

441 **A Additional Related Work**

442 **Bayesian IRL** Ramachandran and Amir [57] first proposed a Bayesian formulation of IRL to solve
 443 the reward ambiguity problem. A MAP inference approach was proposed in [58] and a variational
 444 inference approach was proposed in [59]. Their formulations considers non-entropy-regularized
 445 policies and the dynamics model is fixed during reward inference. In contrast, simultaneous estimation
 446 of reward and dynamics can potentially infer the demonstrator’s biased beliefs about the environment,
 447 which is desirable for psychology and human-robot interaction studies [15, 19, 17]. Despite the
 448 attractiveness, simultaneous estimation is challenging because of the need to invert the agent’s
 449 planning process, especially in continuous domains. Reddy et al. [17] avoids this by representing
 450 agent discrete action policies using neural network-parameterized Q functions and regularizing
 451 the Bellman error to be small over the entire state-action space. This method however cannot be
 452 straightforwardly adapted to the continuous action case. Kwon et al. [60] avoids this by first training
 453 a task-conditioned policy on a distribution of environments with known parameters using meta
 454 reinforcement learning and then use the meta-trained policy to guide inference. This precludes the
 455 method from being used in general settings with unknown task distributions. To our knowledge, our
 456 proposed algorithms are the first to address simultaneous estimation in general environments.

457 **Decision-aware model learning** Decision-aware model learning aims to solve the objective mismatch
 458 problem in model-based RL [61]. Many proposed methods in this class use value-targeted regression
 459 similar to our model loss in (16) [62, 63]. Our analysis and that of Vemula et al. [52] suggest that
 460 value-targeted model objectives may be related to robust objectives. Furthermore, since the set
 461 of value-equivalent models only shrink for increasingly larger set of policy and values [62], using
 462 value-aware model objective alone may not be optimal and additional prediction-based regularizations
 463 may be needed.

464 **B Missing Proofs**

465 **B.1 Proofs For Section 3.1**

466 **Derivation of BTOM Gradients (section 3.1).** Recall the definition of the optimal entropy-
 467 regularized policy and value functions:

$$\begin{aligned}\hat{\pi}(a|s; \theta) &= \frac{\exp(Q_\theta(s, a))}{\sum_{\tilde{a}} \exp(Q_\theta(s, \tilde{a}))} \\ Q_\theta(s, a) &= R_{\theta_1}(s, a) + \gamma \mathbb{E}_{\hat{P}_{\theta_2}(\cdot|s, a)}[V_\theta(s')] \\ V_\theta(s) &= \log \sum_{\tilde{a}} \exp(Q_\theta(s, \tilde{a}))\end{aligned}\tag{20}$$

468 The gradient of the policy log likelihood in terms of the Q function gradient is obtained as follow:

$$\begin{aligned}\nabla_\theta \log \hat{\pi}(a|s; \theta) &= \nabla_\theta Q_\theta(s, a) - \nabla_\theta V_\theta(s) \\ &= \nabla_\theta Q_\theta(s, a) - \frac{1}{Z_\theta} \nabla_\theta \sum_{\tilde{a}} \exp(Q_\theta(s, \tilde{a})) \\ &= \nabla_\theta Q_\theta(s, a) - \frac{1}{Z_\theta} \sum_{\tilde{a}} \exp(Q_\theta(s, \tilde{a})) \nabla_\theta Q_\theta(s, \tilde{a}) \\ &= \nabla_\theta Q_\theta(s, a) - \mathbb{E}_{\tilde{a} \sim \hat{\pi}(\cdot|s; \theta)}[\nabla_\theta Q_\theta(s, \tilde{a})]\end{aligned}\tag{21}$$

469 where $Z_\theta = \sum_{a'} \exp(Q_\theta(s, a'))$ is the normalizer.

470 Recall $\rho_{\hat{P}}^{\hat{\pi}}(\tilde{s}, \tilde{a}|s, a)$ is the discounted state-action occupancy measure starting from pair (s, a) . We
 471 define for any function $f(s, a)$:

$$\mathbb{E}_{\rho_{\hat{P}}^{\hat{\pi}}(\tilde{s}, \tilde{a}|s, a)}[f(s, a)] = \mathbb{E}_{\tau \sim (\hat{P}, \hat{\pi})} \left[\sum_{t=0}^{\infty} \gamma^t f(s, a) \middle| s_0 = s, a_0 = a \right] \quad (22)$$

472 We now derive Q function gradients with respect to the reward parameters θ_1 and dynamics parameters
 473 θ_2 , respectively.

$$\begin{aligned} \nabla_{\theta_1} Q_{\theta}(s, a) &= \nabla_{\theta_1} R_{\theta_1}(s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}_{\theta_2}(\cdot|s, a)}[\nabla_{\theta_1} V_{\theta}(s')] \\ &= \nabla_{\theta_1} R_{\theta_1}(s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}_{\theta_2}(\cdot|s, a), a' \sim \hat{\pi}(\cdot|s'; \theta)}[\nabla_{\theta_1} Q_{\theta}(s', a')] \\ &= \nabla_{\theta_1} R_{\theta_1}(s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}_{\theta_2}(\cdot|s, a), a' \sim \hat{\pi}(\cdot|s'; \theta)} \left[\right. \\ &\quad \left. \nabla_{\theta_1} R_{\theta_1}(s', a') + \gamma \mathbb{E}_{s'' \sim \hat{P}_{\theta_2}(\cdot|s', a'), a'' \sim \hat{\pi}(\cdot|s''; \theta)}[\nabla_{\theta_1} Q_{\theta}(s'', a'')] \right] \\ &= \nabla_{\theta_1} R_{\theta_1}(s, a) + \mathbb{E}_{\tau \sim (\hat{P}, \hat{\pi})} \left[\sum_{h=1}^{\infty} \gamma^h \nabla_{\theta_1} R_{\theta_1}(s_h, a_h) \middle| s_0 = s, a_0 = a \right] \\ &= \mathbb{E}_{\rho_{\hat{P}}^{\hat{\pi}}(\tilde{s}, \tilde{a}|s, a)}[\nabla_{\theta_1} R_{\theta_1}(\tilde{s}, \tilde{a})] \end{aligned} \quad (23)$$

474 In line two we used the result that $\nabla_{\phi} V_{\phi}(s)$ for both $\phi = \theta_1$ and $\phi = \theta_2$ corresponds to the second
 475 term in (21).

$$\begin{aligned} \nabla_{\theta_2} Q_{\theta}(s, a) &= \nabla_{\theta_2} R_{\theta_1}(s, a) + \nabla_{\theta_2} \gamma \mathbb{E}_{s' \sim \hat{P}_{\theta_2}(\cdot|s, a)}[V_{\theta}(s')] \\ &= \gamma \sum_{\tilde{s}} V_{\theta}(\tilde{s}) \nabla_{\theta_2} \hat{P}_{\theta_2}(\tilde{s}|s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}_{\theta_2}(\cdot|s, a), a' \sim \hat{\pi}(\cdot|s'; \theta)}[\nabla_{\theta_2} Q_{\theta}(s', a')] \\ &= \gamma \sum_{\tilde{s}} V_{\theta}(\tilde{s}) \nabla_{\theta_2} \hat{P}_{\theta_2}(\tilde{s}|s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}_{\theta_2}(\cdot|s, a), a' \sim \hat{\pi}(\cdot|s'; \theta)} \left[\right. \\ &\quad \left. \gamma \sum_{\tilde{s}} V_{\theta}(\tilde{s}) \nabla_{\theta_2} \hat{P}_{\theta_2}(\tilde{s}|s', a') + \gamma \mathbb{E}_{s'' \sim \hat{P}_{\theta_2}(\cdot|s', a'), a'' \sim \hat{\pi}(\cdot|s''; \theta)}[\nabla_{\theta_2} Q_{\theta}(s'', a'')] \right] \\ &= \gamma \sum_{\tilde{s}} V_{\theta}(\tilde{s}) \nabla_{\theta_2} \hat{P}_{\theta_2}(\tilde{s}|s, a) + \mathbb{E}_{\tau \sim (\hat{P}, \hat{\pi})} \left[\sum_{h=1}^{\infty} \gamma^{h+1} \sum_{\tilde{s}} V_{\theta}(\tilde{s}) \nabla_{\theta_2} \hat{P}_{\theta_2}(\tilde{s}|s_h, a_h) \middle| s_0 = s, a_0 = a \right] \\ &= \mathbb{E}_{\rho_{\hat{P}}^{\hat{\pi}}(\tilde{s}, \tilde{a}|s, a)} \left[\gamma \sum_{s'} V_{\theta}(s') \nabla_{\theta_2} \hat{P}_{\theta_2}(s'|\tilde{s}, \tilde{a}) \right] \end{aligned} \quad (24)$$

476 We make a quick remark on the identifiability of simultaneous estimation.

477 **Remark B.1.** *Simultaneous reward-dynamics estimation of the form (5) without specific assumptions*
 478 *on the prior $P(\theta)$ is in general unidentifiable.*

479 *Proof.* Let $\mathbf{R} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\mathbf{P} \in \mathbb{R}_+^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$, $\sum_{s'} \mathbf{P}_{ss'}^a = 1$, $\mathbf{Q} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\mathbf{V} \in \mathbb{R}^{|\mathcal{S}|}$ be a set
 480 of Bellman-consistent reward, dynamics, and value functions in matrix form. Let $\mathbf{P}' \neq \mathbf{P}$ be an
 481 alternative dynamics model. We can always find an alternative reward $\mathbf{R}' = \mathbf{R} + \Delta \mathbf{R}$, where:

$$\begin{aligned} \Delta \mathbf{R} &= (\mathbf{Q} - \mathbf{Q}) - \gamma(\mathbf{P}'\mathbf{V} - \mathbf{P}\mathbf{V}) \\ &= -\gamma \Delta \mathbf{P}\mathbf{V} \end{aligned} \quad (25)$$

482 without changing the value functions and optimal entropy-regularized policy. \square

483 Remark B.1 implies that existing simultaneous estimation approaches which do not use explicit or
 484 implicit regularizations, such as the SERD algorithm by [18], cannot in general accurately estimate
 485 expert reward. Paired with theorem 3.1, it shows that these algorithms cannot in general achieve good
 486 performance.

487 B.2 Proofs For Section 3.2

488 Derivation of discounted likelihood (14).

$$\begin{aligned}
 & \mathbb{E}_{P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t \log \hat{\pi}(a_t | s_t; \theta) \right] \\
 &= \mathbb{E}_{P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (Q_{\theta}(s_t, a_t) - V_{\theta}(s_t)) \right] \\
 &= \mathbb{E}_{P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t \left(R_{\theta_1}(s_t, a_t) + \gamma \mathbb{E}_{s' \sim \hat{P}_{\theta_2}(\cdot | s_t, a_t)} [V_{\theta}(s')] \right) \right] - \mathbb{E}_{P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t V_{\theta}(s_t) \right] \\
 &= \mathbb{E}_{P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t R_{\theta_1}(s_t, a_t) \right] - \mathbb{E}_{\mu} [V_{\theta}(s_0)] \\
 &\quad + \mathbb{E}_{P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^{t+1} \mathbb{E}_{s' \sim \hat{P}_{\theta}(\cdot | s_t, a_t)} [V_{\theta}(s')] \right] - \mathbb{E}_{P(\tau)} \left[\sum_{t=1}^{\infty} \gamma^t V_{\theta}(s_t) \right] \\
 &= \mathbb{E}_{P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t R_{\theta_1}(s_t, a_t) \right] - \mathbb{E}_{\mu} [V_{\theta}(s_0)] \\
 &\quad + \mathbb{E}_{P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^{t+1} \mathbb{E}_{s' \sim \hat{P}_{\theta}(\cdot | s_t, a_t)} [V_{\theta}(s')] \right] - \mathbb{E}_{P(\tau)} \left[\sum_{t=0}^{\infty} \gamma^{t+1} \mathbb{E}_{s' \sim P(\cdot | s_t, a_t)} [V_{\theta}(s')] \right] \\
 &= \underbrace{\mathbb{E}_{\rho_{\hat{P}}} [R_{\theta_1}(s_t, a_t)] - \mathbb{E}_{\mu} [V_{\theta}(s_0)]}_{\ell(\theta)} + \underbrace{\gamma \mathbb{E}_{\rho_{\hat{P}}} [\mathbb{E}_{s' \sim \hat{P}_{\theta}(\cdot | s_t, a_t)} V_{\theta}(s')] - \mathbb{E}_{s'' \sim P(\cdot | s_t, a_t)} V_{\theta}(s'')]}_{\mathbf{T1}}
 \end{aligned} \tag{26}$$

489 The following lemma shows that $\mathbf{T1}$ is negligible if the estimated dynamics is accurate under the
 490 expert distribution, which is available from the offline dataset.

491 **Lemma B.2.** Let $\epsilon = \mathbb{E}_{(s,a) \sim P(\tau)} D_{KL}(P(\cdot | s, a) || \hat{P}(\cdot | s, a))$ and $R_{max} = \max_{s,a} |R_{\theta}(s, a)| +$
 492 $\log |\mathcal{A}|$, it holds that

$$|\mathbf{T1}| \leq \frac{\gamma R_{max}}{(1-\gamma)^2} \sqrt{2\epsilon} \tag{27}$$

Proof.

$$\begin{aligned}
 |\mathbf{T1}| &= \left| \sum_{t=0}^{\infty} \gamma^{t+1} \mathbb{E}_{(s_t, a_t) \sim P(\tau)} \left[\sum_{s'} V_{\theta}(s') \left(\hat{P}(s' | s_t, a_t) - P(s' | s_t, a_t) \right) \right] \right| \\
 &\stackrel{(1)}{\leq} \sum_{t=0}^{\infty} \gamma^{t+1} \mathbb{E}_{(s_t, a_t) \sim P(\tau)} \left[\sum_{s'} |V_{\theta}(s')| \left| \hat{P}(s' | s_t, a_t) - P(s' | s_t, a_t) \right| \right] \\
 &\stackrel{(2)}{\leq} \sum_{t=0}^{\infty} \gamma^{t+1} \|V_{\theta}(\cdot)\|_{\infty} \mathbb{E}_{(s_t, a_t) \sim P(\tau)} \left[\left\| \hat{P}(\cdot | s_t, a_t) - P(\cdot | s_t, a_t) \right\|_1 \right] \\
 &\stackrel{(3)}{\leq} \sum_{t=0}^{\infty} \gamma^{t+1} \|V_{\theta}(\cdot)\|_{\infty} \sqrt{2 \mathbb{E}_{(s_t, a_t) \sim P(\tau)} D_{KL}(P || \hat{P})} \\
 &= \frac{\gamma}{1-\gamma} \|V_{\theta}(\cdot)\|_{\infty} \sqrt{2\epsilon}
 \end{aligned}$$

493 where (1) follows from Jensen's inequality, (2) follows from Holder's inequality, and (3) follows
 494 from Pinsker's inequality.

495 Finally, given $\mathcal{H}(\pi(a|s)) = -\sum_a \pi(a|s) \log \pi(a|s) \leq -\sum_a \pi(a|s) \log \frac{1}{|\mathcal{A}|} = \log |\mathcal{A}|$, we have
 496 $\|V_\theta(\cdot)\|_\infty \leq \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t (\max_{s,a} |R_\theta(s,a)| + \log |\mathcal{A}|)] = \frac{R_{\max}}{1-\gamma}$.

497 □

498 B.3 Proofs For Section 3.4

499 We first restate a slight modification of the result from [52], which decomposes the real environment
 500 performance gap between the expert and the learner into their policy and model advantages in the
 501 estimated dynamics:

502 **Lemma B.3.** (Performance difference via advantage in model; Lemma 4.1 in [52]) Let d_P^π denote the
 503 marginal state-action distribution following policy π in environment P . The following relationship
 504 holds:

$$\mathbb{E}_{(s,a) \sim d_P^\pi} [\log \hat{\pi}_{\hat{P}}(a|s)] = \mathbb{E}_{s \sim d_P^\pi} [\mathbb{E}_{a \sim \pi} Q_{\hat{P}}^{\hat{\pi}}(s,a) - V_{\hat{P}}^{\hat{\pi}}(s)] \quad (28)$$

$$= \underbrace{(1-\gamma) \mathbb{E}_{s \sim \mu} [V_P^\pi(s) - V_{\hat{P}}^{\hat{\pi}}(s)]}_{\text{Performance difference in real environment}} \quad (29)$$

$$+ \underbrace{\gamma \mathbb{E}_{(s,a) \sim d_P^\pi} [\mathbb{E}_{s' \sim P} V_{\hat{P}}^{\hat{\pi}}(s') - \mathbb{E}_{s'' \sim \hat{P}} V_{\hat{P}}^{\hat{\pi}}(s'')]}_{\text{Model (dis)advantage under learner distribution}} \quad (30)$$

$$+ \underbrace{\gamma \mathbb{E}_{(s,a) \sim d_P^\pi} [\mathbb{E}_{s' \sim \hat{P}} V_{\hat{P}}^{\hat{\pi}}(s') - \mathbb{E}_{s'' \sim P} V_{\hat{P}}^{\hat{\pi}}(s'')]}_{\text{Model advantage under expert distribution}} \quad (31)$$

505 The performance bound in theorem 3.1 can be obtained from lemma B.3 as follow:

506 **Theorem B.4.** (Restate of theorem 3.1) Let $\epsilon_{\hat{\pi}} = -\mathbb{E}_{(s,a) \sim d_P^\pi} [\log \hat{\pi}_{\hat{P}}(a|s)]$ be the policy estimation
 507 error and $\epsilon_{\hat{P}} = \mathbb{E}_{(s,a) \sim d_P^\pi} D_{KL}[P(\cdot|s,a) \|\hat{P}(\cdot|s,a)]$ be the dynamics estimation error. Let $R_{max} =$
 508 $\max_{s,a} |R_\theta(s,a)| + \log |\mathcal{A}|$. Assuming bounded expert-learner marginal state-action density ratio
 509 $\left\| \frac{d_{\hat{P}}^{\hat{\pi}}(s,a)}{d_P^\pi(s,a)} \right\|_\infty \leq C$, we have the following (absolute) performance bound for the IRL agent:

$$|J_P(\hat{\pi}) - J_P(\pi)| \leq \frac{1}{1-\gamma} \epsilon_{\hat{\pi}} + \frac{\gamma(C+1)R_{max}}{(1-\gamma)^2} \sqrt{2\epsilon_{\hat{P}}} \quad (32)$$

Proof.

$$\begin{aligned} & |J_P(\hat{\pi}) - J_P(\pi)| \\ & \leq \frac{1}{1-\gamma} \epsilon_{\hat{\pi}} \\ & \quad + \frac{\gamma}{1-\gamma} \mathbb{E}_{(s,a) \sim d_P^\pi} \left[\left| \frac{d_{\hat{P}}^{\hat{\pi}}(s,a)}{d_P^\pi(s,a)} (\mathbb{E}_{s' \sim P} V_{\hat{P}}^{\hat{\pi}}(s') - \mathbb{E}_{s'' \sim \hat{P}} V_{\hat{P}}^{\hat{\pi}}(s'')) \right| \right] \\ & \quad + \frac{\gamma}{1-\gamma} \mathbb{E}_{(s,a) \sim d_P^\pi} [|\mathbb{E}_{s' \sim \hat{P}} V_{\hat{P}}^{\hat{\pi}}(s') - \mathbb{E}_{s'' \sim P} V_{\hat{P}}^{\hat{\pi}}(s'')|] \\ & \leq \frac{1}{1-\gamma} \epsilon_{\hat{\pi}} \quad (33) \\ & \quad + \frac{\gamma}{1-\gamma} \left\| \frac{d_{\hat{P}}^{\hat{\pi}}(\cdot, \cdot)}{d_P^\pi(\cdot, \cdot)} \right\|_\infty \|V_{\hat{P}}^{\hat{\pi}}(\cdot)\|_\infty \mathbb{E}_{(s,a) \sim d_P^\pi} \left[\left\| \hat{P}(\cdot|s,a) - P(\cdot|s,a) \right\|_1 \right] \\ & \quad + \frac{\gamma}{1-\gamma} \|V_{\hat{P}}^{\hat{\pi}}(\cdot)\|_\infty \mathbb{E}_{(s,a) \sim d_P^\pi} \left[\left\| \hat{P}(\cdot|s,a) - P(\cdot|s,a) \right\|_1 \right] \\ & = \frac{1}{1-\gamma} \epsilon_{\hat{\pi}} + \frac{\gamma(C+1)R_{max}}{(1-\gamma)^2} \sqrt{2\epsilon_{\hat{P}}} \end{aligned}$$

510 where the last line uses results from lemma B.2. □

511 C Further Algorithm Details and Pseudo Code

512 We estimate the dynamics gradient in (16) and (18) using the REINFORCE method with baseline:

$$\begin{aligned}\nabla_{\theta_2} EV_{\theta}(s, a) &= \sum_{s'} V_{\theta}(s') \nabla_{\theta_2} \hat{P}_{\theta_2}(s'|s, a) \\ &= \mathbb{E}_{s' \sim \hat{P}(\cdot|s, a)} \left[(V_{\theta}(s') - b(s, a)) \nabla_{\theta_2} \log \hat{P}_{\theta_2}(s'|s, a) \right]\end{aligned}$$

513 Following Rigter et al. [45], we set the baseline to $b(s, a) = Q_{\theta}(s, a) - R_{\theta_1}(s, a)$ to reduce gradient
514 variance and further normalize $V_{\theta}(s') - b(s, a)$ across the mini-batch to stabilize training. In the
515 continuous-control setting, the value function can be estimated as $V_{\theta}(s) = \mathbb{E}_{a \sim \hat{\pi}_{\theta}} [Q_{\theta}(s, a) -$
516 $\log \hat{\pi}(a|s; \theta)]$ with a single sample. We apply this gradient for a fixed number of steps, which is a
517 hyperparameter.

518 Pseudo code for the proposed algorithms are listed in Algorithm 1 and Algorithm 2.

Algorithm 1 Deep Bayesian Theory of Mind (BTOM)

Require: Dataset $\mathcal{D} = \{\tau\}$, dynamics model $\hat{P}_{\theta_2}(s'|s, a)$, reward model $R_{\theta_1}(s, a)$, hyperparameters λ_1, λ_2

- 1: **for** $k = 1 : K$ **do**
- 2: Run MBPO to update learner policy $\hat{\pi}(a|s; \theta)$ and value function $Q_{\theta}(s, a)$ in dynamics \hat{P}
- 3: Sample real trajectory τ_{real} starting from $(s, a) \sim \mathcal{D}$ and following \hat{P} and $\hat{\pi}$
- 4: Sample fake trajectory τ_{fake} starting from $s \sim \mathcal{D}$, $a_{\text{fake}} \sim \hat{\pi}(\cdot|s; \theta)$ and following \hat{P} and $\hat{\pi}$
- 5: Evaluate (15) and take a gradient step
- 6: Evaluate (16) and take a few gradient steps.
- 7: **end for**

Algorithm 2 Robust Theory of Mind (RTOM)

Require: Dataset $\mathcal{D} = \{\tau\}$, dynamics model $\hat{P}_{\theta_2}(s'|s, a)$, reward model $R_{\theta_1}(s, a)$, hyperparameters λ_1, λ_2

- 1: **for** $k = 1 : K$ **do**
- 2: Run MBPO to update learner policy $\hat{\pi}(a|s; \theta)$ and value function $Q_{\theta}(s, a)$ in dynamics \hat{P}
- 3: Sample fake trajectory τ_{fake} starting from $s \sim \mathcal{D}$ and following \hat{P} and $\hat{\pi}$
- 4: Evaluate (17) and take a gradient step
- 5: Evaluate (18) and take a few gradient steps
- 6: **end for**

519 D Implementation Details

520 Our implementation builds on top of the official RAMBO implementation¹ [45].

521 D.1 MuJoCo Benchmarks

522 For the MuJoCo benchmarks described in section 4.2, we follow standard practices in model-based
523 RL.

524 D.1.1 Dynamics Pre-training

525 We use an ensemble of $K = 7$ neural networks where each network outputs the mean and covariance
526 parameters of a Gaussian distribution over the difference between the next state and the current state
527 $\delta = s' - s$:

$$\hat{P}_{\theta_2}^{(k)}(\delta|s, a) = \mathcal{N}(\delta | \mu_{\theta_2}^{(k)}(s, a), \Sigma_{\theta_2}^{(k)}(s, a)) \quad (34)$$

¹<https://github.com/marc-rigter/rambo>

528 Each network is a 4-layer feedforward network with 200 hidden units and Sigmoid linear unit (SiLU)
 529 activation function. For the initial pre-training step, we maximize the likelihood of dataset transitions
 530 using a batch size of 256 and early stop when all models stop improving for more than 1 percent.
 531 We then select the 5 best models in terms of mean-squared-error on a 10 % holdout validation set.
 532 During model rollouts, we randomly pick one of the 5 best models (elites) to sample the next state.

Table 2: Shared hyperparameters across different environments

	Hyperparameter	BTOM	RTOM
SAC + MBPO	critic learning rate	3e-4	3e-4
	actor learning rate	3e-4	3e-4
	discount factor (γ)	0.99	0.99
	soft target update parameter (τ)	5e-3	5e-3
	target entropy	$-\dim(A)$	$-\dim(A)$
	minimum temperature (α)	0.1	0.001
	batch size	256	256
	real ratio	0.5	0.5
	model retain epochs	5	5
	training epochs	500	300
	steps per epoch	1000	1000
Dynamics	# model networks	7	7
	# elites	5	5
	adv. rollout batch size	1000	256
	adv. rollout steps	10	10
	adv. update steps	50	50
	adv. loss weighting (λ_1)	0.01	0.01
	supervised. loss weighting (λ_2)	1	1
	learning rate	1e-4	1e-4
	adv. update steps	50	50
	max reward	10	10
Reward	rollout batch size	1000	64
	rollout steps	40	100
	l2 penalty	1e-3	1e-3
	learning rate	1e-4	1e-4
	update steps	1	1

Table 3: Environment-specific hyperparameters

Environment	Hyperparameter	BTOM
HalfCheetah	model rollout batch size	50000
	model rollout steps	5
	model rollout frequency	250
Hopper	model rollout batch size	10000
	model rollout steps	40
	model rollout frequency	250
Walker2d	model rollout batch size	10000
	model rollout steps	40
	model rollout frequency	250

533 D.1.2 Policy Training

534 Our policy training process follows MBPO [37] which uses SAC with automatic temperature tuning
 535 [64]. Shared hyperparameters across different environments are listed in Table 2 and environment-
 536 specific hyperparameters are listed in Table 3. For the actor and critic, we use feedforward neural
 537 networks with 2 hidden layers of 256 units and ReLU activation. We train the actor and critic
 538 networks using a combination of real and simulated samples. We use a real ratio of 0.5, which is
 539 standard practice in model-based RL and IRL. We found that BTOM requires a higher minimum
 540 temperature to stabilize training, which is set to $\alpha = 0.1$.

541 We found that different MuJoCo environments require different model rollout hyperparameters, simi-
 542 lar to what’s reported in [41]. Specifically, Hopper and Walker2d only work with significantly larger
 543 rollout steps. We decrease their rollout batch size to reduce computational overhead. HalfCheetah on
 544 the other hand works better with smaller rollout steps and larger rollout batch size. In contrast to Lu
 545 et al. [41], we did not use different rollout hyperparameters for different datasets.

546 D.1.3 Reward and Dynamics Training

547 We use 10 random trajectories from the D4RL MuJoCo expert dataset after removing all expert
 548 trajectories that resulted in terminal states.

549 We use the same network architecture as the actor-critic to parameterize the reward function. We
 550 further clip the reward function to a maximum range of ± 10 and apply l2 regularization on all weights
 551 with a penalty of 0.001.

552 As described in the main text, we update the reward function by simulating sample trajectories and
 553 taking a single gradient step. For RTOM, we randomly sample expert trajectory segments of length
 554 “rollout steps” and use the first step as the start of our simulated sample paths.

555 We update the dynamics using on-policy rollouts branched from the dataset state-actions. We use
 556 the same batch size for reward and dynamics rollouts, which is 1000 for BTOM and 256 for RTOM.
 557 Because only the first step in BTOM’s real sample paths come from the dataset, it requires a larger
 558 batch size to iterate more data samples. We also train BTOM for more epochs than RTOM.

559 To compute the dynamics log likelihood in the REINFORCE gradient in (34), we treat the ensemble
 560 as a uniform mixture and compute the likelihood as:

$$\hat{P}_{\theta_2}(\delta|s, a) = \frac{1}{K} \sum_{k=1}^K \hat{P}_{\theta_2}^{(k)}(\delta|s, a) \quad (35)$$

561 We set the dynamics adversarial loss weighting to $\lambda_1 = 0.01$ for both BTOM and RTOM. We found
 562 this to work better than what’s in the official RAMBO implementation, which is $\lambda_1 = 0.0768$. Note
 563 that the RAMBO author reported $\lambda_1 = 3e-4$ in their paper but forget to average their REINFORCE
 564 loss over the mini-batch of size 256 in their implementation, which is instead treated as a sum by
 565 default by TensorFlow. We empirically found that small λ_1 leads to severe model exploitation.