*A Provably Better Offline RL Algorithm*
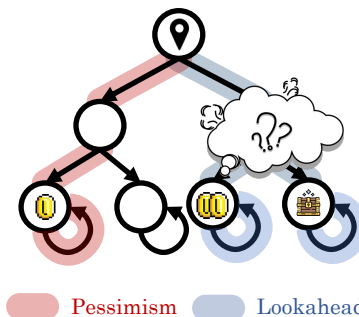
# Don't be Pessimistic Too Early: Look $K$ Steps Ahead!

**Chaoqi Wang**
University of Chicago

**Ziyu Ye**
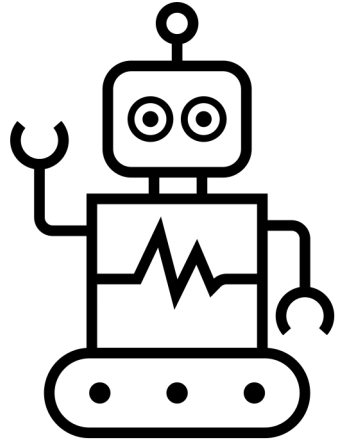University of Chicago

**Kevin Murphy**
Google Research

**Yuxin Chen**
University of Chicago

Pessimism  Lookahead

# Agenda

**1** Introduction
- o  Offline RL & the Pessimism Principle
- o  The Excessive Pessimism Dilemma
- o  Our Contributions

**2** Methods
- o  Formal statement
- o  Lookahead Pessimistic MDP (LP-MDP)

**3** Theoretical Analysis
- o  The Suboptimality Bound
- o  The Effect of Lookahead Horizon

**4** Experimental Results

**5** Takeaways and Future Works
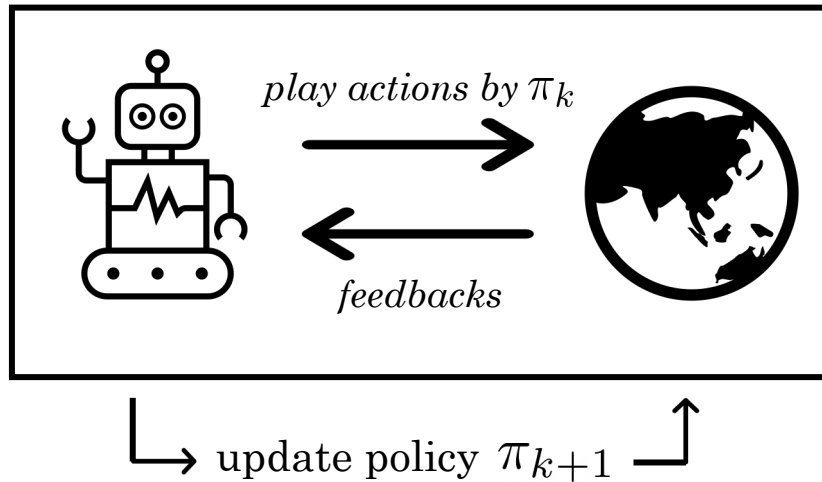
# Part 1
# Introduction

# Reinforcement Learning (RL)

RL is about training agents to **learn** to **make sequential decisions** to achieve **goals**.

# Reinforcement Learning (RL)

RL is about training agents to **learn** to **make sequential decisions** to achieve **goals**.
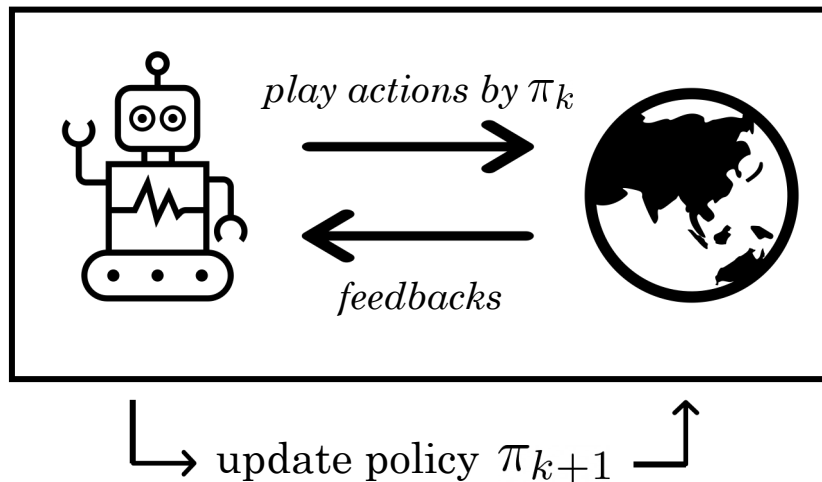
**Classical RL is Online**

# Reinforcement Learning (RL)

RL is about training agents to **learn** to **make sequential decisions** to achieve **goals**.

**Classical RL is Online**



*play actions by* $\pi_k$

*feedbacks*

update policy $\pi_{k+1}$

⚠ requires online interactions



unsafe        time-consuming        costly
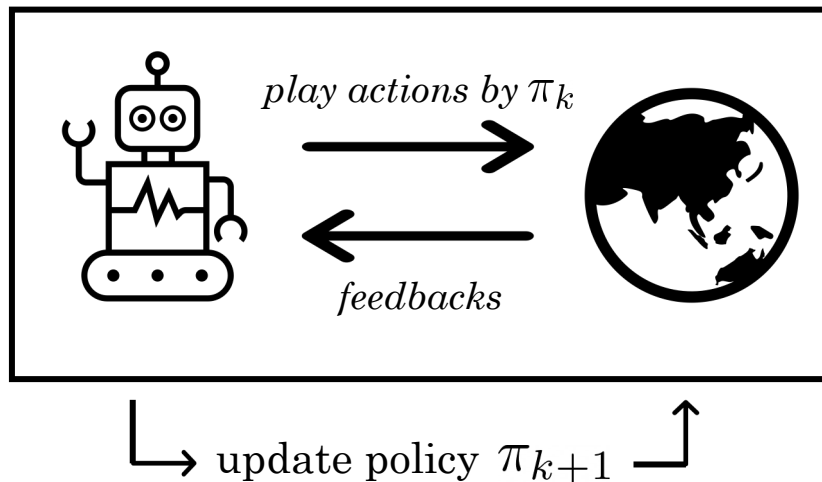
# Reinforcement Learning (RL)

RL is about training agents to **learn** to **make sequential decisions** to achieve **goals**.

**Classical RL is Online**



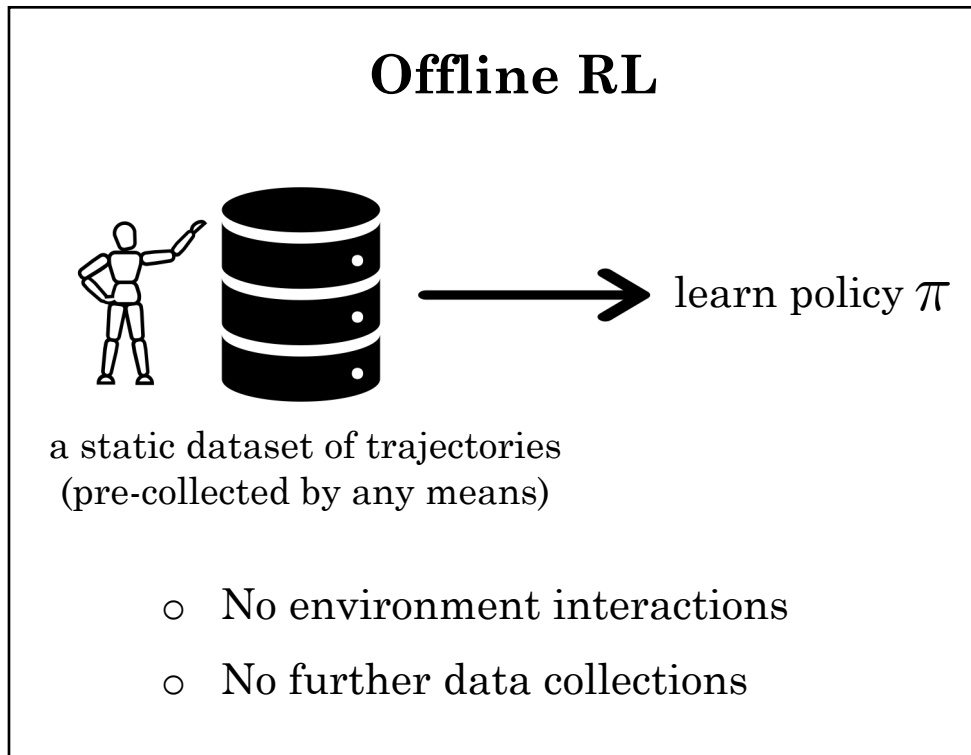⚠ requires online explorations



unsafe    time-consuming    costly

In reality, we often have massive pre-collected data (e.g., by human demonstration).

# Can we still train RL policies *without* any online explorations?

# Offline RL



**Offline RL**

learn policy $\pi$

a static dataset of trajectories
(pre-collected by any means)

○ No environment interactions

○ No further data collections

**Formally…**

$$\mathcal{D} = \{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}_i', r_i)\}$$
$$\mathbf{s} \sim d^{\pi_\beta}(\mathbf{s})$$
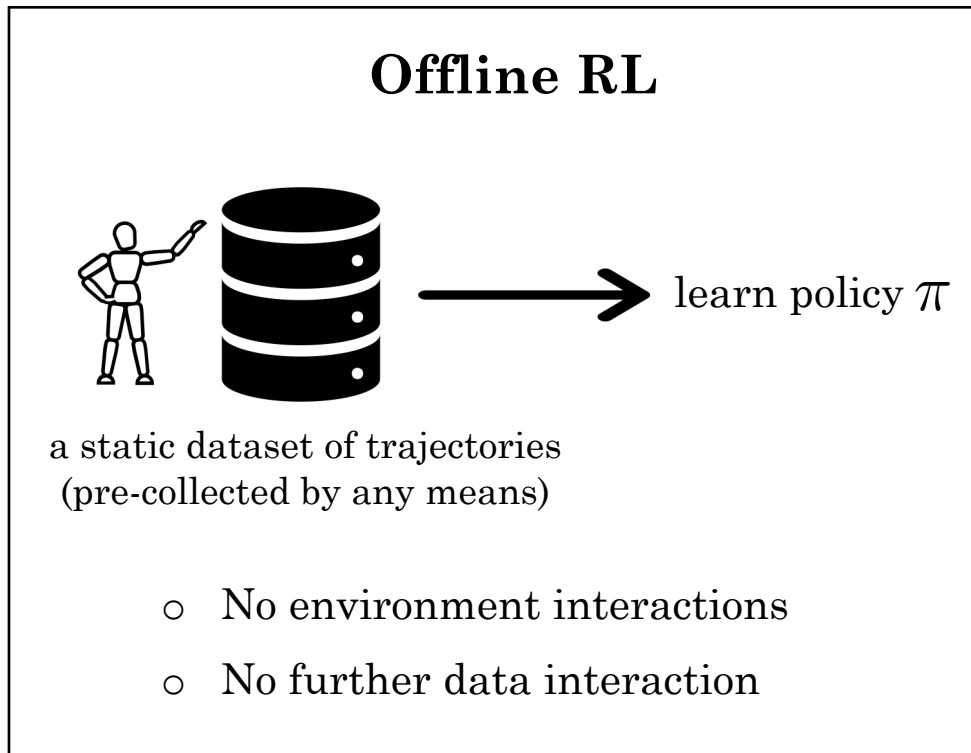$$\mathbf{a} \sim \pi_\beta(\mathbf{a} \mid \mathbf{s})$$
$$\mathbf{s}' \sim p(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})$$
$$r \leftarrow r(\mathbf{s}, \mathbf{a})$$

**Offline RL Objective**

$$\max_\pi \sum_{t=0}^{T} E_{\mathbf{s}_t \sim d^\pi(\mathbf{s}), \mathbf{a}_t \sim \pi(\mathbf{a}|\mathbf{s})} \left[\gamma^t r(\mathbf{s}_t, \mathbf{a}_t)\right]$$

# Offline RL: *Challenges*

## Offline RL



learn policy $\pi$

a static dataset of trajectories
(pre-collected by any means)

- o   No environment interactions
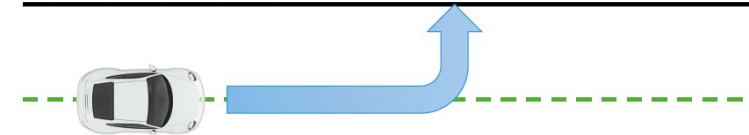- o   No further data interaction

**Fundamental challenge**: *erroneous extrapolation.*

**Training data only contains:**



**What the policy wants to do...**



Good?
Bad?

Online RL can tackle this by trial & error.

How may offline RL deal with potential **out-of-distribution** actions?

# The Pessimism Principle in the Face of Uncertainty

**Key idea**: *avoiding uncertain state & actions* by <u>explicit penalization</u>.

- Wang et al. (2020) regularize the learned policy.
- Kostrikov et al. (2022) penalize the rewards..
- Kidambi et al. (2020) truncate transitions.

# The Pessimism Principle in the Face of Uncertainty



Figure 1. An illustration of P̲essimistic M̲arkov D̲ecision P̲rocess (P-MDP) by Kidambi et al. (2020).
Notice that the **value** of any policy in P-MDP will be the **lower bound** for the true value

**Key idea**: *avoiding uncertain state & actions* by explicit penalization.

- Wang et al. (2020) regularize the learned policy.
- Kostrikov et al. (2022) penalize the rewards..
- Kidambi et al. (2020) truncate transitions.

# What can go wrong for pessimism based algorithms?

# The *Excessive* Pessimism Dilemma

⚠️ The policy may *behave overly conservative*, ends up too *far away* from achievable better performance.

# The *Excessive* Pessimism Dilemma

⚠️ The policy may *behave overly conservative*, ends up too *far away* from achievable better performance.

$r(\text{🪙}) = 1$    $r(\text{🪙🪙}) = 2$    $r(\text{🎁}) = 10$    ☁️ Uncertain region

(Fitted MDP)

# The *Excessive* Pessimism Dilemma

⚠️ The policy may *behave overly conservative*, ends up too *far away* from achievable better performance.

$r(🪙) = 1$          $r(🪙🪙) = 2$          $r(📦) = 10$                    Uncertain region

(Fitted MDP)

Classical Pessimistic MDP
will directly halts here.

No chance to get $> 🪙$ !

# The *Excessive* Pessimism Dilemma

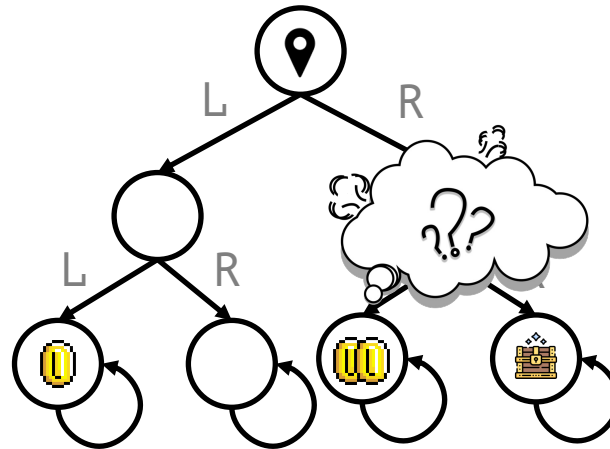⚠ The policy may *behave overly conservative*, ends up too *far away* from achievable better performance.

$r(\text{🪙}) = 1$     $r(\text{🪙🪙}) = 2$     $r(\text{📦}) = 10$     $r(\text{🚫} or \text{⚫}) = 0$     💭 Uncertain region



(Fitted MDP)

Pessimistic
MDP

(a)

$\pi^*_{\text{P-MDP}} = (L, L)$

Halting too early → get at best 🪙.

━━ Possible Path of $\pi^*_{\text{P-MDP}}$

# Can we find a principled way to *modulate* pessimism?

(And to achieve a better performance guarantee... )

# The *Excessive* Pessimism Dilemma: *Mitigating by Lookahead*

⚠️ The policy may *behave overly conservative*, ends up too *far away* from achievable better performance.

$r(🪙) = 1$     $r(🪙🪙) = 2$     $r(🧰) = 10$     $r(🚫 or ⚪) = 0$     💭 Uncertain region

(Fitted MDP)



Pessimistic MDP

Just look one step ahead...

Fortunes (🪙🪙 or 🧰) await for sure!

$\pi^*_{P-MDP} = (L, L)$

Halting too early → get at best 🪙.

━━ Possible Path of $\pi^*_{P-MDP}$

# The *Excessive* Pessimism Dilemma: *Mitigating by* *Lookahead*

⚠️ The policy may *behave overly conservative*, ends up too *far away* from achievable better performance.

$r(\text{🪙}) = 1$   $r(\text{🪙🪙}) = 2$   $r(\text{📦}) = 10$   $r(\text{🚫} or \text{⚪}) = 0$   💭 Uncertain region



(Fitted MDP)

(a)   Pessimistic MDP   Lookahead Pessimistic MDP (Ours)   (b)

$\pi^*_{P-MDP} = (L, L)$

Halting too early ➔ get at best 🪙.

━━ Possible Path of $\pi^*_{P-MDP}$
━━ Possible Path of $\pi^*_{LP-MDP}$

$\pi^*_{LP-MDP} = (R, L \text{ or } R)$

**Look 1 step ahead** ➔ get 🪙🪙 or better!

💡 **Lookahead** enables a *less conservative* policy with *better performance guarantee*!

# Further Insights on Lookahead Pessimism



Possible Path of $\pi^*_{P-MDP}$

Possible Path of $\pi^*_{LP-MDP}$

✅ The lookahead horizon **modulates the pessimism level**.

✅ Lookahead helps circumvent uncertain areas by **path stitching**.

✅ Lookahead **implicitly increases the data coverage**.

# Our Contributions

**Algorithm**: <u>L</u>ookahead <u>P</u>essimistic MDP.

**Theory**: a lower bound monotonically improves with the lookahead horizon $K$.

**Experiments**: solid improvement over baselines on benchmark environments.
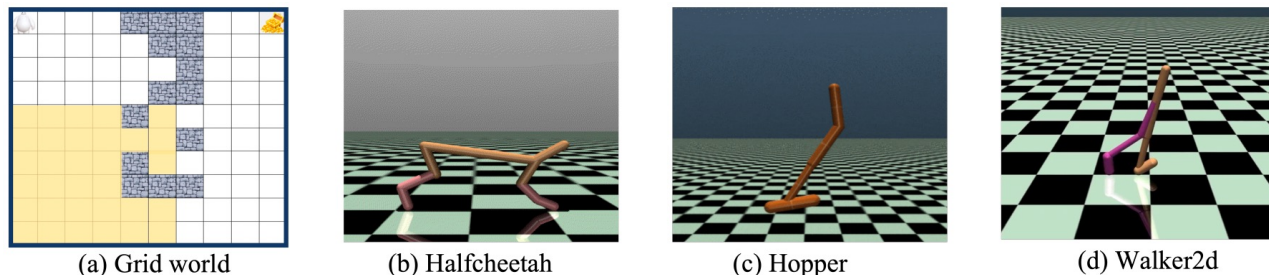
(a) Grid world     (b) Halfcheetah     (c) Hopper     (d) Walker2d

| Dataset | Environment | LP-MDP (Ours) | MOReL (SAC) | MOReL (NPG) | MOPO | CQL | SAC-Off | BEAR | BRAC-p | BRAC-v |
|---|---|---|---|---|---|---|---|---|---|---|
| medium | halfcheetah | 42.6±5.5 | 43.4 | 42.1 | 54.2 | 44.4 | -4.3 | 41.7 | 43.8 | 46.3 |
| medium | hopper | 101.9 ±1.1* | 75.8 | 95.4 | 28.0 | 86.6 | 0.8 | 52.1 | 32.7 | 31.1 |
| medium | walker2d | 64.5±5.1 | 76.8 | 77.8 | 17.8 | 74.5 | 0.9 | 59.1 | 77.5 | 81.1 |
| medium-replay | halfcheetah | 48.5±2.1* | 43.4 | 40.2 | 53.1 | 46.2 | -2.4 | 38.6 | 45.4 | 47.7 |
| medium-replay | hopper | 101.2 ±0.8 | 101.1 | 93.6 | 67.5 | 48.6 | 3.5 | 33.7 | 0.6 | 0.6 |
| medium-replay | walker2d | 82.7 ±5.9* | 46.5 | 49.8 | 39.0 | 32.6 | 1.9 | 19.2 | -0.3 | 0.9 |
| medium-expert | halfcheetah | 51.1±0.9* | 41.6 | 53.3 | 63.3 | 62.4 | 1.8 | 53.4 | 44.2 | 41.9 |
| medium-expert | hopper | 103.7±1.2* | 78.5 | 108.7 | 23.7 | 111.0 | 1.6 | 96.3 | 1.9 | 0.8 |
| medium-expert | walker2d | 81.5±8.5* | 68.0 | 95.6 | 44.6 | 98.7 | -0.1 | 40.1 | 76.9 | 81.6 |
| | Average Scores | 677.7 * | 575.1 | 656.5 | 391.2 | 605.0 | 3.7 | 434.2 | 328.1 | 332.0 |

*Table 1.* Results on D4RL. We report the averaged normalized scores of 3 different random seeds with one standard errors. We highlight the best averaged scores by a blue box. Also, we use * to indicate the tasks where our method has a solid improvement over MOReL-SAC.

---

**Algorithm 1** Policy Learning under Pessimistic MDP with look-ahead (LP$(\xi, \pi, K)$-MDP).

**Require:** Offline data $\mathcal{D}$, threshold $\xi$, look-ahead steps $K$, and learning rate $\alpha$.
1: Fit the dynamics model $\mathcal{M}_{\widehat{p}}$ on $\mathcal{D}$
2: Initialize the policy $\pi_0$
3: **for** $n = 0, 1, 2, \dots$ **do**
4:     Construct the LP$(\xi, \pi, K)$-MDP for $\pi_n$: $\mathcal{M}_{\widehat{p}}^{\pi_n}$
5:     Improve policy: $\pi_{n+1} \leftarrow \texttt{SAC\_Step}(\pi_n, \mathcal{M}_{\widehat{p}}^{\pi_n})$
6: **end for**
7: **return** $\pi_n$

---

**Pessimistic MDP with $K$-step Lookahead**

**Definition 3 (LP$(\xi, \pi, K)$-MDP)** *For any* $(\boldsymbol{s}_t, \boldsymbol{a}_t)$, *and* $\mathcal{M}_{\widehat{p}}$ *and* $\mathcal{M}_p$, *the LP$(\xi, \pi, K)$-MDP of* $\mathcal{M}_p$, *(i.e.,* $\mathcal{M}_{\widehat{p}}^{\pi}$*) is constructed by modifying the transition* $\widehat{p}$ *to be* $\widetilde{p}$ *as:[a]*

**Case 1** *(current transition is* certain*) If* $(\boldsymbol{s}_t, \boldsymbol{a}_t) \notin \mathcal{U}$, *then*

$$\widetilde{p}(\cdot|\boldsymbol{s}_t, \boldsymbol{a}_t) = \widehat{p}(\cdot|\boldsymbol{s}_t, \boldsymbol{a}_t), \qquad (9)$$

**Case 2** *(current transition is* uncertain*) If* $(\boldsymbol{s}_t, \boldsymbol{a}_t) \in \mathcal{U}$, *then*

**Case 2.1** *(all $K$-step look ahead is uncertain)* *If* $(\boldsymbol{s}_t, \boldsymbol{a}_t) \in \mathcal{U}_{-(1:K)}^{\pi}$, *then*

$$\widetilde{p}(\boldsymbol{s}_{t+1} = \mathbf{e}|\boldsymbol{s}_t, \boldsymbol{a}_t) = 1, \qquad (10)$$

**Case 2.2** *(some $k_{\text{th}}$-step look ahead is certain)* *If* $(\boldsymbol{s}_t, \boldsymbol{a}_t) \in \mathcal{U}_k^{\pi}$ *for some* $k \in [K]$, *then we construct a deterministic path such that* $\forall i \in \{1, ..., k-1\}$,

$$\widetilde{p}(\boldsymbol{s}_{t+k} = \boldsymbol{s}^{\star}|\boldsymbol{s}_t, \boldsymbol{a}_t, \pi) = 1, \ \widetilde{r}(\boldsymbol{s}_{t+i}, \cdot) = -R_{\max}. \quad (11)$$

*where* $\boldsymbol{s}^{\star}$ *is defined as:*

$$\boldsymbol{s}^{\star} = \operatorname*{arg\,min}_{\boldsymbol{s}' \in \mathcal{L}_{\mathcal{M}_{\widehat{p}}}^{\pi, k}(\boldsymbol{s}_t, \boldsymbol{a}_t)} V_{\mathcal{M}_{\widehat{p}}^{\pi}}^{\pi}(\boldsymbol{s}'), \qquad (12)$$

[a]The associated reward $\widetilde{r}(\cdot) := r(\cdot)$ unless otherwise stated.

We are all in the gutter, but some of us are *looking at the stars*.

— *Oscar Wilde*

# Part 2
# Methodology

# Preliminaries

## 🔭 **Intuition**

If the current state-action pair is *uncertain* → don't just halt!

Look a few steps ahead → if future states is promising w/ *high certainty*, go for it!

# Preliminaries

### 🔭 Intuition

If the current state-action pair is *uncertain* → don't just halt!

Look a few steps ahead → if future states is promising w/ *high certainty*, go for it!

### Uncertainty Quantification

We first denote the *estimation error* on any pair $(s, a)$ as

$$d(s, a) := \mathbf{d}_{\text{TV}}\left(\widehat{p}(\cdot|s, a)|\, p(\cdot|s, a)\right),$$

which quantifies the total variation distance from the estimated distribution $\widehat{p}$ to the true distribution $p$.

# Preliminaries

 **Intuition**

If the current state-action pair is *uncertain* → don't just halt!

Look a few steps ahead → if future states is promising w/ *high certainty*, go for it!

### Uncertainty Quantification

We first denote the *estimation error* on any pair $(s, a)$ as

$$d(s, a) := \mathbf{d}_{\mathrm{TV}}\left(\widehat{p}(\cdot|s, a) \,\middle|\, p(\cdot|s, a)\right),$$

which quantifies the total variation distance from the estimated distribution $\widehat{p}$ to the true distribution $p$.

### Classical Pessimism asks to halt at ...

**Definition 1 (Uncertain State-Action Set)** $\forall \, \xi \geq 0$, *the set of uncertain state-action pairs is*

$$\mathcal{U}(\xi) := \{(s, a) \in \mathcal{S} \times \mathcal{A} : d(s, a) \geq \xi\}. \tag{7}$$

For simplicity, we drop the dependency on $\xi$ in the notations.

# Partition Uncertain Regions

🔭 We can further **partition** $U$ by properties of some **lookahead sets**.

**Associate Each Pair with Lookahead Sets**

**Definition 2 (Lookahead Set)** *For any $(s_t, a_t) \in \mathcal{S} \times \mathcal{A}$ under MDP $\mathcal{M}_{p'}$ with $p'(\cdot|\cdot, \cdot, \pi)$ denoting the transition distribution relying on $\pi$, the $k_{\text{th}}$-step lookahead set is*
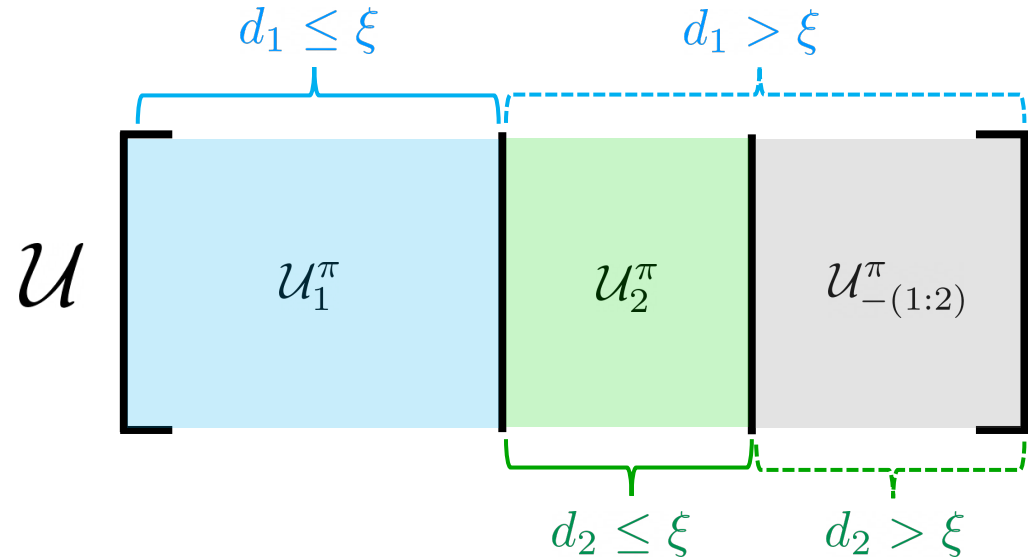
$$\mathcal{L}^{\pi,k}_{\mathcal{M}_{p'}}(s_t, a_t) := \{s \in \mathcal{S}|\ p'(s_{t+k} = s|s_t, a_t, \pi) > 0\}.$$



$$d_1 \leq \xi \qquad d_1 > \xi$$

$\mathcal{U}$

$\mathcal{U}^\pi_1 \qquad \mathcal{U}^\pi_2 \qquad \mathcal{U}^\pi_{-(1:2)}$

$$d_2 \leq \xi \qquad d_2 > \xi$$

**Partition Criteria: Lookahead Certainty**

We say a state-action pair $(s, a)$ is $k_{\text{th}}$-**certain** if for all the states in its $k_{\text{th}}$-step lookahead sets, the fitted dynamics $\mathcal{M}_{\widehat{p}}$ induces only a small estimation error, that is:

$$\forall s' \in \mathcal{L}^{\pi,k}_{\mathcal{M}_{\widehat{p}}}(s, a),\ d(s', \pi(s')) \leq \xi, \qquad (8)$$

where $d(s, \pi(s)) := \max_{a:\pi(a|s)>0} d(s, a)$.

The set $\mathcal{U}$ thus can be partitioned into disjoint subsets by the above lookahead certainty criteria. For all $k \in [1, K]$, we define the subset $\mathcal{U}^\pi_k$ as

$$\mathcal{U}^\pi_k := \{(s, a) \in \mathcal{U} \setminus \cup^{k-1}_{i=1}\mathcal{U}^\pi_i : (s, a) \text{ is } k_{\text{th}}\text{-certain}\}.$$

We further use $\mathcal{U}^\pi_{-(1:K)} := \mathcal{U} \setminus \cup^K_{i=1}\mathcal{U}^\pi_i$ to denote all the remaining state-action pairs in $\mathcal{U}$.

# Construct <u>L</u>ookahead <u>P</u>essimistic <u>MDP</u>

### Pessimistic MDP with $K$-step Lookahead

**Definition 3 (LP$(\xi, \pi, K)$-MDP)** *For any* $(\boldsymbol{s}_t, \boldsymbol{a}_t)$, *and* $\mathcal{M}_{\widehat{p}}$ *and* $\mathcal{M}_p$, *the LP$(\xi, \pi, K)$-MDP of* $\mathcal{M}_p$, *(i.e.,* $\mathcal{M}_{\widetilde{p}}^{\pi}$*) is constructed by modifying the transition* $\widehat{p}$ *to be* $\widetilde{p}$ *as:[a]*

<u>**Case 1**</u> *(current transition is* certain*) If* $(\boldsymbol{s}_t, \boldsymbol{a}_t) \notin \mathcal{U}$, *then*

$$\widetilde{p}(\cdot|\boldsymbol{s}_t, \boldsymbol{a}_t) = \widehat{p}(\cdot|\boldsymbol{s}_t, \boldsymbol{a}_t), \tag{9}$$

<u>**Case 2**</u> *(current transition is* uncertain*) If* $(\boldsymbol{s}_t, \boldsymbol{a}_t) \in \mathcal{U}$, *then*

  <u>**Case 2.1**</u> *(all $K$-step look ahead is uncertain)*
  *If* $(\boldsymbol{s}_t, \boldsymbol{a}_t) \in \mathcal{U}_{-(1:K)}^{\pi}$, *then*

$$\widetilde{p}(\boldsymbol{s}_{t+1} = \mathbf{e}|\boldsymbol{s}_t, \boldsymbol{a}_t) = 1, \tag{10}$$

<u>**Case 2.2**</u> *(some $k_{\text{th}}$-step look ahead is certain)*
*If* $(\boldsymbol{s}_t, \boldsymbol{a}_t) \in \mathcal{U}_k^{\pi}$ *for some* $k \in [K]$, *then we construct a deterministic path such that* $\forall\, i \in \{1, ..., k-1\}$,

$$\widetilde{p}(\boldsymbol{s}_{t+k} = \boldsymbol{s}^{\star}|\boldsymbol{s}_t, \boldsymbol{a}_t, \pi) = 1, \; \widetilde{r}(\boldsymbol{s}_{t+i}, \cdot) = -R_{\max}. \tag{11}$$

*where* $\boldsymbol{s}^{\star}$ *is defined as:*

$$\boldsymbol{s}^{\star} = \underset{\boldsymbol{s}' \in \mathcal{L}_{\mathcal{M}_{\widehat{p}}}^{\pi, k}(\boldsymbol{s}_t, \boldsymbol{a}_t)}{\arg\min} V_{\mathcal{M}_{\widetilde{p}}^{\pi}}^{\pi}(\boldsymbol{s}'), \tag{12}$$

_____
*[a]The associated reward* $\widetilde{r}(\cdot) := r(\cdot)$ *unless otherwise stated.*



**Classical Pessimism Principle**

- Halts at the uncertain state $s_t$

**Lookahead Pessimism**

- Constructs a *less conservative path*: $s_t \rightarrow s^{\star}$
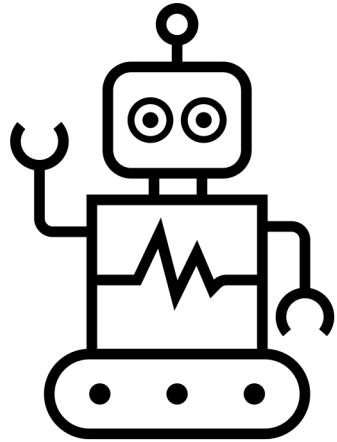- Has a *better worst-case* guarantee

# The Algorithm

---

**Algorithm 1** Policy Learning under Pessimistic MDP with look-ahead (LP($\xi, \pi, K$)-MDP).

---

**Require:** Offline data $\mathcal{D}$, threshold $\xi$, look-ahead steps $K$, and learning rate $\alpha$.

1: Fit the dynamics model $\mathcal{M}_{\widehat{p}}$ on $\mathcal{D}$
2: Initialize the policy $\pi_0$
3: **for** $n = 0, 1, 2, \ldots$ **do**
4:     Construct the LP($\xi, \pi, K$)-MDP for $\pi_n$: $\mathcal{M}_{\widetilde{p}}^{\pi_n}$
5:     Improve policy: $\pi_{n+1} \leftarrow \texttt{SAC\_Step}(\pi_n, \mathcal{M}_{\widetilde{p}}^{\pi_n})$
6: **end for**
7: **return** $\pi_n$

---

# Part 3
# Theoretical Analysis

# The Suboptimality Bound

**Theorem 1 (Performance of the Equilibrium Policy)** *Let $\widetilde{\pi}^\star$ denote the equilibrium policy learned under the $LP(\xi, \widetilde{\pi}^\star, K)$-MDP, and let $\pi^\star$ be the optimal policy under the true MDP $\mathcal{M}_p$. Suppose that for any $(s_t, a_t) \in \mathcal{U}_k^{\widetilde{\pi}^\star}$ with $k \in [1, K]$, $\mathcal{L}_{\mathcal{M}_p}^{\widetilde{\pi}^\star, k}(s_t, a_t) \subseteq \mathcal{L}_{\mathcal{M}_{\widehat{p}}}^{\widetilde{\pi}^\star, k}(s_t, a_t)$ holds, then for any state $s$*

$$V_{\mathcal{M}_p}^{\pi^\star}(s) - V_{\mathcal{M}_p}^{\widetilde{\pi}^\star}(s)$$

$$\leq \underbrace{\frac{4\gamma\xi R_{\max}}{(1-\gamma)^2}}_{(a)} + \underbrace{\frac{2\rho_{s\to\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}^{\pi^\star} \mathbb{E}\left[\gamma^{\mathcal{T}_{s\to\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}^{\pi^\star}}\right] R_{\max}}{1-\gamma}}_{(b)\text{ Incurred by hitting }\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}} \tag{15}$$

$$+ \underbrace{\sum_{k=1}^{K} \rho_{s\to\mathcal{U}_k^{\widetilde{\pi}^\star}}^{\pi^\star} \mathbb{E}\left[\gamma^{\mathcal{T}_{s\to\mathcal{U}_k^{\widetilde{\pi}^\star}}^{\pi^\star}}\right]\left(\sum_{i=1}^{k-1} 2\gamma^i R_{\max} + \gamma^k \Delta_{p,\widetilde{q}}^{\pi^\star, k}(\mathcal{U}_k^{\widetilde{\pi}^\star})\right)}_{(c)\text{ Incurred by hitting }\mathcal{U}_k^{\widetilde{\pi}^\star}\text{ for }k\in[K]}.$$

# The Suboptimality Bound

**Theorem 1 (Performance of the Equilibrium Policy)** *Let $\widetilde{\pi}^\star$ denote the equilibrium policy learned under the $LP(\xi, \widetilde{\pi}^\star, K)$-MDP, and let $\pi^\star$ be the optimal policy under the true MDP $\mathcal{M}_p$. Suppose that for any $(s_t, a_t) \in \mathcal{U}_k^{\widetilde{\pi}^\star}$ with $k \in [1, K]$, $\mathcal{L}_{\mathcal{M}_p}^{\widetilde{\pi}^\star, k}(s_t, a_t) \subseteq \mathcal{L}_{\mathcal{M}_{\widehat{p}}}^{\widetilde{\pi}^\star, k}(s_t, a_t)$ holds, then for any state $s$*

$$V_{\mathcal{M}_p}^{\pi^\star}(s) - V_{\mathcal{M}_p}^{\widetilde{\pi}^\star}(s)$$

hitting the **known region** ⋘

$$\leq \underbrace{\frac{4\gamma\xi R_{\max}}{(1-\gamma)^2}}_{(a)} + \underbrace{\frac{2\rho_{s\to\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}^{\pi^\star} \mathbb{E}\left[\gamma^{\mathcal{T}_{s\to\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}^{\pi^\star}}\right] R_{\max}}{1-\gamma}}_{(b) \text{ Incurred by hitting } \mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}$$

$$+ \underbrace{\sum_{k=1}^{K} \rho_{s\to\mathcal{U}_k^{\widetilde{\pi}^\star}}^{\pi^\star} \mathbb{E}\left[\gamma^{\mathcal{T}_{s\to\mathcal{U}_k^{\widetilde{\pi}^\star}}^{\pi^\star}}\right] \left(\sum_{i=1}^{k-1} 2\gamma^i R_{\max} + \gamma^k \Delta_{p,\widetilde{q}}^{\pi^\star, k}(\mathcal{U}_k^{\widetilde{\pi}^\star})\right)}_{(c) \text{ Incurred by hitting } \mathcal{U}_k^{\widetilde{\pi}^\star} \text{ for } k\in[K]}$$

# The Suboptimality Bound

**Theorem 1 (Performance of the Equilibrium Policy)** *Let $\widetilde{\pi}^\star$ denote the equilibrium policy learned under the $LP(\xi, \widetilde{\pi}^\star, \widetilde{K})$-MDP, and let $\pi^\star$ be the optimal policy under the true MDP $\mathcal{M}_p$. Suppose that for any $(s_t, a_t) \in \mathcal{U}_k^{\widetilde{\pi}^\star}$ with $k \in [1, K]$, $\mathcal{L}_{\mathcal{M}_p}^{\widetilde{\pi}^\star, k}(s_t, a_t) \subseteq \mathcal{L}_{\mathcal{M}_{\widehat{p}}}^{\widetilde{\pi}^\star, k}(s_t, a_t)$ holds, then for any state $s$*

$$V_{\mathcal{M}_p}^{\pi^\star}(s) - V_{\mathcal{M}_p}^{\widetilde{\pi}^\star}(s)$$

hitting the **known region** $\lll$

$$\leq \underbrace{\frac{4\gamma\xi R_{\max}}{(1-\gamma)^2}}_{(a)} + \underbrace{\frac{2\rho_{s\to\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}^{\pi^\star} \mathbb{E}\left[\gamma^{\mathcal{T}_{s\to\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}^{\pi^\star}}\right] R_{\max}}{1-\gamma}}_{\text{(b) Incurred by hitting } \mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}$$

$\ggg$ hitting the **unknown region** & <u>all $K$-step lookahead</u> is **uncertain**

$$+ \underbrace{\sum_{k=1}^{K} \rho_{s\to\mathcal{U}_k^{\widetilde{\pi}^\star}}^{\pi^\star} \mathbb{E}\left[\gamma^{\mathcal{T}_{s\to\mathcal{U}_k^{\widetilde{\pi}^\star}}^{\pi^\star}}\right]\left(\sum_{i=1}^{k-1} 2\gamma^i R_{\max} + \gamma^k \Delta_{p,\widetilde{q}}^{\pi^\star, k}(\mathcal{U}_k^{\widetilde{\pi}^\star})\right)}_{\text{(c) Incurred by hitting } \mathcal{U}_k^{\widetilde{\pi}^\star} \text{ for } k \in [K]}$$

# The Suboptimality Bound

**Theorem 1 (Performance of the Equilibrium Policy)** *Let $\widetilde{\pi}^\star$ denote the equilibrium policy learned under the $LP(\xi, \widetilde{\pi}^\star, \widetilde{K})$-MDP, and let $\pi^\star$ be the optimal policy under the true MDP $\mathcal{M}_p$. Suppose that for any $(s_t, a_t) \in \mathcal{U}_k^{\widetilde{\pi}^\star}$ with $k \in [1, K]$, $\mathcal{L}_{\mathcal{M}_p}^{\widetilde{\pi}^\star, k}(s_t, a_t) \subseteq \mathcal{L}_{\mathcal{M}_{\widehat{p}}}^{\widetilde{\pi}^\star, k}(s_t, a_t)$ holds, then for any state $s$*

$$V_{\mathcal{M}_p}^{\pi^\star}(s) - V_{\mathcal{M}_p}^{\widetilde{\pi}^\star}(s)$$

hitting the **known region** ⫷

$$\leq \underbrace{\frac{4\gamma\xi R_{\max}}{(1-\gamma)^2}}_{(a)} + \underbrace{\frac{2\rho_{s \to \mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}^{\pi^\star} \mathbb{E}\left[\gamma^{\mathcal{T}_{s \to \mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}^{\pi^\star}}\right] R_{\max}}{1-\gamma}}_{\text{(b) Incurred by hitting } \mathcal{U}_{-(1:K)}^{\widetilde{\pi}^\star}}$$

hitting the **unknown region** ⫸
& <u>all $K$-step lookahead is **uncertain**</u>

$$+ \underbrace{\sum_{k=1}^{K} \rho_{s \to \mathcal{U}_k^{\widetilde{\pi}^\star}}^{\pi^\star} \mathbb{E}\left[\gamma^{\mathcal{T}_{s \to \mathcal{U}_k^{\widetilde{\pi}^\star}}^{\pi^\star}}\right]\left(\sum_{i=1}^{k-1} 2\gamma^i R_{\max} + \gamma^k \Delta_{p,\widetilde{q}}^{\pi^\star, k}(\mathcal{U}_k^{\widetilde{\pi}^\star})\right)}_{\text{(c) Incurred by hitting } \mathcal{U}_k^{\widetilde{\pi}^\star} \text{ for } k \in [K]}$$

hitting the **unknown region** ⫸
& <u>some $k$-step lookahead is **certain**</u>

# The Suboptimality Bound

**Theorem 1 (Performance of the Equilibrium Policy)** *Let $\widetilde{\pi}^{\star}$ denote the equilibrium policy learned under the $LP(\xi, \widetilde{\pi}^{\star}, \widetilde{K})$-MDP, and let $\pi^{\star}$ be the optimal policy under the true MDP $\mathcal{M}_p$. Suppose that for any $(s_t, a_t) \in \mathcal{U}_k^{\widetilde{\pi}^{\star}}$ with $k \in [1, K]$, $\mathcal{L}_{\mathcal{M}_p}^{\widetilde{\pi}^{\star}, k}(s_t, a_t) \subseteq \mathcal{L}_{\mathcal{M}_{\widehat{p}}}^{\widetilde{\pi}^{\star}, k}(s_t, a_t)$ holds, then for any state $s$*

$$V_{\mathcal{M}_p}^{\pi^{\star}}(s) - V_{\mathcal{M}_p}^{\widetilde{\pi}^{\star}}(s)$$

hitting the **known region** ⋘

$$\leq \underbrace{\frac{4\gamma\xi R_{\max}}{(1-\gamma)^2}}_{(a)} + \underbrace{\frac{2\rho_{s\to\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^{\star}}}^{\pi^{\star}} \mathbb{E}\left[\gamma^{\mathcal{T}_{s\to\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^{\star}}}^{\pi^{\star}}}\right] R_{\max}}{1-\gamma}}_{(b)\ \text{Incurred by hitting }\mathcal{U}_{-(1:K)}^{\widetilde{\pi}^{\star}}}$$

⋙ hitting the **unknown region** & all $K$-step lookahead is **uncertain**

$$+ \underbrace{\sum_{k=1}^{K} \rho_{s\to\mathcal{U}_k^{\widetilde{\pi}^{\star}}}^{\pi^{\star}} \mathbb{E}\left[\gamma^{\mathcal{T}_{s\to\mathcal{U}_k^{\widetilde{\pi}^{\star}}}^{\pi^{\star}}}\right]\left(\sum_{i=1}^{k-1} 2\gamma^i R_{\max} + \gamma^k \Delta_{p,\widetilde{q}}^{\pi^{\star}, k}(\mathcal{U}_k^{\widetilde{\pi}^{\star}})\right)}_{(c)\ \text{Incurred by hitting }\mathcal{U}_k^{\widetilde{\pi}^{\star}}\ \text{for }k\in[K]}$$

⋙ hitting the **unknown region** & some $k$-step lookahead is **certain**

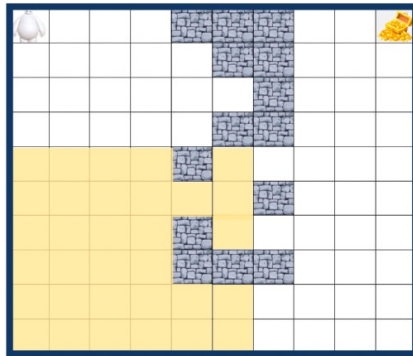💡 Monotonically improves with $K$ → guaranteed **better lower bound** than existing work!
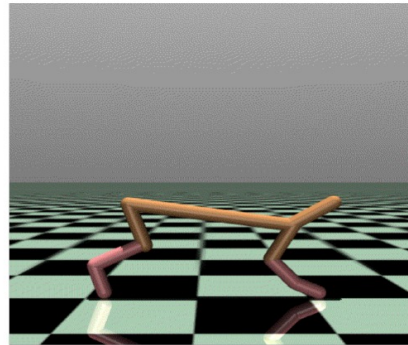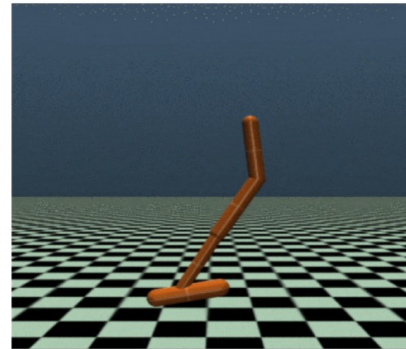
# Part 4
# Experimental Results

# Datasets

Popular benchmark suite used in many papers (Kidambi et al., 2020, etc.).
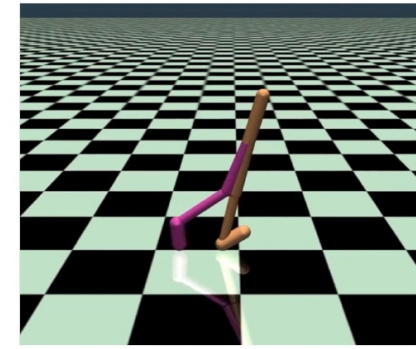


| (a) Grid world | (b) Halfcheetah | (c) Hopper | (d) Walker2d |

*Figure 4.* Visualization of considered tasks. (a) The grid world task is adapted from Eysenbach et al. (2022) with increased difficulty. The agent starts where the baymax is located (top left corner). The reward for hitting the treasure is $+50$, $+1$ for yellow-shaded grids, and $+0.5$ for other grids. The walls cannot be crossed. For (b), (c) and (d), the tasks come from the D4RL benchmark (Fu et al., 2020).

# Results: Performance Improvement

| Dataset | Environment | LP-MDP (Ours) | MOReL (SAC) | MOReL (NPG) | MOPO | CQL | SAC-Off | BEAR | BRAC-p | BRAC-v |
|---|---|---|---|---|---|---|---|---|---|---|
| medium | halfcheetah | $42.6\pm5.5$ | 43.4 | 42.1 | 54.2 | 44.4 | -4.3 | 41.7 | 43.8 | 46.3 |
| medium | hopper | $101.9\pm1.1^*$ | 75.8 | 95.4 | 28.0 | 86.6 | 0.8 | 52.1 | 32.7 | 31.1 |
| medium | walker2d | $64.5\pm5.1$ | 76.8 | 77.8 | 17.8 | 74.5 | 0.9 | 59.1 | 77.5 | 81.1 |
| medium-replay | halfcheetah | $48.5\pm2.1^*$ | 43.4 | 40.2 | 53.1 | 46.2 | -2.4 | 38.6 | 45.4 | 47.7 |
| medium-replay | hopper | $101.2\pm0.8$ | 101.1 | 93.6 | 67.5 | 48.6 | 3.5 | 33.7 | 0.6 | 0.6 |
| medium-replay | walker2d | $82.7\pm5.9^*$ | 46.5 | 49.8 | 39.0 | 32.6 | 1.9 | 19.2 | -0.3 | 0.9 |
| medium-expert | halfcheetah | $51.1\pm0.9^*$ | 41.6 | 53.3 | 63.3 | 62.4 | 1.8 | 53.4 | 44.2 | 41.9 |
| medium-expert | hopper | $103.7\pm1.2^*$ | 78.5 | 108.7 | 23.7 | 111.0 | 1.6 | 96.3 | 1.9 | 0.8 |
| medium-expert | walker2d | $81.5\pm8.5^*$ | 68.0 | 95.6 | 44.6 | 98.7 | -0.1 | 40.1 | 76.9 | 81.6 |
| Average Scores | | $677.7^*$ | 575.1 | 656.5 | 391.2 | 605.0 | 3.7 | 434.2 | 328.1 | 332.0 |

Table 1. Results on D4RL. We report the averaged normalized scores of 3 different random seeds with one standard errors. We highlight the best averaged scores by a blue box . Also, we use * to indicate the tasks where our method has a solid improvement over MOReL-SAC.

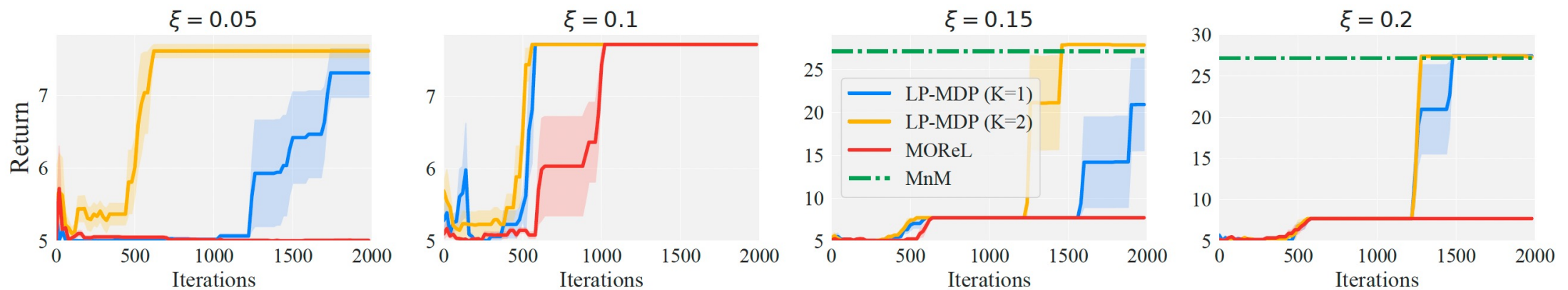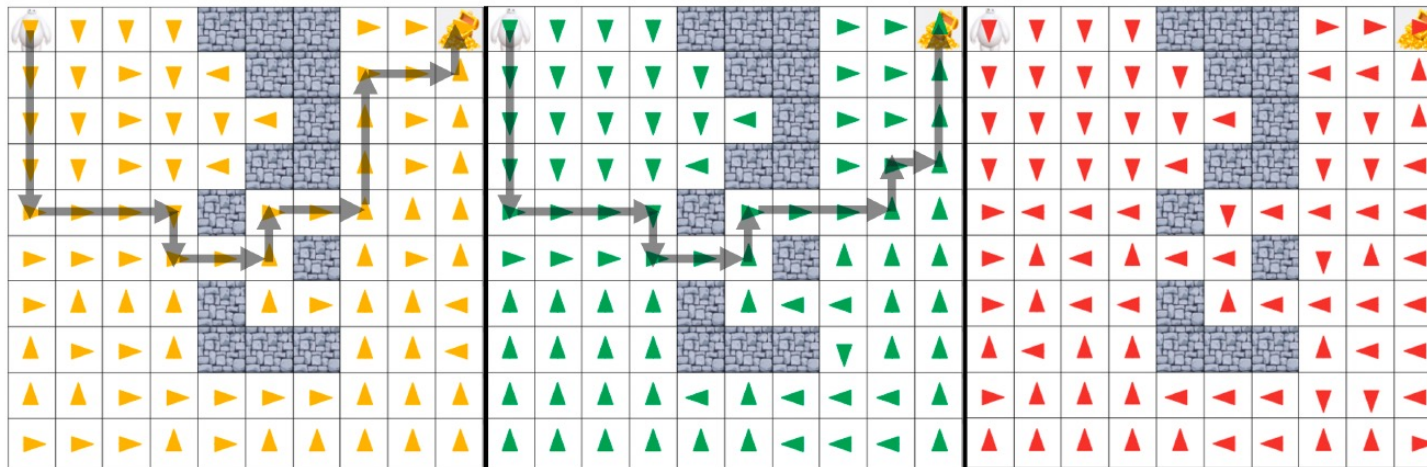# Results: The Effect of Lookahead Horizon *K*



*Figure 5.* Results on the gridworld task. The dashed green curve is the performance of the expert policy obtained with MnM (Eysenbach et al., 2022). The shaded region is the one standard error of three trials with different random seeds.

# Policy Behaviors



*Figure 6.* A visualization of the policies. The yellow policy (left) is the one learnt under LP-MDP with $k = 2$ and $\xi = 0.15$. In the middle, the green policy is the expert policy learnt using MnM (Eysenbach et al., 2022). Lastly, the red policy corresponds to the one learnt using MOReL. The grey trajectories are possible paths from the starting grid to the treasure by following the learned policies.
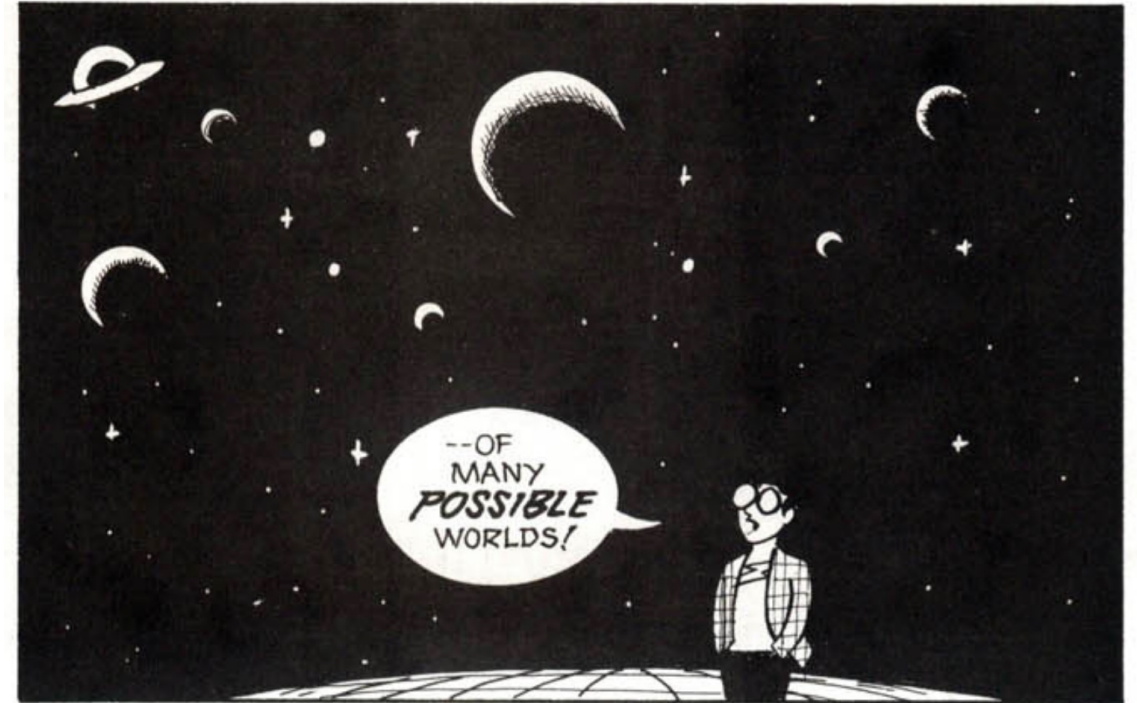
# Summary

**Key Takeaways**

- Don't be pessimistic too early; be <u>far-sighted</u>!

- Classical pessimism principle can result in rather sub-optimal policy.

- <u>Model lookahead</u> helps modulate pessimism, giving better guarantee.

**Future directions**

- Offline-to-Online RL.

- RL as sequence modeling problem (*i.e.*, return/goal conditioned imitation learning).