

Optimistic Active Exploration of Dynamical Systems (Sukhija et al., NeurIPS 2023)

Chung-En Tsai Ben Bullinger

Paper Presentation,

Foundations of Reinforcement Learning (263-5255-00L), Spring 2025



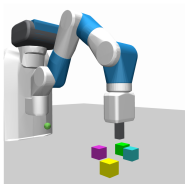
Outline

1. Introduction
2. Problem Formulation
3. Related Work and Contributions
4. Algorithm
5. Theoretical Guarantees
6. Empirical Evaluation
7. Conclusions

Outline

1. Introduction
2. Problem Formulation
3. Related Work and Contributions
4. Algorithm
5. Theoretical Guarantees
6. Empirical Evaluation
7. Conclusions

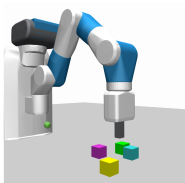
Motivation: Beyond Task-Specific Models



Traditional Model-Based RL:

- Learn a model for one specific reward
- Biased exploration of the state-action space
- May not generalize to new tasks

Motivation: Beyond Task-Specific Models



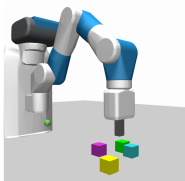
Traditional Model-Based RL:

- Learn a model for one specific reward
- Biased exploration of the state-action space
- May not generalize to new tasks

A More Ambitious Goal:

- Learn the underlying **dynamical system**
- Global exploration
- **Zero-shot** transfer to new tasks

Motivation: Beyond Task-Specific Models



Traditional Model-Based RL:

- Learn a model for one specific reward
- Biased exploration of the state-action space
- May not generalize to new tasks

A More Ambitious Goal:

- Learn the underlying **dynamical system**
- Global exploration
- **Zero-shot** transfer to new tasks

*Can we **efficiently** learn accurate models that generalize to **any** downstream task?*

Outline

1. Introduction
- 2. Problem Formulation**
3. Related Work and Contributions
4. Algorithm
5. Theoretical Guarantees
6. Empirical Evaluation
7. Conclusions

Problem Formulation (1/2)

Given an **unknown** discrete-time dynamical system $\mathbf{f}^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S} \subseteq \mathbb{R}^d$:

$$\mathbf{s}_{t+1} = \mathbf{f}^*(\mathbf{s}_t, \mathbf{a}_t) + \mathbf{w}_t, \quad \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

Goal: Design algorithms for learning $\hat{\mathbf{f}} \approx \mathbf{f}^*$ in episodic setting that is good for any reward functions.

Problem Formulation (1/2)

Given an **unknown** discrete-time dynamical system $f^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S} \subseteq \mathbb{R}^d$:

$$s_{t+1} = f^*(s_t, a_t) + w_t, \quad w_t \sim \mathcal{N}(0, I).$$

Goal: Design algorithms for learning $\hat{f} \approx f^*$ in episodic setting that is good for any reward functions.

$$\mathcal{D} = \left\{ \begin{array}{c} \pi_1 \\ \text{---} \\ \pi_n \end{array} \right\}$$

Problem Formulation (2/2)

Our problem can be divided into two parts:

- **Exploration (our focus!)**: Collecting data by executing policies.
- **Estimation**: Building model from data.

Problem Formulation (2/2)

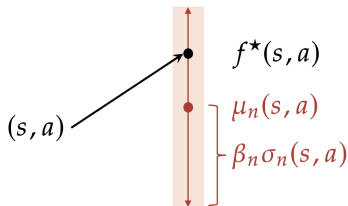
Our problem can be divided into two parts:

- **Exploration (our focus!)**: Collecting data by executing policies.
- **Estimation**: Building model from data.

Well-Calibration Assumption

Given n data, we can construct confidence intervals:

$$|f_i^*(s, a) - \mu_{n,i}(s, a)| \leq \beta_n \sigma_{n,i}(s, a) \quad \forall i \in [d], \text{ w.h.p.,}$$



Outline

1. Introduction
2. Problem Formulation
- 3. Related Work and Contributions**
4. Algorithm
5. Theoretical Guarantees
6. Empirical Evaluation
7. Conclusions

Related Work: Active Learning

	Bandit Optimization	Our Problem
Goal	$\max_x f^*(x)$	Find $\hat{f} \approx f^*$
f^*	Objective	Dynamics
State space	No	Continuous
Planning	No	Yes

N. Srinivas et al. "Gaussian process optimization in the bandit setting: No regret and experimental design." ICML 2010. (The 1st presentation today!)

Related Work: Active Learning

	Bandit Optimization	Our Problem
Goal	$\max_x f^*(x)$	Find $\hat{f} \approx f^*$
f^*	Objective	Dynamics
State space	No	Continuous
Planning	No	Yes
Technique	Optimism in the face of uncertainty	

N. Srinivas et al. “Gaussian process optimization in the bandit setting: No regret and experimental design.” ICML 2010. (The 1st presentation today!)

Related Work: Maximum Entropy RL

What about collecting data by maximum entropy or entropy-regularized RL?

$$\max_{\pi} H(d_{\mu}^{\pi}), \quad \max_{\pi} E \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) - \beta \log \pi(a_t | s_t)) \mid s_0 \sim \mu, \pi \right].$$

- **No guarantee** for arbitrary reward functions in downstream tasks.

E. Hazan et al. “Provably efficient maximum entropy exploration.” ICML 2019. (Previous presentation.)
B. Eysenbach and S. Levine. “Maximum entropy RL (provably) solves some robust RL problems.” ICLR 2022.

Related Work: Reward-Free RL

- Study this problem in tabular settings or impose structural assumptions on values function and transitions.
- No practical algorithms.

C. Jin et al. “Reward-free exploration for reinforcement learning.” ICML 2020.

R. Wang et al. “On reward-free reinforcement learning with linear function approximation.” NeurIPS 2020.

J. Chen et al. “On the statistical efficiency of reward-free exploration in non-linear RL.” NeurIPS 2022.

Contributions

- Propose OpAX, a practical algorithm for learning non-linear dynamics with **continuous** state-action spaces and without structural assumptions.
- In this setting, OpAX is the **first** algorithm to have a theoretical guarantee for a rich family of nonlinear dynamics.

Outline

1. Introduction
2. Problem Formulation
3. Related Work and Contributions
- 4. Algorithm**
5. Theoretical Guarantees
6. Empirical Evaluation
7. Conclusions

OpAX: Overview

*Learn dynamics by **maximizing information gain**, not rewards*

OpAX: Overview

Learn dynamics by *maximizing information gain*, not rewards

	Traditional Model-Based RL	OpAX
Objective	$\max_{\pi} E[R(\pi)]$	$\max_{\pi} I(f^*; \mathcal{D})$
Exploration	Reward-driven	Uncertainty-driven
Convergence to f^*	No guarantees	$\sigma_n(s, a) \rightarrow 0$ (GP)
Generalization	No guarantees	Zero-shot (GP)

OpAX (1/4): Information-Theoretic Objective

Ideal Goal: Maximize mutual information between dynamics f^* and collected data

OpAX (1/4): Information-Theoretic Objective

Ideal Goal: Maximize mutual information between dynamics f^* and collected data

Greedy Episodic Approximation: At episode n , choose policy maximizing expected information gain:

$$\pi_n^* = \arg \max_{\pi \in \Pi} \mathbb{E}_{\tau^\pi} [\mathcal{I}(f_{\tau^\pi}^*; \tau^\pi \mid \mathcal{D}_{1:n-1})]$$

where:

- $f_{\tau^\pi}^* = (f^*(s_0, a_0), \dots, f^*(s_{T-1}, a_{T-1}))$: true dynamics along trajectory
- $\mathcal{D}_{1:n-1}$: data from previous episodes

OpAX (1/4): Information-Theoretic Objective

Ideal Goal: Maximize mutual information between dynamics f^* and collected data

Greedy Episodic Approximation: At episode n , choose policy maximizing expected information gain:

$$\pi_n^* = \arg \max_{\pi \in \Pi} \mathbb{E}_{\tau^\pi} [\mathcal{I}(f_{\tau^\pi}^*; \tau^\pi \mid \mathcal{D}_{1:n-1})]$$

where:

- $f_{\tau^\pi}^* = (f^*(s_0, a_0), \dots, f^*(s_{T-1}, a_{T-1}))$: true dynamics along trajectory
- $\mathcal{D}_{1:n-1}$: data from previous episodes

Challenge

This objective is **intractable** in general. Need tractable approximation.

OpAX (2/4): Tractable Upper Bound

Lemma (Information Gain Upper Bound)

For Gaussian noise $w_t \sim \mathcal{N}(0, I)$ and epistemic uncertainty σ_{n-1} :

$$I(\mathbf{f}_{\tau}^*; \tau^\pi \mid \mathcal{D}_{1:n-1}) \leq \frac{1}{2} \sum_{t=0}^{T-1} \sum_{j=1}^d \log(1 + \sigma_{n-1,j}^2(\mathbf{s}_t, \mathbf{a}_t))$$

OpAX (2/4): Tractable Upper Bound

Lemma (Information Gain Upper Bound)

For Gaussian noise $w_t \sim \mathcal{N}(0, I)$ and epistemic uncertainty σ_{n-1} :

$$I(f_{\tau}^*; \tau^\pi \mid \mathcal{D}_{1:n-1}) \leq \frac{1}{2} \sum_{t=0}^{T-1} \sum_{j=1}^d \log(1 + \sigma_{n-1,j}^2(s_t, a_t))$$

Tractable Exploration Objective: Maximize the upper bound

$$\pi_n^* = \arg \max_{\pi \in \Pi} \mathbb{E}_{\tau^\pi} \left[\sum_{t=0}^{T-1} \sum_{j=1}^d \log(1 + \sigma_{n-1,j}^2(s_t, a_t)) \right]$$

subject to: $s_{t+1} = f^*(s_t, a_t) + w_t$.

OpAX (2/4): Tractable Upper Bound

Lemma (Information Gain Upper Bound)

For Gaussian noise $w_t \sim \mathcal{N}(0, I)$ and epistemic uncertainty σ_{n-1} :

$$I(f_{\tau}^*; \tau^\pi \mid \mathcal{D}_{1:n-1}) \leq \frac{1}{2} \sum_{t=0}^{T-1} \sum_{j=1}^d \log(1 + \sigma_{n-1,j}^2(s_t, a_t))$$

Tractable Exploration Objective: Maximize the upper bound

$$\pi_n^* = \arg \max_{\pi \in \Pi} \mathbb{E}_{\tau^\pi} \left[\sum_{t=0}^{T-1} \sum_{j=1}^d \log(1 + \sigma_{n-1,j}^2(s_t, a_t)) \right]$$

subject to: $s_{t+1} = f^*(s_t, a_t) + w_t$.

Collect data where model is most uncertain \Rightarrow maximal information gain

OpAX (3/4): Hallucination Policy

Problem: The planning objective requires knowing f^* , but f^* is **unknown**!

¹K. Chua et al. “Deep reinforcement learning in a handful of trials using probabilistic dynamics models.” NeurIPS 2018.

²M. Simchowitz et al. “Naive exploration is optimal for online lqr.” ICML 2020.

OpAX (3/4): Hallucination Policy

Problem: The planning objective requires knowing f^* , but f^* is **unknown**!

Naive Approach: Use mean estimate μ_{n-1} instead of f^*

¹K. Chua et al. “Deep reinforcement learning in a handful of trials using probabilistic dynamics models.” NeurIPS 2018.

²M. Simchowitz et al. “Naive exploration is optimal for online lqr.” ICML 2020.

OpAX (3/4): Hallucination Policy

Problem: The planning objective requires knowing f^* , but f^* is **unknown**!

Naive Approach: Use mean estimate μ_{n-1} instead of f^*

- ✗ Susceptible to model biases¹
- ✗ Provably optimal only for linear systems²

¹K. Chua et al. “Deep reinforcement learning in a handful of trials using probabilistic dynamics models.” NeurIPS 2018.

²M. Simchowitz et al. “Naive exploration is optimal for online lqr.” ICML 2020.

OpAX (3/4): Hallucination Policy

Problem: The planning objective requires knowing f^* , but f^* is **unknown**!

Naive Approach: Use mean estimate μ_{n-1} instead of f^*

- ✗ Susceptible to model biases¹
- ✗ Provably optimal only for linear systems²

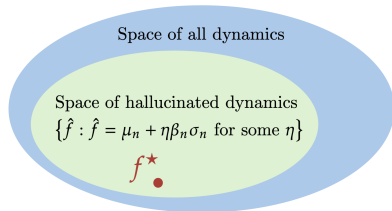
OpAX Solution:

Optimism in the face of uncertainty

Use a **hallucination policy** η to select optimistic dynamics

$$\hat{f}_j(s, a) = \mu_{n-1,j}(s, a) + \beta_{n-1} \sigma_{n-1,j}(s, a) \cdot \eta_j(s)$$

where $\eta : \mathcal{S} \rightarrow [-1, 1]^d$.



¹K. Chua et al. "Deep reinforcement learning in a handful of trials using probabilistic dynamics models." NeurIPS 2018.

²M. Simchowitz et al. "Naive exploration is optimal for online lqr." ICML 2020.

OpAX (4/4): Optimistic Planning

OpAX Optimistic Objective

$$\pi_n, \eta_n = \arg \max_{\pi \in \Pi, \eta \in \Xi} \mathbb{E}_{\tau^{\pi, \eta}} \left[\sum_{t=0}^{T-1} \sum_{j=1}^d \log (1 + \sigma_{n-1,j}^2(\hat{s}_t, \pi(\hat{s}_t))) \right]$$

subject to:

$$\hat{s}_{t+1} = \mu_{n-1}(\hat{s}_t, \pi(\hat{s}_t)) + \beta_{n-1} \sigma_{n-1}(\hat{s}_t, \pi(\hat{s}_t)) \cdot \eta(\hat{s}_t) + w_t$$

where $\Xi = \{\eta : \mathcal{S} \rightarrow [-1, 1]^d\}$ is the space of **hallucination policies**.

- $\eta(\hat{s}_t)$ picks the most optimistic transition within confidence bounds

OpAX (4/4): Optimistic Planning

OpAX Optimistic Objective

$$\pi_n, \eta_n = \arg \max_{\pi \in \Pi, \eta \in \Xi} \mathbb{E}_{\tau^{\pi, \eta}} \left[\sum_{t=0}^{T-1} \sum_{j=1}^d \log (1 + \sigma_{n-1, j}^2(\hat{s}_t, \pi(\hat{s}_t))) \right]$$

subject to:

$$\hat{s}_{t+1} = \mu_{n-1}(\hat{s}_t, \pi(\hat{s}_t)) + \beta_{n-1} \sigma_{n-1}(\hat{s}_t, \pi(\hat{s}_t)) \cdot \eta(\hat{s}_t) + w_t$$

where $\Xi = \{\eta : \mathcal{S} \rightarrow [-1, 1]^d\}$ is the space of **hallucination policies**.

- $\eta(\hat{s}_t)$ picks the most optimistic transition within confidence bounds

Implementation

Solve as optimal control with extended action space (π, η) .

Algorithm

OpAX: OPTIMISTIC ACTIVE EXPLORATION

Init: Statistical model $(\mu_0, \sigma_0, \beta_0)$

for episode $n = 1, \dots, N$ **do**

$$\pi_n = \arg \max_{\pi \in \Pi} \max_{\eta \in \Xi} \mathbb{E} \left[\sum_{t=0}^{T-1} \sum_{j=1}^d \log \left(1 + \sigma_{n-1,j}^2(s_t, a_t) \right) \right] \quad \blacktriangleright \text{Prepare policy}$$

$\mathcal{D}_n \leftarrow \text{ROLLOUT}(\pi_n)$ \blacktriangleright Measure

Update $(\mu_n, \sigma_n, \beta_n) \leftarrow \mathcal{D}_{1:n}$ \blacktriangleright Update model

end for

Algorithm

OpAX: OPTIMISTIC ACTIVE EXPLORATION

Init: Statistical model $(\mu_0, \sigma_0, \beta_0)$

for episode $n = 1, \dots, N$ **do**

$$\pi_n = \arg \max_{\pi \in \Pi} \max_{\eta \in \Xi} \mathbb{E} \left[\sum_{t=0}^{T-1} \sum_{j=1}^d \log \left(1 + \sigma_{n-1,j}^2(s_t, a_t) \right) \right] \quad \blacktriangleright \text{Prepare policy}$$

$\mathcal{D}_n \leftarrow \text{ROLLOUT}(\pi_n)$ \blacktriangleright Measure

Update $(\mu_n, \sigma_n, \beta_n) \leftarrow \mathcal{D}_{1:n}$ \blacktriangleright Update model

end for

Key Components:

- **Information reward:** $\log(1 + \sigma^2)$
- **Optimistic planning:** Max over η

Algorithm

OpAX: OPTIMISTIC ACTIVE EXPLORATION

Init: Statistical model $(\mu_0, \sigma_0, \beta_0)$

for episode $n = 1, \dots, N$ **do**

$$\pi_n = \arg \max_{\pi \in \Pi} \max_{\eta \in \Xi} \mathbb{E} \left[\sum_{t=0}^{T-1} \sum_{j=1}^d \log \left(1 + \sigma_{n-1,j}^2(s_t, a_t) \right) \right] \quad \blacktriangleright \text{Prepare policy}$$

$\mathcal{D}_n \leftarrow \text{ROLLOUT}(\pi_n)$ \blacktriangleright Measure

Update $(\mu_n, \sigma_n, \beta_n) \leftarrow \mathcal{D}_{1:n}$ \blacktriangleright Update model

end for

Key Components:

- **Information reward:** $\log(1 + \sigma^2)$
- **Optimistic planning:** Max over η

Implementation:

- MPC for short horizons
- SAC for long horizons

Outline

1. Introduction
2. Problem Formulation
3. Related Work and Contributions
4. Algorithm
- 5. Theoretical Guarantees**
6. Empirical Evaluation
7. Conclusions

Assumptions

Regularity of the Dynamics

f^* lies in a RKHS with kernel k and $\|f^*\| \leq 1$.

N. Srinivas et al. “Gaussian process optimization in the bandit setting: No regret and experimental design.” ICML 2010. (The 1st presentation today!)

Assumptions

Regularity of the Dynamics

f^* lies in a RKHS with kernel k and $\|f^*\| \leq 1$.

Lemma

$GP(0, k)$ is well-calibrated.

Recall

$$\gamma_n(k) = \max_{\mathcal{D}_1, \dots, \mathcal{D}_n} \log \det(I + K_n).$$

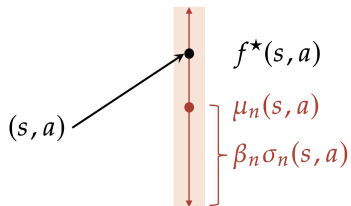
N. Srinivas et al. “Gaussian process optimization in the bandit setting: No regret and experimental design.” ICML 2010. (The 1st presentation today!)

Main Guarantee

Theorem (Informal)

With high probability,

$$\max_{\pi} \mathbb{E}_{\pi} \left[\underbrace{\max_t \|\sigma_n(s_t, a_t)\|_2^2}_{\text{maximum uncertainty along the trajectory}} \right] \leq \text{poly}(T) \cdot \frac{\gamma_n(k)}{\sqrt{n}}, \quad \forall n \in \mathbb{N}.$$

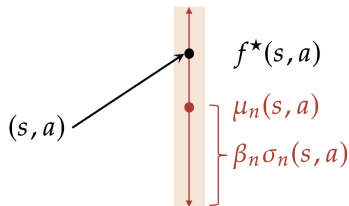


Main Guarantee

Theorem (Informal)

With high probability,

$$\max_{\pi} \mathbb{E}_{\pi} \left[\underbrace{\max_t \|\sigma_n(s_t, a_t)\|_2^2}_{\text{maximum uncertainty along the trajectory}} \right] \leq \text{poly}(T) \cdot \frac{\gamma_n(k)}{\sqrt{n}}, \quad \forall n \in \mathbb{N}.$$



For linear and RBF kernels, $\gamma_n \approx \text{polylog}(nT)$ and $\lim_{n \rightarrow \infty} \sigma_n(s, a) \rightarrow 0$.

Zero-Shot Performance

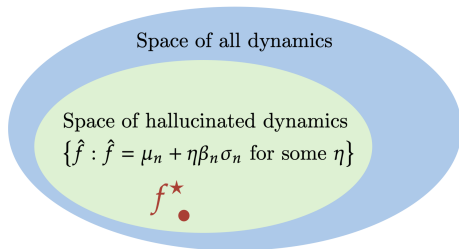
Theorem (Informal)

Consider an RL problem with a known bounded reward:

$$\arg \max_{\pi} R(\pi), \quad s_{t+1} = f^*(s_t, a_t) + w_t.$$

If $n \gtrsim \text{poly}(d, T, \varepsilon^{-1}, \gamma_n(k))$, then $\hat{\pi}$ is ε -optimal, where

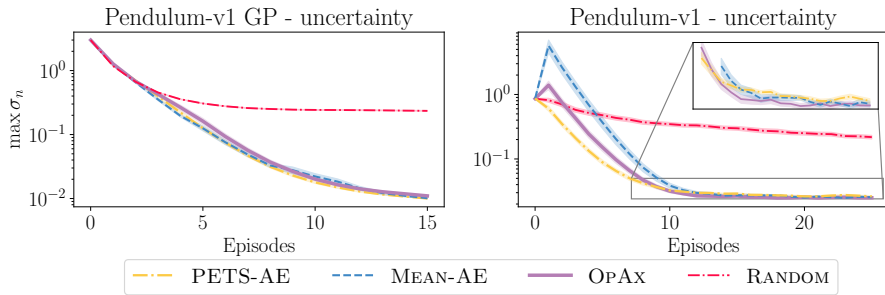
$$\hat{\pi} \in \arg \max_{\pi} \min_{\eta} R(\pi, \eta), \quad s_{t+1} = \mu_n(s_t, a_t) + \eta(s_t) \beta_n \sigma_n(s_t, a_t) + w_t.$$



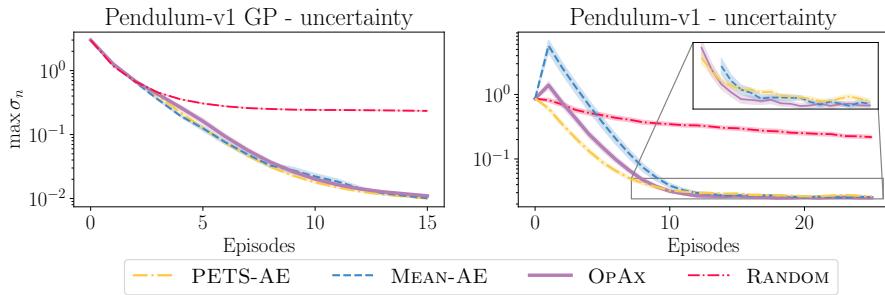
Outline

1. Introduction
2. Problem Formulation
3. Related Work and Contributions
4. Algorithm
5. Theoretical Guarantees
- 6. Empirical Evaluation**
7. Conclusions

Experiments: Epistemic Uncertainty Reduction



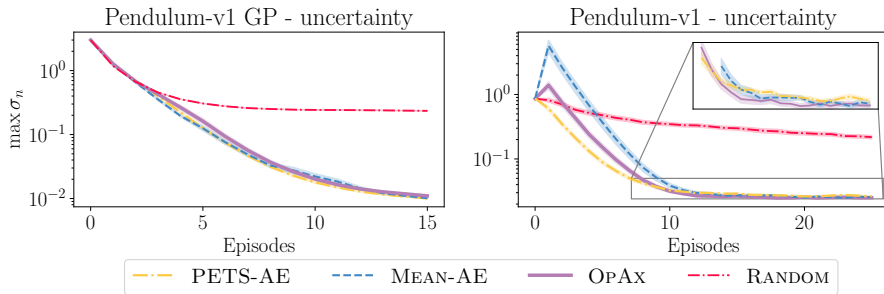
Experiments: Epistemic Uncertainty Reduction



Key Results:

- **Efficient exploration:** Active methods reduces uncertainty faster than random baseline
- **Optimism:** OpAX provides edge over mean planning

Experiments: Epistemic Uncertainty Reduction



Key Results:

- **Efficient exploration:** Active methods reduces uncertainty faster than random baseline
- **Optimism:** OpAX provides edge over mean planning

Experiments validate theoretical results

Experiments: Downstream Tasks

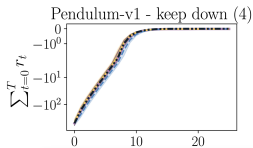
- OpAX performs **on par** with H-UCRL¹ on trained tasks
- OpAX **outperforms** H-UCRL on some unseen tasks

¹S. Curi et al. “Efficient model-based reinforcement learning through optimistic policy search and planning.” NeurIPS 2020.

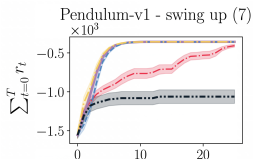
Experiments: Downstream Tasks

- OpAX performs **on par** with H-UCRL¹ on trained tasks
- OpAX **outperforms** H-UCRL on some unseen tasks

Evaluation (■ : H-UCRL, ■ : OpAX, ■ : Random):



(a) Trained task



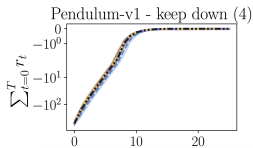
(b) Zero-shot task

¹S. Curi et al. "Efficient model-based reinforcement learning through optimistic policy search and planning." NeurIPS 2020.

Experiments: Downstream Tasks

- OpAX performs **on par** with H-UCRL¹ on trained tasks
- OpAX **outperforms** H-UCRL on some unseen tasks

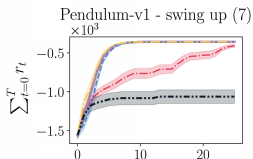
Evaluation (■ : H-UCRL, ■ : OpAX, ■ : Random): **Zero-shot swing up task:**



(a) Trained task

(play)

(c) H-UCRL



(b) Zero-shot task

(play)

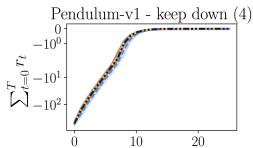
(d) OpAx

¹S. Curi et al. "Efficient model-based reinforcement learning through optimistic policy search and planning." NeurIPS 2020.

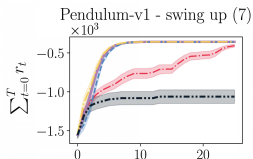
Experiments: Downstream Tasks

- OpAX performs **on par** with H-UCRL¹ on trained tasks
- OpAX **outperforms** H-UCRL on some unseen tasks

Evaluation (■ : H-UCRL, ■ : OpAX, ■ : Random): **Zero-shot swing up task:**



(a) Trained task



(b) Zero-shot task

(play)

(c) H-UCRL

More experiments in the paper.

(play)

(d) OpAx

¹S. Curi et al. "Efficient model-based reinforcement learning through optimistic policy search and planning." NeurIPS 2020.

Outline

1. Introduction
2. Problem Formulation
3. Related Work and Contributions
4. Algorithm
5. Theoretical Guarantees
6. Empirical Evaluation
- 7. Conclusions**

Conclusions

	Traditional Model-Based RL	OpAX
Objective	$\max_{\pi} E[R(\pi)]$	$\max_{\pi} I(f^*; \mathcal{D})$
Exploration	Reward-driven	Uncertainty-driven
Convergence to f^*	No guarantees	$\sigma_n(s, a) \rightarrow 0$ (GP)
Generalization	No guarantees	Zero-shot (GP)

OpAX

Combines **active learning** with RL and achieves **theoretical guarantees** for learning many continuous nonlinear systems.



Open Challenges and Future Work

Key Limitations:

- **Computational cost:** Solving (π_n, η_n) is expensive
- **Model requirements:** Needs well-calibrated uncertainty
- **Known noise:** Guarantees only hold for (sub-)Gaussian noise.

Promising Directions:

- **Robust models:** Handle model misspecification
- **Beyond MDPs:** POMDPs, continuous time
- **Practical deployment:** Real-world robotics