

Model-based Reinforcement Learning (MBRL)

Networked Intelligence for Comprehensive Efficiency (NICE) Lab
College of Information Science and Electronic Engineering
Zhejiang University
<http://nice.rongpeng.info/>



Sep. 1, 2025

Content



1 Basic of model-based RL

- Model-based vs. Model-free
- What is a Model?

2 Evolution of Model-based RL

- Distributional shift in model-based RL
- History about Model-based Reinforcement Learning
- General Dyna-Style Algorithms
- Uncertainty in model-based RL
- Example: PETs

3 Strengths and Weaknesses of Model-based RL

4 Recent Technic to the Challenges of MBRL

- Error Accumulation
- OOD and Quantify Uncertainty

5 Conclusion of MBRL

Basic of model-based RL



Overview: Model-based vs. Model-free

RL systems can make decisions in two fundamental ways:

Diagram of model-free reinforcement learning

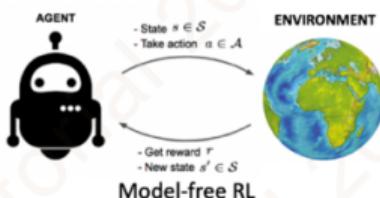
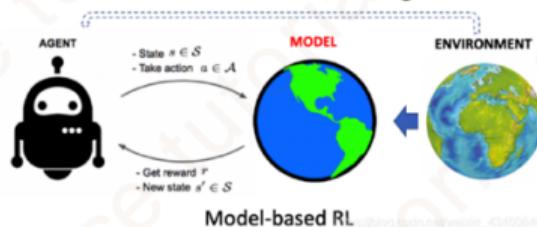


Diagram of model-based reinforcement learning



■ Model-free methods

- ▶ Disregard explicit modeling of the environment
- ▶ Policies are learned directly from interaction and experience.
- ▶ No requirement to understand the internal dynamics of the environment.
- ▶ Reference: BAIR Blog (2019).



Overview——Model-based vs. Model-free

RL systems can make decisions in two fundamental ways:

Diagram of model-free reinforcement learning

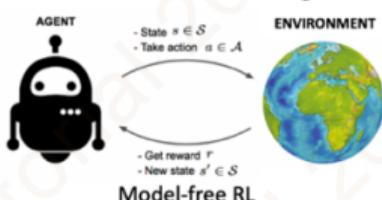
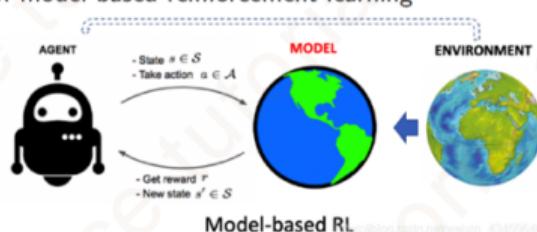


Diagram of model-based reinforcement learning



■ Model-based methods

- A **model** is something that mimics the behavior of the environment, or more generally, allows inferences to be made about how the environment will behave. [Sutton RL Book]
- A model provides an approximation of the environment's dynamics—an internal representation of states, actions, and transitions, rewards.
- Common terms for such models include:
 - ▶ Predictive model
 - ▶ World model
 - ▶ Environment (M) model
 - ▶ Learned dynamics function



What is a Model?

- Sequential decision-making problems are typically formulated as a **Markov Decision Process (MDP)**: $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, P, R\}$, where \mathcal{S} denotes the state space, \mathcal{A} the action space, P the state transition dynamics, and R the reward function.
- **MBRL** is generally composed of two key components: **Learning** and **Planning**.
 - **Learning**: Estimating the environment model $\mathcal{M} = (P, R)$ from collected experience.
 - **Planning**: Leveraging the learned model to simulate interactions and optimize decision-making policies.
- The model consists of:
 - A state transition probability function $f: \mathcal{S} \times \mathcal{A} \rightarrow p(\mathcal{S})$, which predicts the next state.
 - A reward function $R: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$, which assigns a scalar feedback signal.

$$p_{\theta}(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T) = p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

Policy model

$$\max_{\theta} E_{\tau \sim p_{\theta}(\tau)} [\sum_t r(\mathbf{s}_t, \mathbf{a}_t)]_{\text{policy}}$$

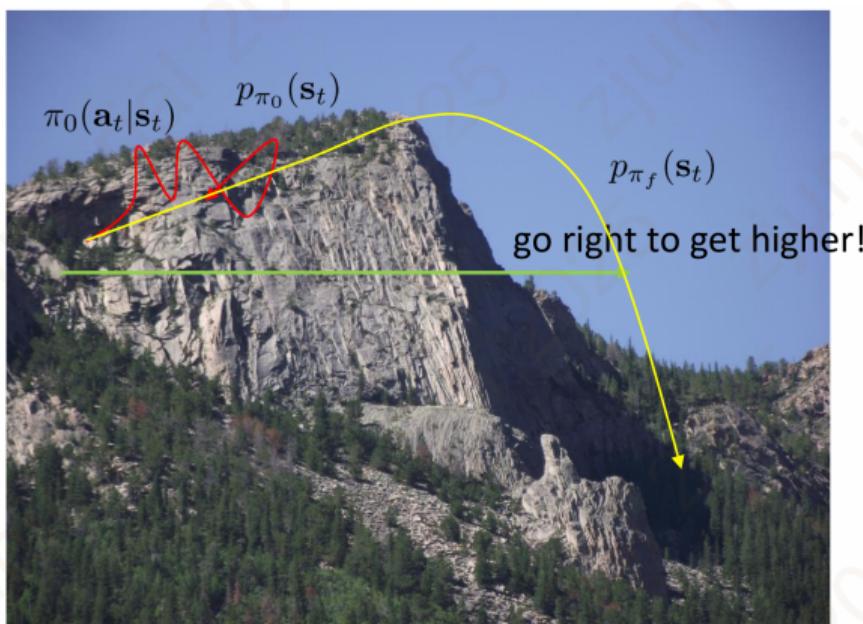
Reward

Figure: RL objective function and the role of the model in simulating environment dynamics.

Evolution of Model-based RL



Distribution Mismatch Problem



- 1 Run base policy $\pi_0(a_t|s_t)$ (e.g., random policy) to collect $\mathcal{D} = \{(s, a, s')_i\}$
- 2 Learn dynamics model $f(s, a)$ to minimize $\sum_i \|f(s_i, a_i) - s'_i\|^2$
- 3 Plan through $f(s, a)$ to choose actions

$$p_{\pi_f}(s_t) \neq p_{\pi_0}(s_t)$$

- Distribution mismatch problem becomes exacerbated as we use more expressive model classes.

Model-based Reinforcement Learning v1.0



Can we make $p_{\pi_0}(s_t) = p_{\pi_f}(s_t)$?

We need to collect data from $p_{\pi_f}(s_t)$.

Model-based reinforcement learning version 1.0:

- 1 Run base policy $\pi_0(a_t|s_t)$ (e.g., random policy) to collect $\mathcal{D} = \{(s, a, s')_i\}$.
- 2 Learn dynamics model $f(s, a)$ to minimize $\sum_i \|f(s_i, a_i) - s'_i\|^2$.
- 3 Plan through $f(s, a)$ to choose actions.
- 4 Execute those actions and add the resulting data $\{(s, a, s')_j\}$ to \mathcal{D} .

Can We Do Better?-Model Predictive Control(MPC)



Replanning helps mitigate model errors.

model-based reinforcement learning version 1.5:

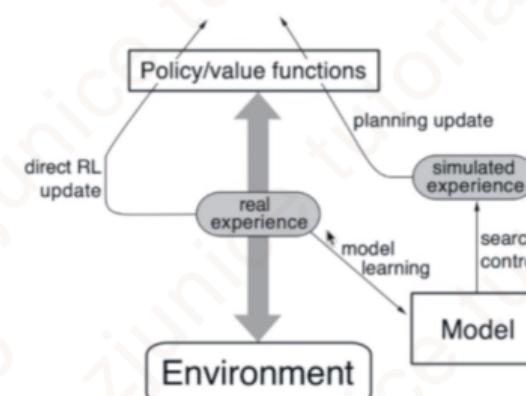
every N steps

1. run base policy $\pi_0(\mathbf{a}_t|\mathbf{s}_t)$ (e.g., random policy) to collect $\mathcal{D} = \{(\mathbf{s}, \mathbf{a}, \mathbf{s}')_i\}$
2. learn dynamics model $f(\mathbf{s}, \mathbf{a})$ to minimize $\sum_i \|f(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{s}'_i\|^2$
3. plan through $f(\mathbf{s}, \mathbf{a})$ to choose actions
4. execute the first planned action, observe resulting state \mathbf{s}' (MPC)
5. append $(\mathbf{s}, \mathbf{a}, \mathbf{s}')$ to dataset \mathcal{D}

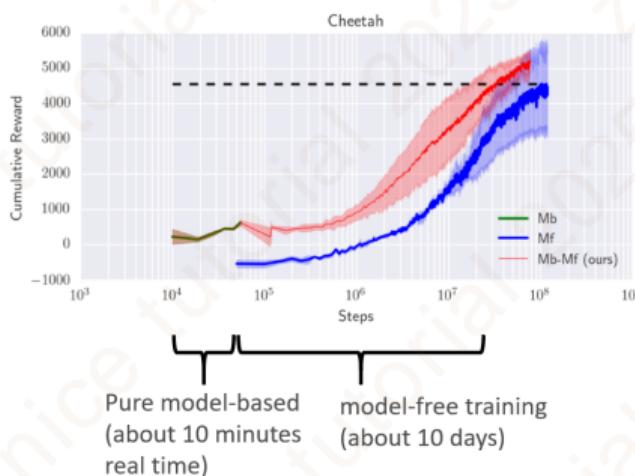


General Dyna-Style Algorithms

- **Model-based / Model-free Hybrid:** Early methods combined strengths of both paradigms to improve performance.
- **Roles of models in Dyna-style algorithms:**
 - 1 Learn an environment model; use simulated experience to update value or policy (*indirect RL*).
 - 2 Directly update value or policy from real experience (*direct RL*).
- Dyna integrates real and simulated experience, enabling both direct and planning updates, thereby improving **sample efficiency**.
- **Core Idea:** Replace part of real experience with model-generated transitions to reduce interaction cost and accelerate learning.



Example: Neural Network Dynamics for Model-based RL with Model-free Fine-tuning



Comparison of pure model-based, model-free, and hybrid (Mb-Mf) training.

Why the performance gap?

- UCB team evaluated a **hybrid model-based + model-free** approach against pure model-free RL.
- Method: *Probabilistic Neural Network + MPC + TRPO*.
- **Observation:** Pure model-based achieves rapid early gains (brown curve), but quickly plateaus.
- **Reason:** Limited real data and restricted model capacity cause **overfitting**, leading to stagnation in cumulative reward.
- The **hybrid Mb-Mf RL approach** = model-based for fast start + model-free for strong finish



How can Uncertainty Estimation Help?

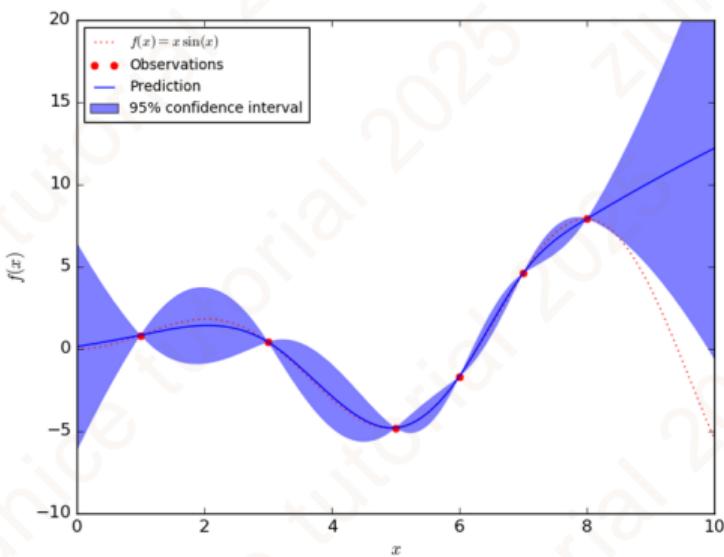


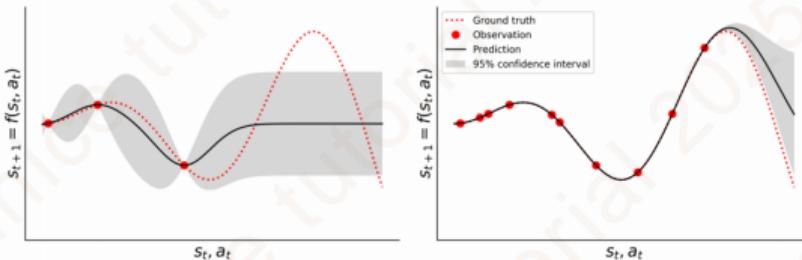
Illustration: predictions with 95% confidence interval.

- Need to explore to improve performance.
- Expected value \neq pessimistic value.
- Expected value \neq optimistic value.
- But expected value is often a good start for decision-making.

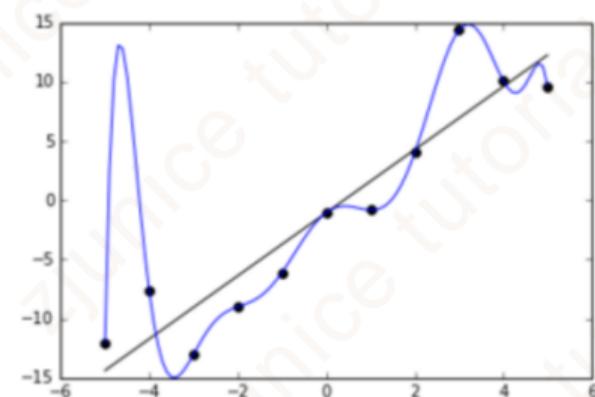


Classification of Uncertainty

- **Epistemic (Model) Uncertainty:** Caused by limited or insufficient training data, leading to an inaccurate approximation of the environment. It can be reduced by collecting more data or improving the model.
- **Aleatoric (Statistical) Uncertainty:** Inherent randomness in the data (e.g., noise or stochasticity). It cannot be reduced even with unlimited data.



Epistemic uncertainty: Caused by intrinsic noise; cannot be reduced by adding data.



Aleatoric uncertainty: Few observations (left) yield wide prediction intervals; More data (right) reduces uncertainty.



Estimating Model Uncertainty

Deterministic NN (point estimate):

$$s_{t+1} = f_\theta(s_t, a_t), \quad \mathcal{L}_{\text{det}}(\theta) = \sum_{n=1}^N \|s_{n+1} - f_\theta(s_n, a_n)\|^2 \quad (1)$$

Uncertainty: none (residuals reflect only data noise).

Probabilistic NN (stochastic dynamics):

$$p_\theta(s_{t+1} | s_t, a_t) = \mathcal{N}(\mu_\theta, \Sigma_\theta), \quad \mathcal{L}_{\text{prob}}(\theta) = - \sum_{n=1}^N \log p_\theta(s_{n+1} | s_n, a_n) \quad (2)$$

Uncertainty: captures **aleatoric** (data) uncertainty via the predicted variance.

Model	Aleatoric uncertainty	Epistemic uncertainty
<i>Baseline Models</i>		
Deterministic NN (D)	No	No
Probabilistic NN (P)	Yes	No
Deterministic ensemble NN (DE)	No	Yes
Gaussian process baseline (GP)	Homoscedastic	Yes
<i>Our Model</i>		
Probabilistic ensemble NN (PE)	Yes	Yes

Figure: Uncertainties captured by different models.



Estimating Model Uncertainty

- Ensemble NN (bootstrap / independent nets)

$$\hat{f}_\theta(s, a) = \frac{1}{B} \sum_{b=1}^B f_{\theta_b}(s, a),$$

$\text{Var}\left[\{f_{\theta_b}(s, a)\}_{b=1}^B\right] \Rightarrow \text{epistemic uncertainty}$

Idea: mean gives prediction; dispersion across members estimates model uncertainty.

- Takeaways

- Deterministic: no uncertainty (only residual noise).
- Probabilistic: models **aleatoric** via variance.
- Ensemble: estimates **epistemic** via member variance.

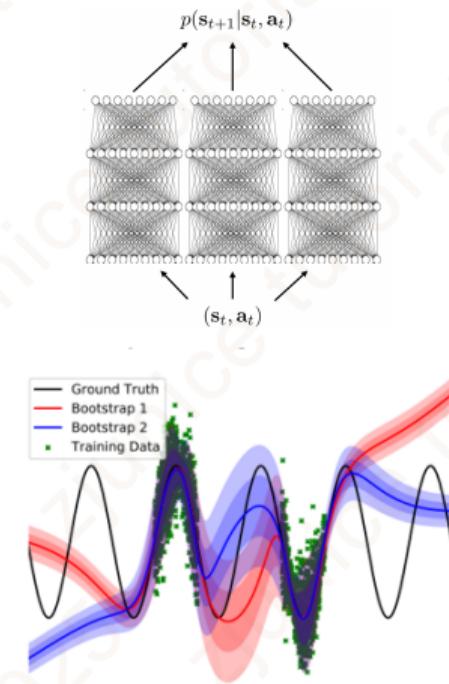


Figure: Predicted variance (aleatoric).

Example: Probabilistic Ensembles with Trajectory Sampling (PETS)

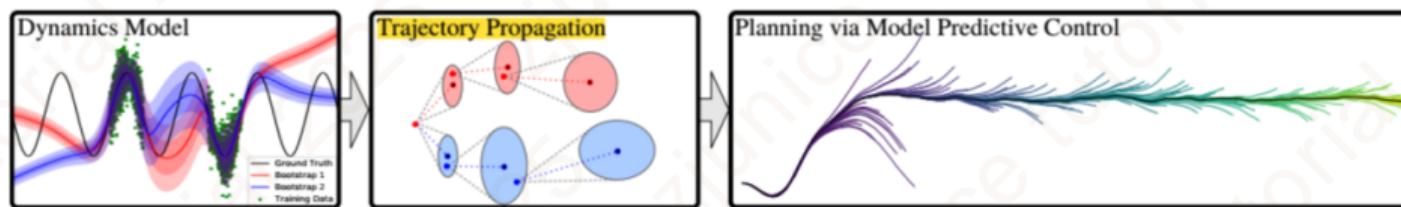


Figure: PETS framework: Dynamics Model → Trajectory Propagation → MPC Planning

- PETS combines the strengths of existing methods, using an **ensemble model** for prediction and a **probabilistic dynamics model** to output distributions. Effectively separates the two major types of uncertainty in MBRL: **Aleatoric** and **Epistemic**.
- Compared with model-free RL, PETS achieves significantly higher sample efficiency. On MuJoCo benchmarks, PETS reaches asymptotic performance comparable to SAC and PPO with far fewer samples.

Chua K, Calandra R, et al., "Deep reinforcement learning in a handful of trials using probabilistic dynamics models," *NeurIPS*, 2018.

R. Wen, J. Huang, et al., "Multi-Agent Probabilistic Ensembles With Trajectory Sampling for Connected Autonomous Vehicles," *IEEE T-VT*, vol. 73, no. 11, pp.

16076–16091, Nov. 2024.

Strengths and Weaknesses of Model-based RL

Strengths and Weaknesses of Model-based RL



■ Strengths

- Reduced reliance on real environment interactions:
 - ▶ Higher **sample efficiency**, crucial for practical deployment.
 - ▶ Faster training and adaptation compared to purely model-free methods.

■ Weaknesses

- **Error Accumulation:** Prediction errors compound during multi-step rollouts, often leading to inferior final performance.
- **Out-of-Distribution (OOD):** Models face difficulty when planning over unseen state-action pairs, causing unreliable predictions.
- **Uncertainty Estimation:** Accurately quantifying both **epistemic** (model) and **aleatoric** (data) uncertainty remains challenging, limiting robustness and safety.

Recent Technic to the Challenges of MBRL



Error Accumulation in Model-Based RL

Cause:

- Dynamics models inevitably have approximation errors.
- Small errors accumulate and amplify over long-horizon rollouts.
- Planning heavily relies on predicted trajectories \Rightarrow errors grow exponentially.

Solutions:

- **Short-horizon / Branched Rollouts** Start rollouts from real states and limit rollout length (e.g., *MBPO*).
- **Ensemble Models** Use multiple models and average predictions to reduce bias (e.g., *PETS*).
- **Uncertainty-based Truncation** Stop rollouts once prediction uncertainty becomes too high (e.g., *UA-MBRL*).
- **Model Regularization / Better Dynamics Models** Improve generalization and reduce error propagation.



Example: Model-Based Policy Optimization (MBPO)

- **Key Idea:** Start rollouts from real states with short horizons (branched rollout) to reduce model bias and error accumulation.

- **Main Components:**

- 1 **Predictive model:** Use a model ensemble to estimate transition dynamics.

$$p_{\theta}(s_{t+1} \mid s_t, a_t) = \mathcal{N}(\mu_{\theta}(s_t, a_t), \Sigma_{\theta}(s_t, a_t))$$

- 2 **Policy optimization:** Apply model-free policy gradient (e.g., SAC).
 - 3 **Data generation:** From real states, perform *k-step branched rollouts* to augment training data.
- **Trade-off:** Rollout length k balances sample efficiency and policy performance.



Algorithm 2 Model-Based Policy Optimization with Deep Reinforcement Learning

```

1: Initialize policy  $\pi_{\phi}$ , predictive model  $p_{\theta}$ , environment dataset  $\mathcal{D}_{\text{env}}$ , model dataset  $\mathcal{D}_{\text{model}}$ 
2: for  $N$  epochs do
3:   Train model  $p_{\theta}$  on  $\mathcal{D}_{\text{env}}$  via maximum likelihood
4:   for  $E$  steps do
5:     Take action in environment according to  $\pi_{\phi}$ ; add to  $\mathcal{D}_{\text{env}}$ 
6:     for  $M$  model rollouts do
7:       Sample  $s_t$  uniformly from  $\mathcal{D}_{\text{env}}$ 
8:       Perform  $k$ -step model rollout starting from  $s_t$  using policy  $\pi_{\phi}$ ; add to  $\mathcal{D}_{\text{model}}$ 
9:     for  $G$  gradient updates do
10:    Update policy parameters on model data:  $\phi \leftarrow \phi - \lambda_{\pi} \hat{\nabla}_{\phi} J_{\pi}(\phi, \mathcal{D}_{\text{model}})$ 

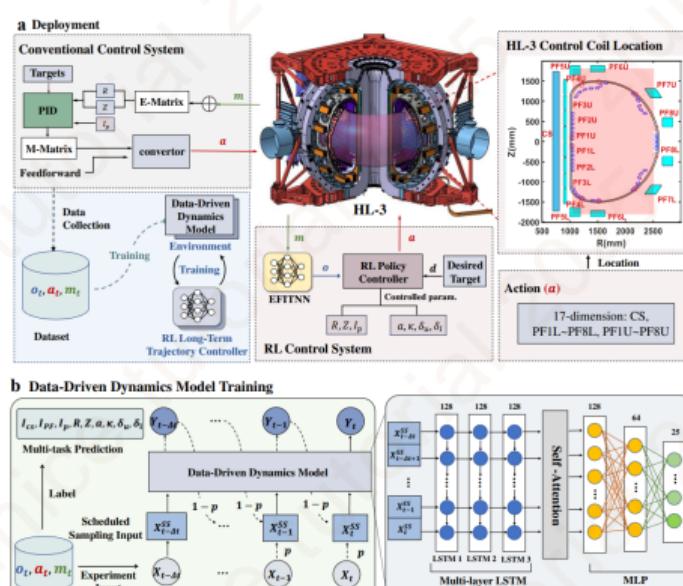
```

Branched rollout: augmenting real experience with short model rollouts.

J. M. F. J. et al., "When to trust your model: Model-based policy optimization", Advances in neural information processing systems, 2019, 32.



Example: Data-Driven RL Control for HL-3 Tokamak



N. Wu, Z. Yang, et al., "High-Fidelity Data-Driven Dynamics Model for Reinforcement Learning-based Control in HL-3 Tokamak," Commun. Phys., Early access.

Algorithm:

Data-driven dynamics model + RL controller for plasma current & shape control.

How It Mitigates Error Accumulation:

- **High-fidelity dynamics model:** replaces costly physics simulators.
- **Scheduled Sampling:** blend real and model inputs, reduce exposure bias.
- **Multi-task prediction + Self-attention:** capture correlations & long-term dependencies.
- **Model ensemble & regularization:** smooth prediction errors.

Result: Stable 400ms, 1kHz trajectory control on HL-3, robust to parameter variations.

Out-of-Distribution (OOD) Problem in Model-Based RL



Cause:

- Policies may visit state-action pairs not covered by the training dataset.
- Dynamics models are unreliable in unseen regions \Rightarrow inaccurate predictions.
- Leads to unstable learning and degraded policy performance.

Solutions:

- **Policy Constraints** Constrain learned policy to stay close to the behavior policy (e.g., *BCQ*, *BRAC*).
- **Reward Penalty / Conservative Modeling** Penalize OOD regions to ensure safety (e.g., *MOPO*, *MOReL*, *MOBILE*).
- **Adversarial / Discriminator-based Methods** Train a discriminator to detect in-distribution vs. OOD states (e.g., *MOREC*).
- **Conservative Value Learning** Explicitly lower Q-values in uncertain regions to avoid overestimation (e.g., *CQL*).



Offline vs. Online RL

Reinforcement Learning with Online Interactions



Offline Reinforcement Learning

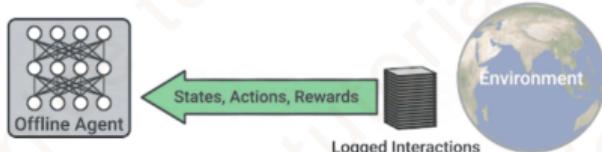


Figure: Illustration of Online vs. Offline RL

Aspect	Offline RL	Online RL
Data Source	Pre-collected and fixed datasets; no interaction with the environment during training.	Real-time interaction with the environment; data dynamically collected during training.
Training Method	Relies only on offline data for model learning and policy optimization; requires accurate dynamics model.	Combines model simulation with real interaction; real-time data helps correct model errors and improve strategies.



Model-based Offline RL

- Offline MBRL allows the dynamics model to generate synthetic data, extending adaptability and exploration beyond the limited offline dataset.
- However, due to the limited coverage of offline data, learned models often deviate from the true environment. Without constraints, model-generated samples may cause **out-of-distribution (OOD)** issues.

Key Question: How to optimize policies using imperfect models?

Solution: Quantify model uncertainty to mitigate OOD problems.

Algorithm	Innovation
MOReL (2020)	Constructs a pessimistic MDP (P-MDP) from offline data, penalizing uncertain transitions.
MOPO (2020)	Adds uncertainty-based penalties to the reward function to prevent exploration into OOD regions.
MOBILE (2023)	Improves MOPO by introducing ensemble-based Bellman error penalties .
MOREC (2023)	Inspired by adversarial learning, adds reward penalties proportional to model uncertainty.



MOPO: Model-based Offline Policy Optimization

- **Contribution:** MOPO introduces a **penalty-based reward modification** into model-based offline RL. The goal is to construct robust policies across all regions, not just the support of the behavior policy.
- **Key Idea:** Add penalty terms to **capture model uncertainty**, preventing policies from exploiting model bias. This achieves a trade-off between generalization and robustness.
- **Dynamics Model:** Learned as an **ensemble of neural networks**, each outputting a Gaussian distribution over next states and rewards:

$$\hat{T}_{\theta,\phi}(s_{t+1}, r \mid s_t, a_t) = \mathcal{N}(\mu_\theta(s_t, a_t), \Sigma_\phi(s_t, a_t))$$

- **Reward Penalty:** Penalized reward is defined as:

$$\tilde{r}(s, a) = \hat{r}(s, a) - \lambda \max_{i=1, \dots, N} \|\Sigma_\phi^i(s, a)\|_F$$

where λ is a hyperparameter controlling the penalty magnitude.



MOREL: Model-Based Offline Reinforcement Learning

Key Idea:

- Learn a **pessimistic MDP (P-MDP)** from offline datasets.
- P-MDP provides a lower bound on true MDP performance, mitigating **model-exploitation**.
- Divide state-action space into **known and unknown regions**.
- Penalize unknown regions with **large negative rewards** to ensure safe policies.

Unknown State-Action Detector (USAD):

$$U_{\text{practical}}(s, a) = \begin{cases} \text{Known,} & \text{disc}(s, a) \leq \text{threshold} \\ \text{Unknown,} & \text{disc}(s, a) > \text{threshold} \end{cases}$$

$$\text{disc}(s, a) = \max_{i,j} \|f_{\phi_i}(s, a) - f_{\phi_j}(s, a)\|_2$$

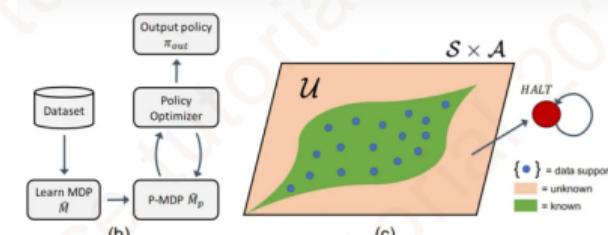


Illustration of MOREL: dataset → learned MDP → P-MDP with known/unknown partition.

K. R. R. A. et al., “Morel: Model-based offline reinforcement learning,” Advances in neural information processing systems, 2020, 33: 21810–21823.

MOBILE: Model-Bellman Inconsistency Penalized Offline Policy Optimization



Motivation If the **Bellman target** for $Q(s, a)$ varies significantly across dynamics models, the action a at state s is unreliable and should be penalized.

Setup Learn an **ensemble** of N dynamics models $\{\hat{T}_{\theta_i}\}_{i=1}^N$ from the offline dataset. Define the model-based Bellman operator under policy π :

$$\hat{\mathcal{T}}^\pi Q_\psi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \hat{T}_\theta(\cdot|s, a), a' \sim \pi(\cdot|s')}[Q_\psi(s', a')].$$

Bellman Inconsistency (Intuition)

$$\epsilon(s, a) = |\hat{\mathcal{T}}^\pi Q_\psi(s, a) - \mathcal{T}^\pi Q_\psi(s, a)|$$

True \mathcal{T}^π is unknown, but inconsistency can be *estimated* by the variation across ensemble models.

S. Y. Z. J. et al., "Model-Bellman inconsistency for model-based offline reinforcement learning," in *International Conference on Machine Learning (ICML)*, PMLR, 2023, pp. 33177–33194.

MOBILE: Model-Bellman Inconsistency Penalized Offline Policy Optimization



Model-Bellman Inconsistency (MBI) Use the standard deviation of Bellman targets across models:

$$U(s, a) := \text{Std}\left(\hat{T}_{\theta_i}^{\pi} Q_{\psi}(s, a)\right) = \sqrt{\frac{1}{N} \sum_{i=1}^N \left(\hat{T}_{\theta_i}^{\pi} Q_{\psi}(s, a) - \overline{\hat{T}^{\pi} Q_{\psi}(s, a)}\right)^2}.$$

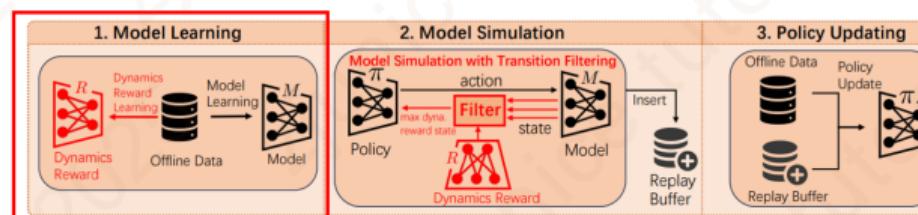
Penalty Shaping Modify target value with penalty term:

$$\hat{y}^{\text{MOBIP}}(s, a) = \hat{T}^{\pi} Q_{\psi}(s, a) - \beta U(s, a),$$

where $\beta > 0$ controls the trade-off between exploitation and conservatism.



MOREC: Model-based Offline RL with Reward Consistency



Step 1: Model Learning

- **Dynamics Model:** $P: \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}'$ Learn dynamics model from offline data as the RL environment:

$$\max_{\theta} L_M(\theta) := \mathbb{E}_{(s, a, s') \sim \mathcal{D}} [\log P_{\theta}(s' | s, a)]$$

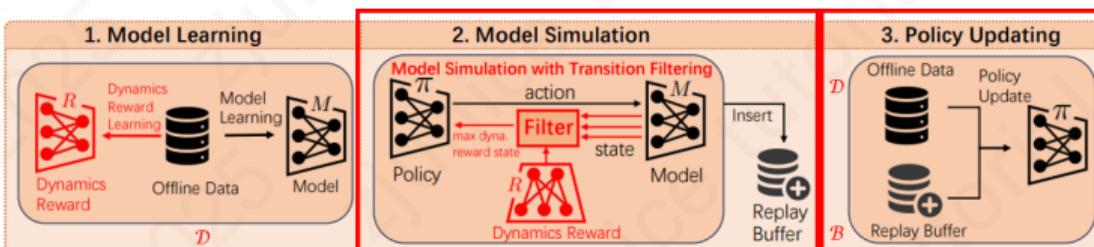
- **Discriminator:** $D: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow (0, 1)$ – With P fixed: maximize distinction between real and model-generated transitions. – With D fixed: minimize divergence to align with offline data.

$$\max_D f(D) = \max_D \min_P \left\{ \mathbb{E}_{(s, a, s') \sim \mathcal{D}} [\log D(s, a, s')] + \mathbb{E}_{(s, a, s') \sim \mathcal{P}} [\log (1 - D(s, a, s'))] \right\}$$

- **Dynamics Reward:** $r^D: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ Define reward based on discriminator confidence:

$$r^D(s, a, s') = -\log(1 - D(s, a, s'))$$

Reward reflects how consistent a transition is with offline data.



Step 2: Model Simulation (choose transitions s' based on $r^D(s, a, s')$)

- 1 Randomly select an initial state s_t from offline dataset \mathcal{D} , and sample action $a_t \sim \pi(\cdot | s_t)$.
- 2 The dynamics model outputs transition distribution $P(s_{t+1} | s_t, a_t)$, using an ensemble $P = \{P_{\theta_m}(s' | s, a), m \in M\}$ with M models.
- 3 At each step, each model P_{θ_m} generates N candidate next states s_{t+1}^m .
- 4 Reassign probabilities over all $M \times N$ candidate next states:

$$P(m, n) = \frac{\exp(r^D(s_t, a_t, s_{t+1}^{m,n})/\kappa)}{\sum_{m \in M} \sum_{n \in N} \exp(r^D(s_t, a_t, s_{t+1}^{m,n})/\kappa)}$$

- 5 Sample s_{t+1} based on $P(m, n)$, obtain trajectory $\{s_t, a_t, r_t, s_{t+1}, \dots\}$, and store it in replay buffer \mathcal{B}_0 .

Step 3: Policy Updating

Sample from offline data \mathcal{D} and replay buffer \mathcal{B} to update policy π and Q using SAC/PPO/TD3, etc.



Avoid explicit uncertainty estimation

RANMBO: Robust Adversarial Model-Based Offline RL

- Introduces **Robust Adversarial RL(RARL)** on top of MBRL.
- Models the **worst-case transition dynamics** within a bounded set.
- Learns policies that remain safe and stable without subjective uncertainty quantification.

COMBO: Conservative Offline MBRL

- Extends **Conservative Q-learning** to model-based RL.
- Penalizes OOD state-action pairs by assigning **lower Q-values**.
- Uses **interpolation mechanism**: mixes real and model-generated rollouts to stabilize learning.
- Yields more reliable value estimates without explicit uncertainty modeling.

Key Idea: Both methods bypass subjective uncertainty estimation, instead relying on **adversarial dynamics (RANMBO)** or **conservative value regularization (COMBO)** to handle OOD.

R. Model-based, L. B., et al., "Rambo-rl: Robust adversarial model-based offline reinforcement learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022, pp. 16082–16097.

Y. Tang, K. Agarwal, et al., "COMBO: Conservative offline model-based policy optimization," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021, pp. 28954–28967.

Conclusion of MBRL



Conclusion of MBRL

Key Difference from Model-Free RL: Model-based RL explicitly learns a full **MDP** model $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ from data, while model-free RL bypasses model learning.

Model Learning: Neural networks, ensembles, Bayesian models, VAEs, GANs.

Action Selection:

- **Planning:** Use the learned model to predict future states and optimize decisions (e.g., DP, MPC).
- **Hybrid learning:** Combine real and model-generated trajectories for policy updates.

Strengths:

- High **sample efficiency**, suitable for costly or risky environments.
- Strong **generalization**, transferable across diverse tasks.

Weaknesses:

- **Error accumulation:** Small model errors compound during long rollouts.
- Vulnerable to **OOD states**, leading to unreliable or unstable policies.
- Trade-off between **exploration and exploitation**.

Thank you

Networked Intelligence for Comprehensive Efficiency (NICE) Lab
College of Information Science and Electronic Engineering

Zhejiang University
<https://nice.rongpeng.info>