## 395 B Implementation

396 COPlanner framework is versatile and applicable to any dyna-style MBRL algorithm. In this section,  
 397 we are going to introduce the implementation of two algorithms we used for experiment in Section 5:  
 398 COPlanner-MBPO for proprioceptive control and COPlanner-DreamerV3 for visual control.

### 399 B.1 COPlanner-MBPO

400 MBPO [13] trains an ensemble of probabilistic neural networks [2] as dynamics model. It utilises  
 401 negative log-likelihood loss to update each network in the ensemble:

$$\mathcal{L}(\theta) = \sum_{n=1}^N [\mu_{\theta}^b(s_n, a_n) - s_{n+1}]^T \Sigma_{\theta}^{b-1}(s_n, a_n) [\mu_{\theta}^b(s_n, a_n) - s_{n+1}] + \log \det \Sigma_{\theta}^b(s_n, a_n) \quad (6)$$

402 For the policy component, MBPO adopts soft actor-critic [6]. We combine COPlanner with MBPO,  
 403 the pseudocode is shown in Algorithm 3.

### 404 B.2 COPlanner-DreamerV3

405 DreamerV3 [10] is a dyna-style MBRL method that solves long-horizon tasks from visual inputs  
 406 purely by latent imagination. Its world model consists of an image encoder, a Recurrent State-Space

---

**Algorithm 3** COPlanner-MBPO

---

**Require:** interaction epochs  $I$ , rollout horizon  $H_r$ , planning horizon  $H_p$ , number of candidate actions  $K$ , conservative rate  $\alpha_c$ , optimistic rate  $\alpha_o$

- 1: Initialize policy  $\pi_\phi$ , dynamics model ensemble  $\hat{T}_\theta = \{\hat{T}_\theta^1, \dots, \hat{T}_\theta^i\}$ , real sample buffer  $\mathcal{D}_e$ , model sample buffer  $\mathcal{D}_m$
- 2: **for**  $I$  epochs **do**
- 3:   **for**  $t = 1$  to  $T$  **do**
- 4:     // *Optimistic environment exploration*
- 5:     Select action with optimistic rate  $a_t = \text{UP-MPC}(\pi_\phi, s_t, \hat{T}_\theta, K, H_p, \alpha_o)$
- 6:     Interact with the real environment with  $a_t$ , add real sample  $(s_t, a_t, r_t, s_{t+1})$  to real sample buffer  $\mathcal{D}_e$
- 7:     Train dynamics model  $\hat{T}_\theta$  via Equation 6
- 8:     **for**  $M$  model rollouts **do**
- 9:       Sample initial rollout states from real sample buffer  $\mathcal{D}_e$
- 10:       **for**  $h = 0$  to  $H_r - 1$  **do**
- 11:         // *Conservative model rollouts*
- 12:          $\hat{a}_h = \text{UP-MPC}(\pi_\phi, \hat{s}_h, \hat{T}_\theta, K, H_p, -\alpha_c)$  (Select action with conservative rate), rollout learned dynamics model and add to model sample buffer  $\mathcal{D}_m$
- 13:       **for**  $G$  gradient updates **do**
- 14:         Update current policy  $\pi_\phi$  using model-generated samples from model sample buffer  $\mathcal{D}_m$

---

407 Model (RSSM) [8] to learn the dynamics, and predictors for the image, reward, and discount factor.  
408 The world model components are:

$$\begin{aligned} \text{Recurrent model:} & \quad h_t = f_\phi(h_{t-1}, z_{t-1}, a_{t-1}) \\ \text{Representation model:} & \quad z_t \sim q_\phi(z_t | h_t, x_t) \\ \text{Transition predictor:} & \quad \hat{z}_t \sim p_\phi(\hat{z}_t | h_t) \\ \text{Image predictor:} & \quad \hat{x}_t \sim p_\phi(\hat{x}_t | h_t, z_t) \\ \text{Reward predictor:} & \quad \hat{r}_t \sim p_\phi(\hat{r}_t | h_t, z_t) \\ \text{Discount predictor:} & \quad \hat{\gamma}_t \sim p_\phi(\hat{\gamma}_t | h_t, z_t) \end{aligned}$$

409 where the recurrent model, the representation model, and the transition predictor are components of  
410 RSSM. The loss function for the world model learning is:

$$\begin{aligned} \mathcal{L}(\phi) = \mathbb{E}_{q_\phi(z_{1:T} | a_{1:T}, x_{1:T})} & \left[ \sum_{t=1}^T (-\ln p_\phi(x_t | h_t, z_t) - \ln p_\phi(r_t | h_t, z_t) - \ln p_\phi(\gamma_t | h_t, z_t)) \right. \\ & + \beta_1 \max(1, \text{KL}[sg(q_\phi(z_t | h_t, x_t)) || p_\phi(z_t | h_t)]) \\ & \left. + \beta_2 \max(1, \text{KL}[q_\phi(z_t | h_t, x_t) || sg(p_\phi(z_t | h_t))]) \right], \end{aligned} \quad (7)$$

411 where sg means stop gradient. Besides, DreamerV3 also use actor-critic framework as their policy.  
412 In particular, they leverage a stochastic actor that chooses actions and a deterministic critic. The  
413 actor and critic are trained cooperatively. The actor goal is to output actions leading to states that  
414 maximize the critic output, while the critic aims to accurately estimate the sum of future rewards that  
415 the actor can achieve from each imagined state (or model rollout state). For more training details  
416 about DreamerV3, please refer to their original paper [10].

417 To estimate model uncertainty in COPlanner-DreamerV3, we train an ensemble of one-step predictive  
418 models  $\hat{T}_\theta = \{\hat{T}_\theta^1, \dots, \hat{T}_\theta^i\}$ , each of these models takes a latent stochastic state  $z_t$  and action  $a_t$  as  
419 input and predicts the next latent deterministic recurrent states  $h_t$ . The ensemble is trained using  
420 MSE loss. During model rollouts, we use the world model to generate trajectories, and the one-step  
421 model ensemble to evaluate the uncertainty of sample at each rollout step. Here we provide the  
422 pseudocode of COPlanner-DreamerV3 in Algorithm 4.

---

**Algorithm 4** COPlanner-DreamerV3

---

**Require:** Rollout horizon  $H_r$ , planning horizon  $H_p$ , number of candidates actions  $K$ , conservative rate  $\alpha_c$ , optimistic rate  $\alpha_o$

- 1: Initialize real sample buffer  $\mathcal{D}_e$  with  $S$  random seed episodes.
- 2: Initialize policy  $\pi_\psi$ , critic  $v_\xi$ , one-step model ensemble  $\hat{T}_\theta = \{\hat{T}_\theta^1, \dots, \hat{T}_\theta^i\}$ , world model parameter  $\phi$
- 3: **while** not converged **do**
- 4:   **for** update step  $c = 1..C$  **do**
- 5:     Draw  $\mathcal{B}$  data sequences  $\{(a_t, x_t, r_t)\}_{t=k}^{k+L} \sim \mathcal{D}_e$
- 6:     Compute a latent stochastic states  $z_t \sim q_\phi(z_t|h_t, x_t)$
- 7:     Update world model parameter  $\phi$  via Equation 7
- 8:     *// Conservative model rollouts*
- 9:     Imagine trajectories  $\{(z_\tau, a_\tau)\}_{\tau=t}^{t+H_r}$  from each  $z_t$  with  $a_\tau = \text{UP-MPC}(\pi_\psi, z_\tau, \hat{T}_\theta, K, H_p, -\alpha_c)$ .
- 10:     Update  $v_\xi$  and  $\pi_\psi$  using imagined trajectories.
- 11:    **for** time step  $t = 1..T$  **do**
- 12:     Compute  $z_t \sim q_\phi(z_t|h_t, x_t)$
- 13:     *// Optimistic environment exploration*
- 14:     Select action  $a_t = \text{UP-MPC}(\pi_\psi, z_t, \hat{T}_\theta, K, H_p, \alpha_o)$ .
- 15:     Interact with the real environment and obtain  $(x_t, a_t, r_t, x_t + 1)$
- 16:     Add experience to  $\mathcal{D}_e \leftarrow \mathcal{D}_e \cup \{(x_t, a_t, r_t, x_t + 1)\}_{t=1}^T$

---

## 423 C Hyperparameters

424 In this section, we provide the specific parameters used in each task in our experiments.

Table 1: Hyperparameters of COPlanner-MBPO on proprioceptive control DMC.

Parameter	Value
Conservative rate $\alpha_c$	2
Optimistic rate $\alpha_o$	1
Action candidate number $K$	5
Planning horizon $H_p$	5
Real ratio	0.5 Reacher-xx 0.8 Finger-spin 0 Others
Rollout horizon $H_r$	[20, 150, 1, 1] Finger-spin [20, 150, 1, 4] Others

425

### 426 C.1 Proprioceptive control DMC and MuJoCo

427 We use COPlanner-MBPO in all proprioceptive control tasks. For the dynamics model ensemble, we  
428 adopted the same setup as MBPO [13] original paper, with an ensemble size of 7 and an elite number  
429 of 5, which means each time we select the best five out of seven neural networks for model rollouts.  
430 Each network in the ensemble is MLP with 4 hidden layers of size 200, using ReLU as the activation  
431 function. We train the dynamics model every 250 interaction steps with the environment. The actor  
432 and critic structures are both MLP with 4 hidden layers. In proprioceptive control DMC, the hidden  
433 layer size of actor and critic is 512, and updated 10 times each environment step, while in MuJoCo  
434 the hidden layer size is 256, and they are updated 20 times each environment step. The batch size for  
435 model training and policy training is both 256. The learning rate for model training is 1e-3, while the  
436 learning rate for policy training is 3e-4.

437 In MBPO, the authors use samples from both the real sample buffer and the model sample buffer to  
 438 train the policy, and the ratio of the two is referred to as the real ratio. In addition, MBPO has a unique  
 439 mechanism for the rollout horizon  $H_r$ , which linearly increases with the increase of environment  
 440 epochs, with each environment epoch including 1000 environment steps.  $[a, b, x, y]$  denotes a  
 441 thresholded linear function, *i.e.* at epoch  $e$ , rollout horizon is  $h = \min(\max(x + \frac{e-a}{b-a}(y-x), x), y)$ .  
 442 The settings for conservative rate  $\alpha_c$ , optimistic rate  $\alpha_o$ , action candidate number  $K$ , planning horizon  
 443  $H_p$  and the above two parameters in different environments are provided in Table 1 and 2.

Table 2: Hyperparameters of COPlanner-MBPO on MuJoCo.

Parameter	Value
Conservative rate $\alpha_c$	0.1 Hopper 2 Walker 0.5 Ant 1 Humanoid
Optimistic rate $\alpha_o$	0.05 Hopper 1 Walker, Humanoid 0.1 Ant
Action candidate number $K$	5
Planning horizon $H_p$	5
Real ratio	0.05
Rollout horizon $H_r$	[20, 100, 1, 4] Hopper 1 Walker [20, 150, 1, 15] Ant [20, 300, 1, 15] Humanoid

444

## 445 C.2 Visual control DMC

446 In Visual control DMC, we use the COPlanner-DreamerV3 method. We keep all parameters consistent  
 447 with the DreamerV3 original paper [10], except for our newly introduced conservative rate  $\alpha_c$ ,  
 448 optimistic rate  $\alpha_o$ , action candidate number  $K$ , and planning horizon  $H_p$ . In Table 3, we provide the  
 449 specific settings of conservative rate  $\alpha_c$ , optimistic rate  $\alpha_o$ , action candidate number  $K$ , and planning  
 450 horizon  $H_p$  for each task. It’s worth noting that, although using a conservative rate of 0.5 can perform  
 451 well, we find that for the two tasks in Quadruped, using a conservative rate of 2 yields the best sample  
 452 efficiency and performance. For other parameters, please refer to the original DreamerV3 paper. For  
 453 the one-step predictive model ensemble, we use a model ensemble with ensemble size of 8. Each  
 454 network in the ensemble is MLP with 5 hidden layers of size 1024.

Table 3: Hyperparameters of COPlanner-DreamerV3 on visual control DMC. We keep all other hyperparameters consistent with the DreamerV3 original paper.

Parameter	Value
Conservative rate $\alpha_c$	0.5
Optimistic rate $\alpha_o$	1
Action candidate number $K$	4
Planning horizon $H_p$	4

455

## 456 D More experiments

### 457 D.1 Comparison with more proprioceptive control MBRL methods

458 In this section, we compared our approach with more proprioceptive control MBRL methods on  
459 MuJoCo tasks. In addition to the three baseline methods from Section 5.1, MBPO [13], P2P-MPC  
460 [33], and MEEE [34], we introduced two more baselines: PDML [31], a method that dynamically  
461 adjusts the weights of each sample in the real sample buffer to enhance the prediction accuracy of the  
462 learned dynamics model for the current policy, thereby significantly improving the performance of  
463 MBPO. And MoPAC [19], a method that also uses policy-guided MPC to reduce model bias. Unlike  
464 our approach, MoPAC’s policy-guided MPC is solely used for multi-step prediction during rollout  
465 based on total reward to select actions. It does not incorporate a measure of model uncertainty, and  
466 therefore, cannot achieve the optimistic exploration and conservative rollouts of COPlanner. The  
467 experiment results are shown in Table 4.

468 As can be seen from Table 4, our method still holds a significant advantage, achieving the best  
469 performance in three tasks (Hopper, Walker2d, and Ant). In Humanoid task, it is only surpassed by  
470 PDML but is substantially better than the other methods. It’s worth mentioning that our approach is  
471 orthogonal to PDML, and they can be combined. We believe that by integrating COPlanner with  
472 PDML, the performance can be further enhanced.

Table 4: Comparison of different MBRL methods on proprioceptive control MuJoCo tasks. Performance is averaged over 8 random seeds.

	Hopper (70k)	Walker2d (150k)	Ant (150k)	Humanoid (150k)
Ours	<b>3325.6 ± 153.7</b>	<b>4402.8 ± 376.5</b>	<b>5142.3 ± 138.3</b>	4994.3 ± 449.4
PDML	3274.2 ± 224.1	4378.5 ± 248.9	4992.5 ± 365.1	<b>5396.7 ± 391.3</b>
MBPO	2844.6 ± 158.0	4221.1 ± 281.1	2311.1 ± 252.5	1706.0 ± 976.3
P2P-MPC	2316.8 ± 459.9	4151.7 ± 516.9	4681.7 ± 591.6	3706.1 ± 1360.4
MEEE	3076.4 ± 165.3	3873.8 ± 549.6	3932.8 ± 352.7	654.2 ± 94.7
MoPAC	3174.2 ± 233.8	2893.6 ± 472.6	4382.5 ± 301.7	1084.6 ± 573.2

### 473 D.2 Comparison with more visual control MBRL methods

474 In this section, we conducted comparisons with more MBRL methods that use latent dynamics  
475 models for visual control on 8 tasks from visual DMC. In addition to DreamerV3 [10] and LEVA  
476 [18], we introduced two more baselines. The first is TDMPC [11]. TDMPC learns a task-oriented  
477 latent dynamics model and uses this model for planning. During the planning process, TDMPC also  
478 learns a policy to sample a small number of actions, thereby accelerating MPC. The second is PlaNet  
479 [8]. PlaNet uses the RSSM latent model, which is the same as the Dreamer series [7, 9, 10], and  
480 directly uses this model to perform MPC in the latent space to select actions. The experiment results  
481 are shown in Table 5. From the results, it is evident that our method has a significant advantage over  
482 all the baselines.

### 483 D.3 Experiments combined with DreamerV2

484 We also combine COPlanner with DreamerV2 [9] for experimentation, with the results shown in  
485 Figure 7. After integrating with DreamerV2, our method also achieves a significant improvement in  
486 both sample efficiency and performance.

### 487 D.4 Model error and rollout uncertainty analysis

488 In this section, we will investigate the impact of COPlanner on model learning and model rollouts. We  
489 provide the curves of how model prediction error and rollout uncertainty change as the environment  
490 step increases in Figure 8. We conduct experiments on two proprioceptive control DMC tasks



Table 5: Performance comparison of different MBRL methods on visual DMC tasks at 1 million environment steps.

	Hopper-stand	Hopper-hop	Quadruped-walk	Quadruped-run
Ours	<b>916.2 ± 19.0</b>	<b>406.6 ± 105.6</b>	<b>541.8 ± 113.6</b>	<b>443.4 ± 39.3</b>
DreamerV3	908.3 ± 21.7	277.9 ± 163.6	445.3 ± 129.8	354.9 ± 69.0
LEXA	569.6 ± 56.2	209.1 ± 126.3	235.3 ± 117.7	254.5 ± 61.0
TDMPC	821.6 ± 70.8	189.2 ± 19.7	427.8 ± 50.2	393.8 ± 40.9
PlaNet (5m)	5.96	0.37	238.90	280.45

	Acrobot-swingup	Cartpole-swingup-sparse	Finger-turn-easy	Finger-turn-hard
Ours	<b>332.8 ± 37.0</b>	<b>781.8 ± 24.5</b>	<b>724.9 ± 126.3</b>	<b>414.2 ± 166.6</b>
DreamerV3	264.8 ± 44.8	647.0 ± 193.1	545.8 ± 108.8	243.9 ± 180.9
LEXA	126.5 ± 25.0	95.3 ± 35.5	666.6 ± 36.6	362.5 ± 79.5
TDMPC	227.5 ± 16.9	668.3 ± 49.1	703.8 ± 65.2	402.7 ± 112.6
PlaNet (5m)	3.21	0.64	451.22	312.55

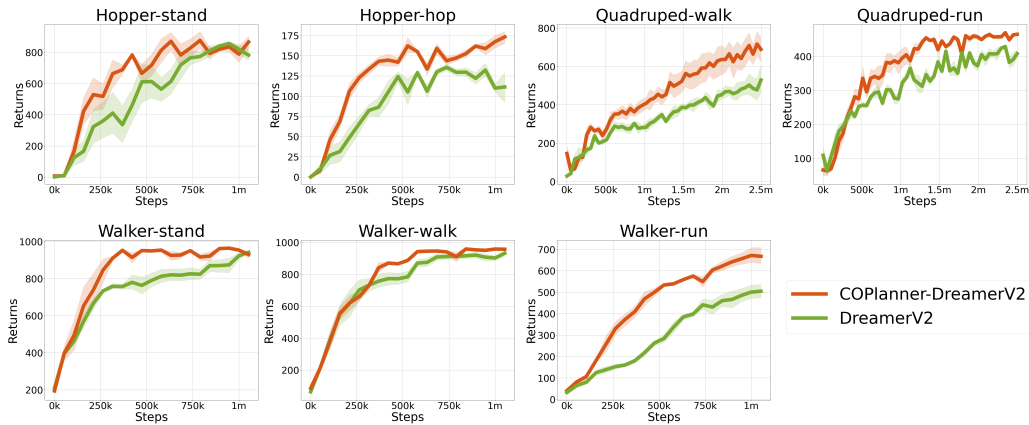


Figure 7: Experiment results of COPlanner-DreamerV2 on 7 visual control DMC tasks. The results are averaged over 4 random seeds, and shaded regions correspond to the 95% confidence interval among seeds. During evaluation, for each seed of each method, we test for up to 1000 steps in the test environment and perform 10 evaluations to obtain an average value. The evaluation interval is every 1000 environment steps.

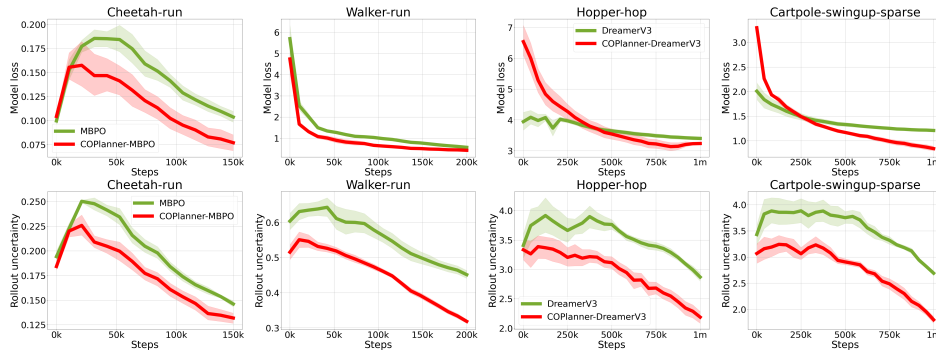


Figure 8: Model learning loss and rollout uncertainty curves for COPlanner and two other model-based RL baselines. The left four are proprioceptive control DMC tasks, and the right four are visual control DMC tasks.

491 (Cheetah-run and Walker-run) and two visual control DMC tasks (Hopper-hop and Cartpole-swingup-  
 492 sparse).

493 (1) In proprioceptive control DMC, we use the MSE loss between model prediction and ground truth  
 494 next state to evaluate model prediction error during training, while in visual control DMC, we use  
 495 the KL divergence between the latent dynamics prediction and the next stochastic representation to  
 496 compute latent model prediction error. We observe that after integrating COPlanner in proprioceptive  
 497 control tasks, the model prediction error is significantly reduced. In more complex visual control  
 498 tasks, due to obtaining more diverse samples through exploration in the early stages of training, the  
 499 model prediction error of COPlanner is higher than the baseline (DreamerV3). However, as training  
 500 progresses, the model prediction error rapidly decreases, becoming significantly lower than the  
 501 model prediction error of DreamerV3. This allows the model to fully learn from the diverse samples,  
 502 leading to an improvement in policy performance. (2) For the evaluation of rollouts uncertainty, we  
 503 calculate the model disagreement for each sample in the model-generated replay buffer used for policy  
 504 training using dynamics model ensemble. We find that COPlanner significantly reduces rollout  
 505 uncertainty due to conservative rollouts, suggesting that the impact of model errors on policy learning  
 506 is minimized. This experiment further demonstrates that the success of COPlanner is attributed to  
 507 both optimistic exploration and conservative rollouts.

### 508 D.5 Hyperparameter study

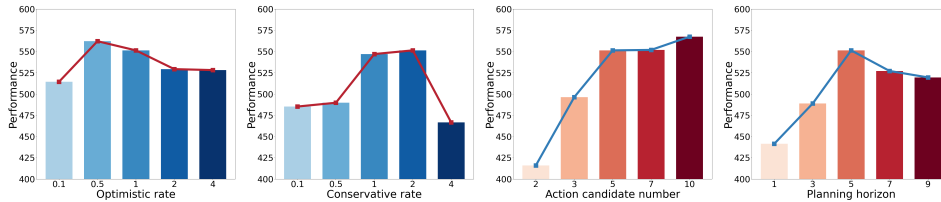


Figure 9: Ablation studies of COPlanner’s different hyperparameters. Experiments are conducted using COPlanner-MBPO on Walker-run tasks of proprioceptive control DMC. The results are averaged over 4 random seeds. From left to right, the results are for different parameters of optimistic rate  $\alpha_o$ , conservative rate  $\alpha_c$ , action candidate number  $K$ , and planning horizon  $H_p$ .

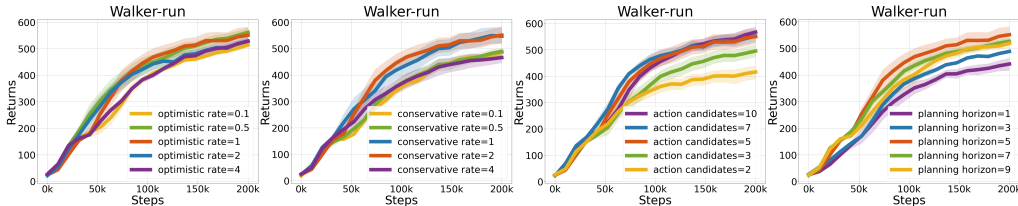


Figure 10: Performance curves of COPlanner’s hyperparameter study. Experiments are conducted using COPlanner-MBPO on Walker-run tasks of proprioceptive control DMC. The results are averaged over 4 random seeds.

509 In this section, we conducted hyperparameter studies to investigate the impact of different hyper-  
 510 parameters on COPlanner. We performed experiments on the Walker-run task of proprioceptive  
 511 control DMC using COPlanner-MBPO. The original hyperparameter settings are: optimistic rate  $\alpha_o$   
 512 is 1, conservative rate  $\alpha_c$  is 2, action candidate number  $K$  is 5, and planning horizon  $H_p$  is 5. When  
 513 conducting ablation experiments for each hyperparameter, other parameters remain unchanged. The  
 514 results are shown together in Figure 9 and more detailed curves are given in Figure 10.

515 **Optimistic rate  $\alpha_o$ :** we observe that the best  $\alpha_o$  lies between 0.5 to 1. When the  $\alpha_o$  is too large,  
 516 COPlanner tends to excessively explore high uncertainty areas while neglecting rewards, leading  
 517 to a decrease in sample efficiency and performance. On the other hand, when the  $\alpha_o$  is too small,  
 518 COPlanner fails to achieve the desired exploration effect.

519 **Conservative rate  $\alpha_c$ :** the optimal range for the  $\alpha_c$  is between 1 and 2. A too large  $\alpha_c$  may lead to  
 520 overly conservative selection of low-reward actions, while a too small  $\alpha_c$  would be unable to make  
 521 model rollouts avoid model uncertain areas.

522 **Action candidate number  $K$ :** we find that  $K$  has a significant impact on sample efficiency and per-  
 523 formance. When  $K$  is set to 2, the improvement of COPlanner over MBPO in terms of performance  
 524 and sample efficiency is relatively limited. This is reasonable because if there are only a few action  
 525 candidates, our selection space is very limited, and even with the use of uncertainty bonus and penalty  
 526 to select actions, there may not be much difference. When  $K$  increases to more than 5, the effect  
 527 of COPlanner becomes very stable, and more candidates do not bring noticeable improvements in  
 528 performance and sample efficiency.

529 **Planning horizon  $H_p$ :** when  $H_p$  is 1, we find that COPlanner’s improvement on performance and  
 530 sample efficiency is relatively limited. This also confirms what we mentioned in Section 1: only  
 531 considering the current step while ignoring the long-term uncertainty impact cannot completely  
 532 avoid model errors, as samples with low current model uncertainty might still lead to future rollout  
 533 trajectories falling into model uncertain regions. As the planning horizon gradually increases,  
 534 performance and sample efficiency also rise. When the planning horizon is too long ( $H_p$  equals to 7  
 535 or 9), it is possible that due to the bottleneck of the model planning capability, most action candidates’  
 536 corresponding trajectories fall into model uncertain areas, leading to a slight decline in performance  
 537 and sample efficiency.

### 538 D.6 Ablation study of model uncertainty estimation methods

539 We conduct an ablation study on the Hopper-hop task in visual control DMC to evaluate different  
 540 uncertainty estimation methods. We adopt two methods, RE3 [26] and MADE [36], which are used  
 541 to estimate intrinsic rewards in pixel input, to replace the disagreement in calculating  $u(s, a)$  in  
 542 Equation 4 and 5. The results are shown in Figure 11. We find that the performance achieved using  
 543 these two methods is similar to that of disagreement. This demonstrate that using disagreement to  
 544 calculate uncertainty is not the primary reason for the observed performance improvement.

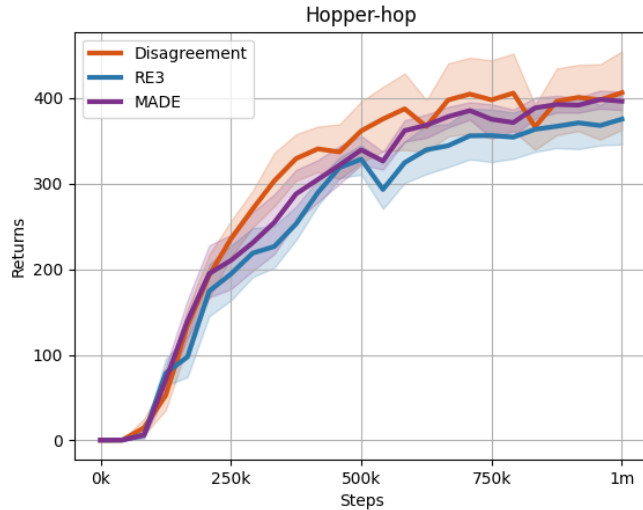


Figure 11: Ablation study of different uncertainty estimation methods.

### 545 D.7 Computational time consumption of COPlanner

546 We provide a comparison of the computational time consumption between the baseline methods and  
 547 COPlanner across different domains in Table 6. All timings are reported using a single NVIDIA  
 548 2080ti GPU.

Table 6: Average time consumption (h).

	MuJoCo	Proprioceptive DMC	Visual DMC
COPlanner-MBPO	41.2	11.3	N/A
MBPO	33.78	10.6	N/A
COPlanner-DreamerV3	N/A	N/A	17.9
DreamerV3	N/A	N/A	13.1

## 549 D.8 Ablation study

550 In this section, we aim to investigate the impact of different components within COPlanner on  
 551 the sample efficiency and performance. We conduct experiments on two proprioceptive control  
 552 DMC tasks (Walker-stand and Walker-run) using MBPO as baseline and two visual control DMC  
 553 tasks (Hopper-hop and Cartpole-swingup-sparse) with DreamerV3 as baseline. The results are  
 554 demonstrated in Figure 12. Due to page limitations, we provide the ablation study on various  
 555 hyperparameters of COPlanner in Appendix D.5 and the ablation study of uncertainty estimation  
 556 methods in Appendix D.6.

557 **From this ablation study, we can see that effectively combining optimistic exploration and**  
 558 **conservative rollouts is necessary to achieve the best results.** We find that when only using  
 559 optimistic exploration (COPlanner w. Explore only), the sample efficiency and performance in all  
 560 tasks are significantly improved, which highlights the importance of expanding the model. When  
 561 only using conservative rollouts (COPlanner w. Rollout only), there is some improvement in sample  
 562 efficiency and performance but to a lesser extent. In more complex visual control tasks, only using  
 563 conservative rollouts may lead to over-conservatism, resulting in an inability to learn an effective  
 564 policy in sparse reward environments (as observed with a broken seed in Cartpole-swingup-sparse) or  
 565 a decrease in sample efficiency during the early stages of learning (Hopper-hop). This is reasonable  
 566 because conservative rollouts may avoid high uncertainty and high reward areas to ensure the stability  
 567 of policy updates. Moreover, without efficiently expanding the model, it is challenging to find better  
 568 solutions using only conservative rollouts in complex visual control tasks. Experimental results show  
 569 that both optimistic exploration and conservative rollouts are crucial, and using either one individually  
 570 can lead to an improvement in performance. When combining the two (as in COPlanner), we can  
 571 achieve the best results, further demonstrating the effectiveness of our method.

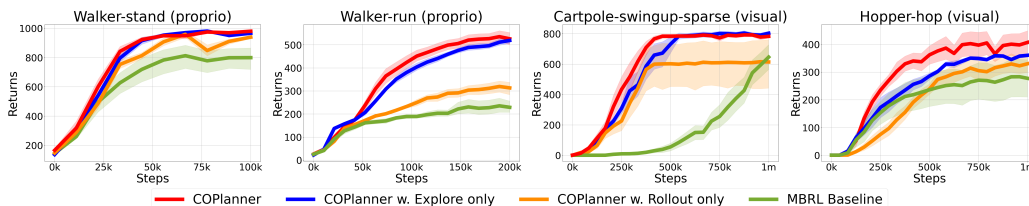


Figure 12: Ablation studies of optimistic exploration and conservative rollouts on different tasks using different mbrl baselines. In the first two proprioceptive control tasks we use MBPO as baseline. For the last two visual control tasks we employ DreamerV3. The results are averaged over 8 random seeds. We can observe that the best results are achieved when combining optimistic exploration and conservative rollouts. The benefit is more pronounced in more-challenging visual tasks.