

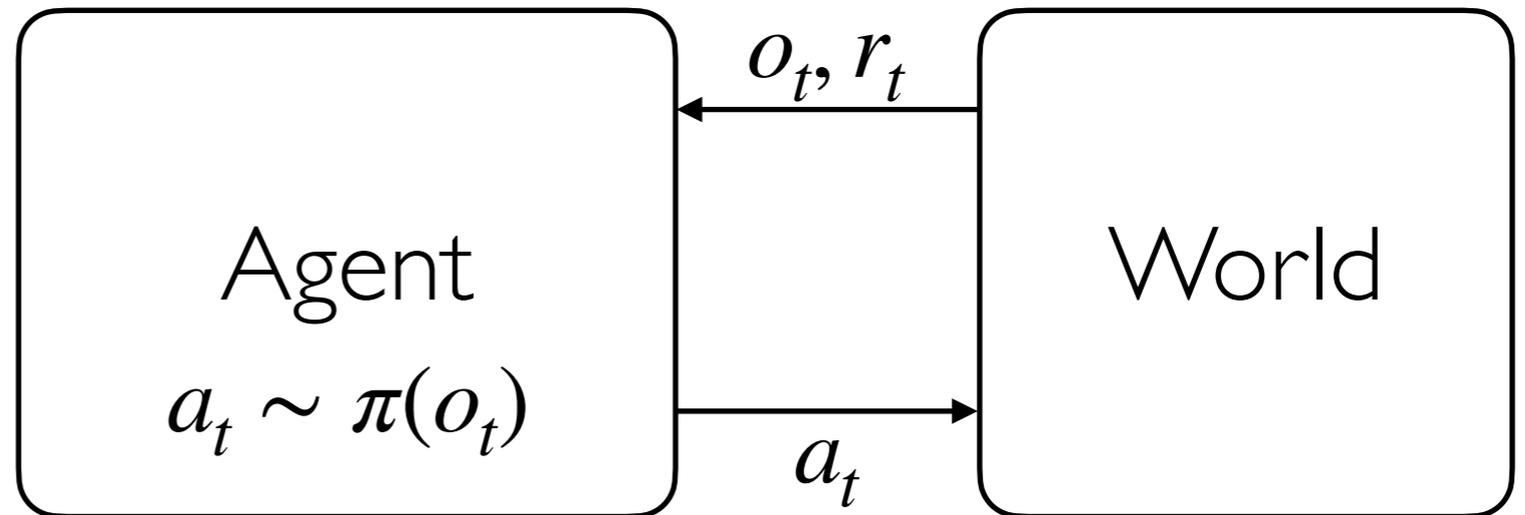
Learning Models for RL

Saurabh Gupta

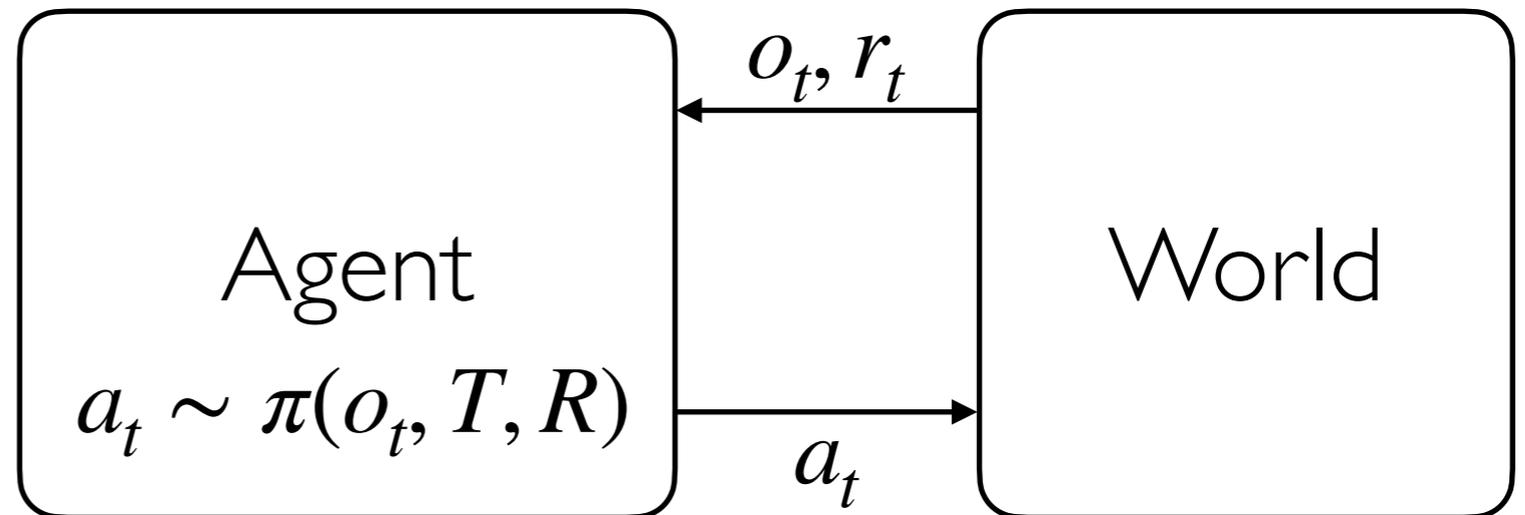
Solving MDPs

Policy: $a_t \sim \pi(o_t)$

Most General Case



More Specific Case



Fully Observed System

$$o_t = s_t$$

Known Transition Function

$$s_{t+1} \sim T(s_t, a_t)$$

Known Reward Function

$$R(s_{t+1}, s_t, a_t)$$

Solving a RL Problem

Better Reward Signals

Better Optimization

Convert into a Supervised Training Problem

Solve a Related but Supervision-rich Problem

Build Models and Plan with Them

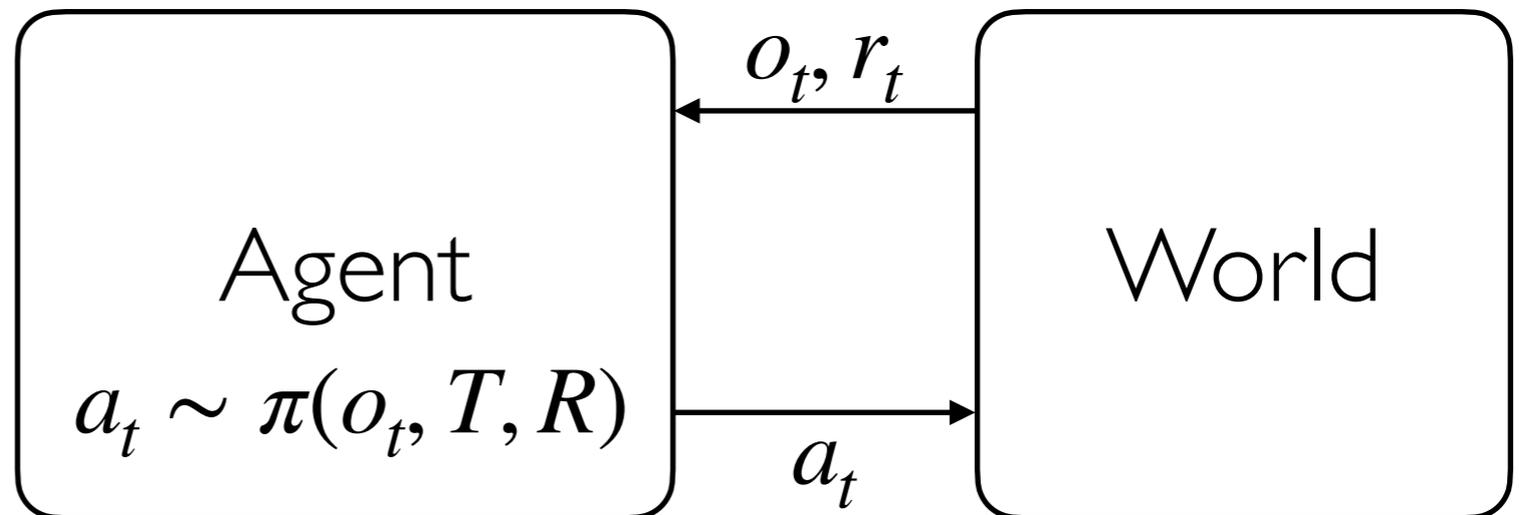
Model-free RL
with sparse
rewards

Known reward,
known model.
Model-based RL



Build Models and Plan with Them

Relax Assumptions For Model-Based RL



Fully Observed System

$$o_t = s_t$$

Known Transition Function

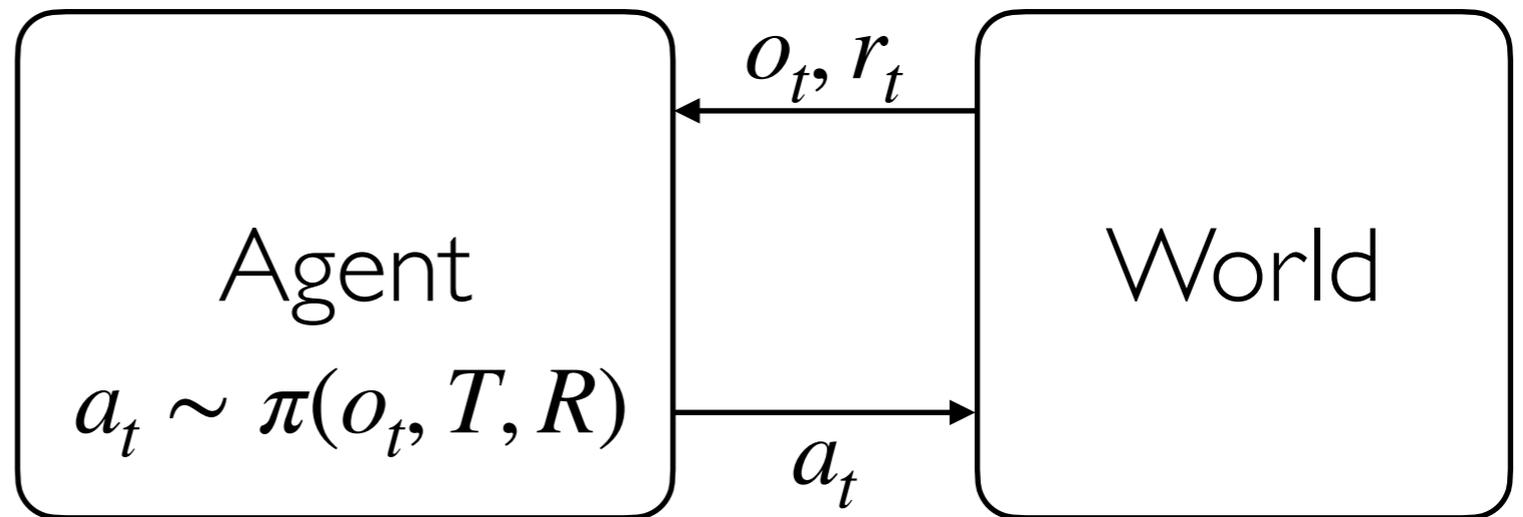
$$s_{t+1} \sim T(s_t, a_t)$$

Known Reward Function

$$R(s_{t+1}, s_t, a_t)$$

Build Models and Plan with Them

Relax Assumptions For Model-Based RL



Fully Observed System $o_t = s_t$

~~Known Transition Function $s_{t+1} \sim T(s_t, a_t)$~~

Learn the Transition Function $s_{t+1} \sim \hat{T}(s_t, a_t)$

Known Reward Function $R(s_{t+1}, s_t, a_t)$

Build Models and Plan with Them

Relax Assumptions For Model-Based RL

Learn the Transition Function $s_{t+1} \sim \hat{T}(s_t, a_t)$

1. Collect data for learning the model:
 1. Initialize system and execute random actions
 2. Collect dataset of triples (s_{t+1}, s_t, a_t)
2. Fit a *function* \hat{T} that predicts s_{t+1} given s_t and a_t
 - Gaussian Process
 - Neural Network
3. Plan using function \hat{T}

Build Models and Plan with Them

Relax Assumptions For Model-Based RL

Learn the Transition Function $s_{t+1} \sim \hat{T}(s_t, a_t)$

Optimizer

Random Shooting

First Order

Function \hat{T}

Gaussian Process

PILCO

Neural Network

PETS

Deep PILCO

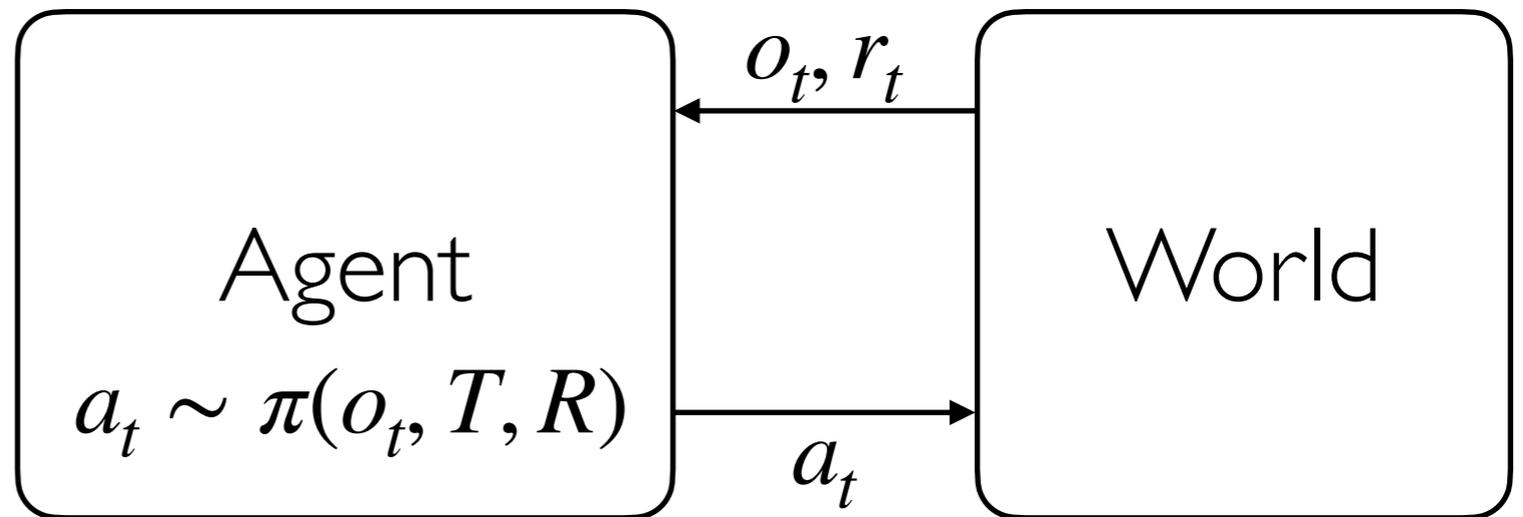
[PILCO] M. Deisenroth et al. **PILCO: A Model-based and Data-Efficient Approach to Policy Search**. ICML 2011

[PETS] K. Chua et al. **Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models**. NIPS 2018

[Deep PILCO] Y. Gal et al. **Improving PILCO with Bayesian neural network dynamics models**. ICMLW 2016.

Build Models and Plan with Them

Relax Assumptions For Model-Based RL



Fully Observed System $o_t = s_t$

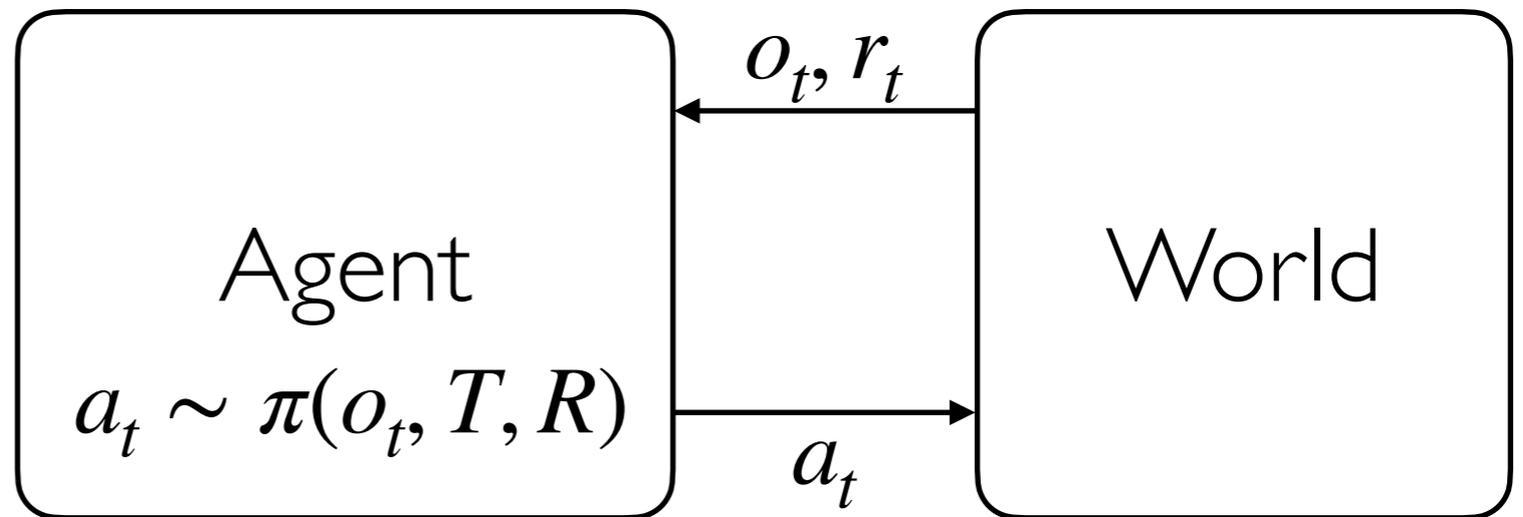
~~Known Transition Function $s_{t+1} \sim T(s_t, a_t)$~~

Learn the Transition Function $s_{t+1} \sim \hat{T}(s_t, a_t)$

Known Reward Function $R(s_{t+1}, s_t, a_t)$

Build Models and Plan with Them

Relax Assumptions For Model-Based RL



~~Fully Observed System $o_t = s_t$~~

~~Known Transition Function $s_{t+1} \sim T(s_t, a_t)$~~

~~Learn the Transition Function $s_{t+1} \sim \hat{T}(s_t, a_t)$~~

Learn Transition Function (using observations) $o_{t+1} \sim \hat{T}(o_t, a_t)$

~~Known Reward Function $R(s_{t+1}, s_t, a_t)$~~

Reward Function (using observations) $R(o_{t+1}, o_t, a_t)$

Build Models and Plan with Them

Relax Assumptions For Model-Based RL

Learn Transition Function (using observations) $o_{t+1} \sim \hat{T}(o_t, a_t)$

1. Collect data for learning the model:
 1. Initialize system and execute random actions
 2. Collect dataset of triples (o_{t+1}, o_t, a_t)
2. Fit a function \hat{T} that predicts o_{t+1} given o_t and a_t
 - ~~Gaussian Process~~
 - Neural Network
3. Plan using function \hat{T}

Build Models and Plan with Them

Relax Assumptions For Model-Based RL

Learn the Transition Function $s_{t+1} \sim \hat{T}(s_t, a_t)$

Optimizer

Random Shooting

First Order

Function \hat{T}

Gaussian Process

PILCO

Neural Network

PETS

Deep PILCO

[PILCO] M. Deisenroth et al. **PILCO: A Model-based and Data-Efficient Approach to Policy Search**. ICML 2011

[PETS] K. Chua et al. **Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models**. NIPS 2018

[Deep PILCO] Y. Gal et al. **Improving PILCO with Bayesian neural network dynamics models**. ICMLW 2016.

PILCO

Algorithm 1 PILCO

- 1: **init:** Sample controller parameters $\theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.
Apply random control signals and record data.
 - 2: **repeat**
 - 3: Learn probabilistic (GP) dynamics model, see Sec. 2.1, using all data.
 - 4: Model-based policy search, see Sec. 2.2–2.3.
 - 5: **repeat**
 - 6: Approximate inference for policy evaluation, see Sec. 2.2: get $J^\pi(\theta)$, Eqs. (10)–(12), (24).
 - 7: Gradient-based policy improvement, see Sec. 2.3: get $dJ^\pi(\theta)/d\theta$, Eqs. (26)–(30).
 - 8: Update parameters θ (e.g., CG or L-BFGS).
 - 9: **until** convergence; **return** θ^*
 - 10: Set $\pi^* \leftarrow \pi(\theta^*)$.
 - 11: Apply π^* to system (single trial/episode) and record data.
 - 12: **until** task learned
-

Gaussian Processes

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right)$$

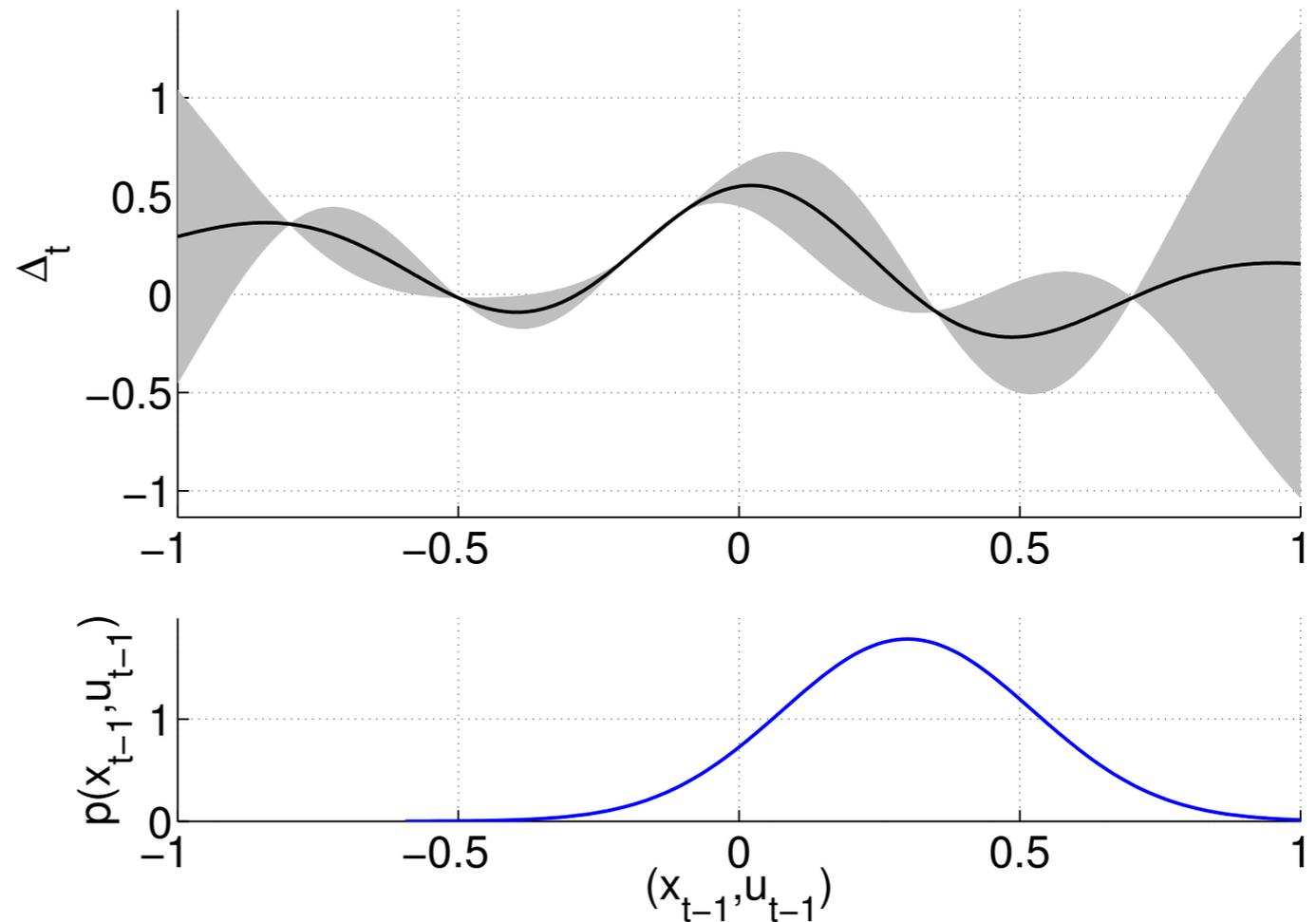
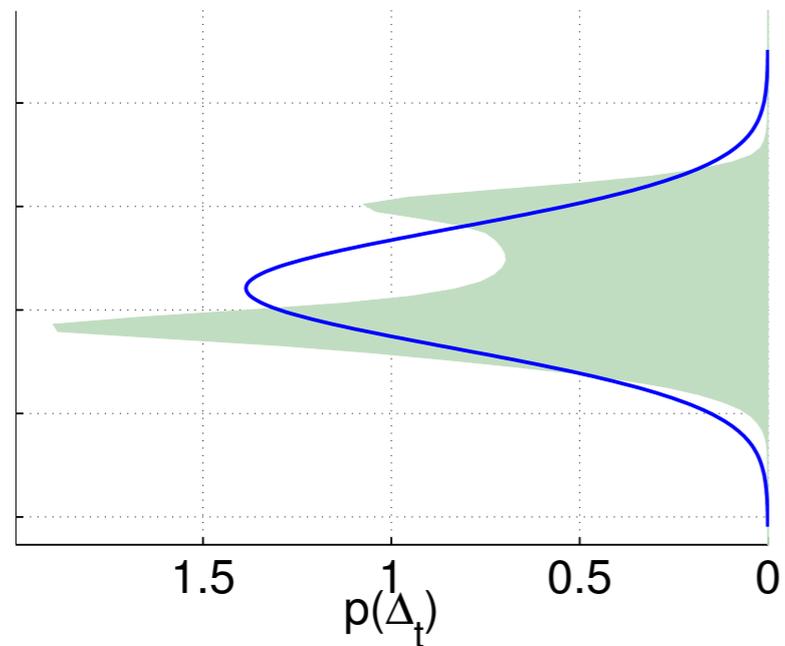
$$\mathbf{f}_* | X_*, X, \mathbf{f} \sim \mathcal{N} \left(K(X_*, X) K(X, X)^{-1} \mathbf{f}, \right. \\ \left. K(X_*, X_*) - K(X_*, X) K(X, X)^{-1} K(X, X_*) \right)$$

<http://chifeng.scripts.mit.edu/stuff/gp-demo/>

Policy Evaluation

$$\mathbb{E}_{\mathbf{x}_t} [c(\mathbf{x}_t)] = \int c(\mathbf{x}_t) \mathcal{N}(\mathbf{x}_t | \mu_t, \Sigma_t) d\mathbf{x}_t,$$

Given $p(x_0)$ and policy $\pi(\theta)$, want to compute $p(x_t)$.



Policy Improvement

$$\mathbb{E}_{\mathbf{x}_t} [c(\mathbf{x}_t)] = \int c(\mathbf{x}_t) \mathcal{N}(\mathbf{x}_t | \mu_t, \Sigma_t) d\mathbf{x}_t ,$$

Compute gradient with respect to θ
using chain rule.

PILCO

Algorithm 1 PILCO

- 1: **init:** Sample controller parameters $\theta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.
Apply random control signals and record data.
 - 2: **repeat**
 - 3: Learn probabilistic (GP) dynamics model, see Sec. 2.1, using all data.
 - 4: Model-based policy search, see Sec. 2.2–2.3.
 - 5: **repeat**
 - 6: Approximate inference for policy evaluation, see Sec. 2.2: get $J^\pi(\theta)$, Eqs. (10)–(12), (24).
 - 7: Gradient-based policy improvement, see Sec. 2.3: get $dJ^\pi(\theta)/d\theta$, Eqs. (26)–(30).
 - 8: Update parameters θ (e.g., CG or L-BFGS).
 - 9: **until** convergence; **return** θ^*
 - 10: Set $\pi^* \leftarrow \pi(\theta^*)$.
 - 11: Apply π^* to system (single trial/episode) and record data.
 - 12: **until** task learned
-

Build Models and Plan with Them

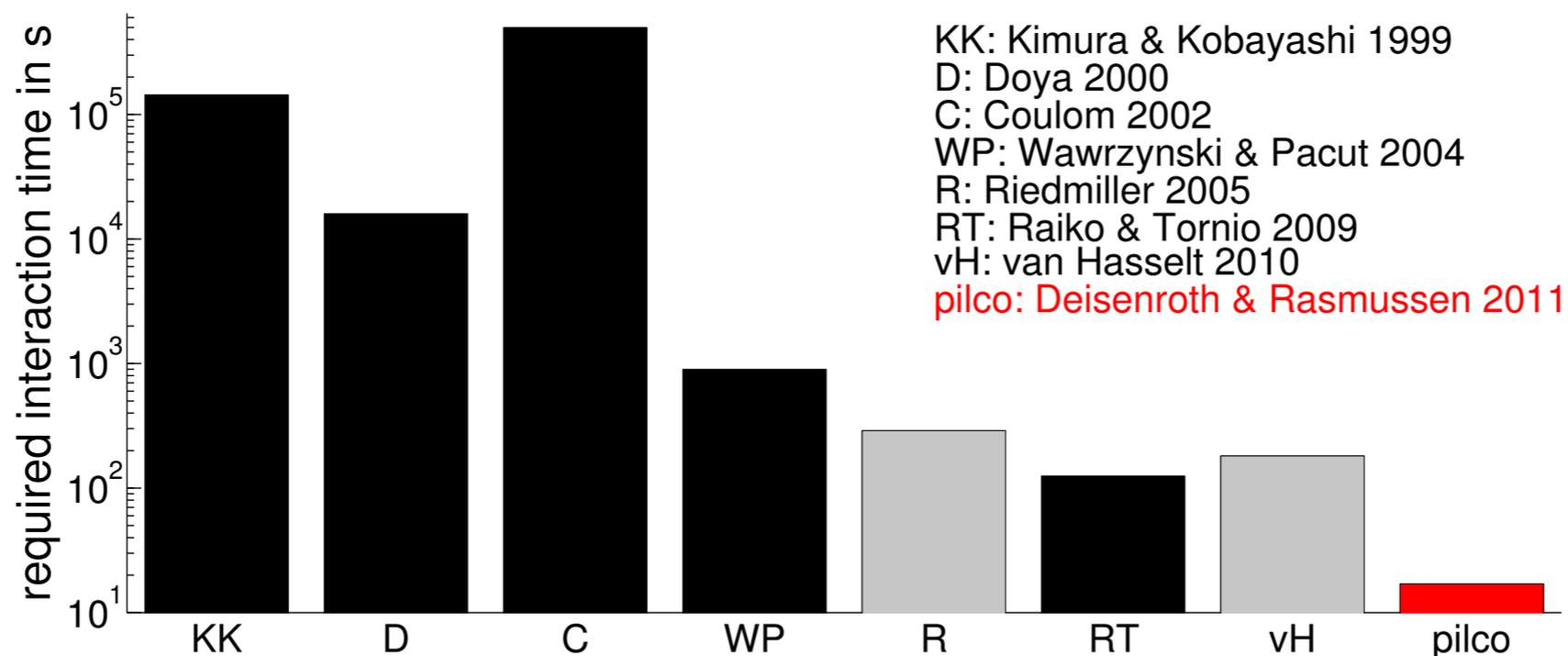
PILCO - Inverting a pendulum



PILCO

Table 1. PILCO's data efficiency scales to high dimensions.

	cart-pole	cart-double-pole	unicycle
state space	\mathbb{R}^4	\mathbb{R}^6	\mathbb{R}^{12}
# trials	≤ 10	20–30	≈ 20
experience	≈ 20 s	≈ 60 s– 90 s	≈ 20 s– 30 s
parameter space	\mathbb{R}^{305}	\mathbb{R}^{1816}	\mathbb{R}^{28}

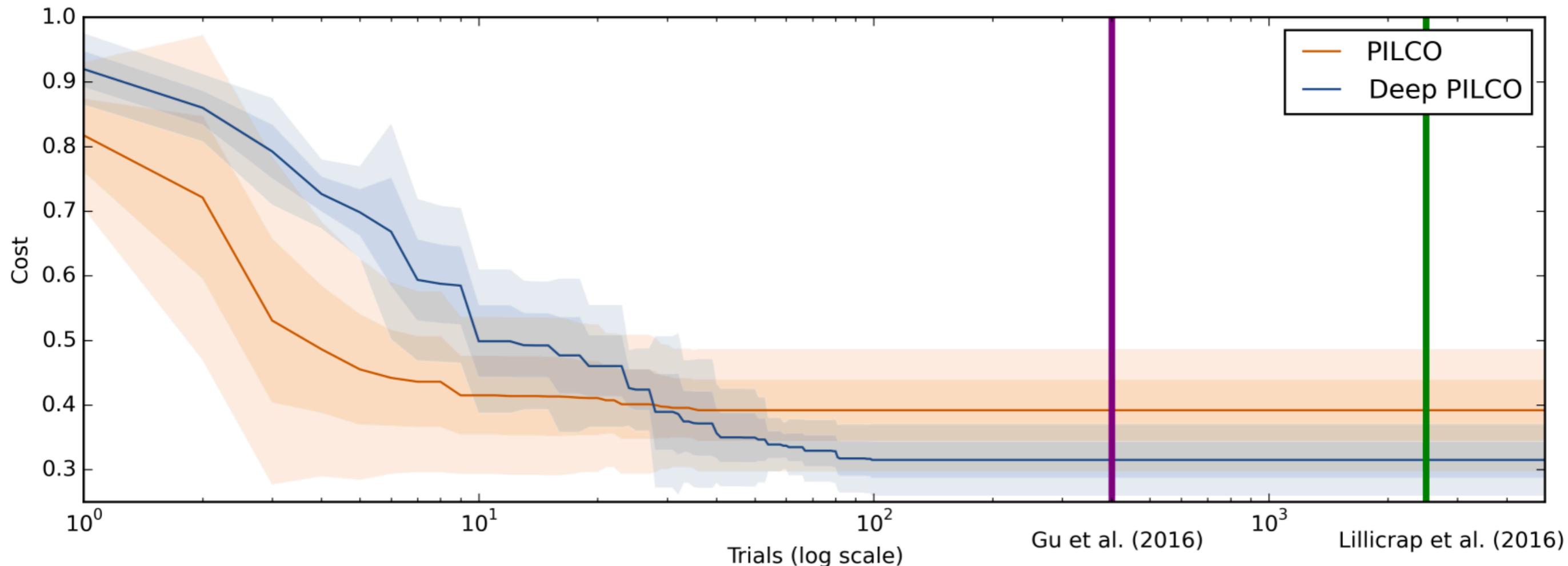


Build Models and Plan with Them

PILCO - Inverting a pendulum

Improving PILCO with Bayesian Neural Network Dynamics Models

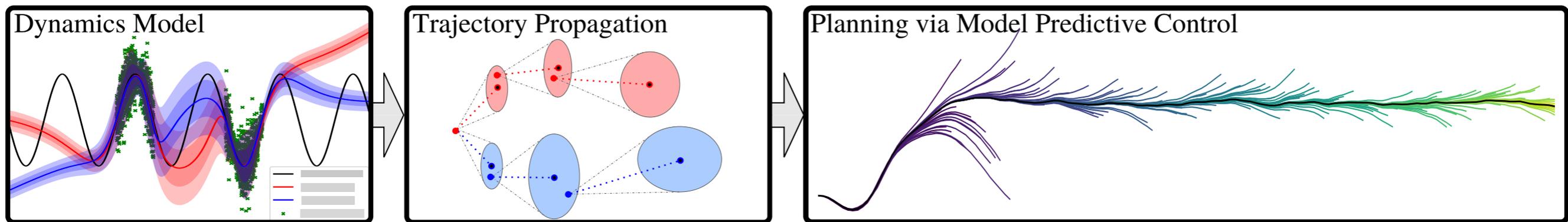
Yarin Gal and Rowan Thomas McAllister and Carl Edward Rasmussen¹



[PILCO] M. Deisenroth et al. **PILCO: A Model-based and Data-Efficient Approach to Policy Search.** ICML 2011

PETS

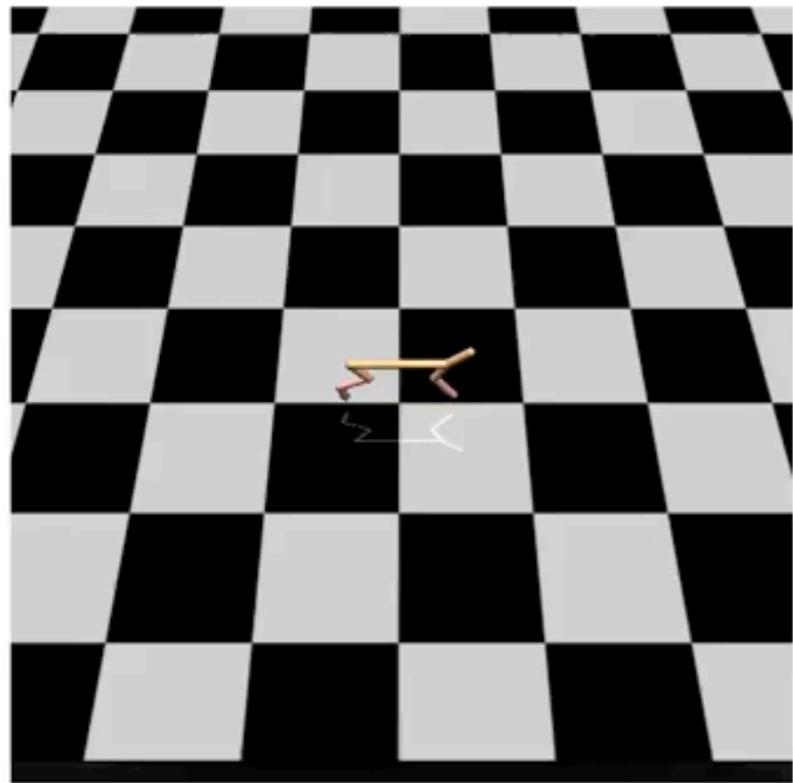
- Employ Neural Network based probabilistic models for Model-based RL
- Use an ensemble of probabilistic models
- Act at test time using CEM method



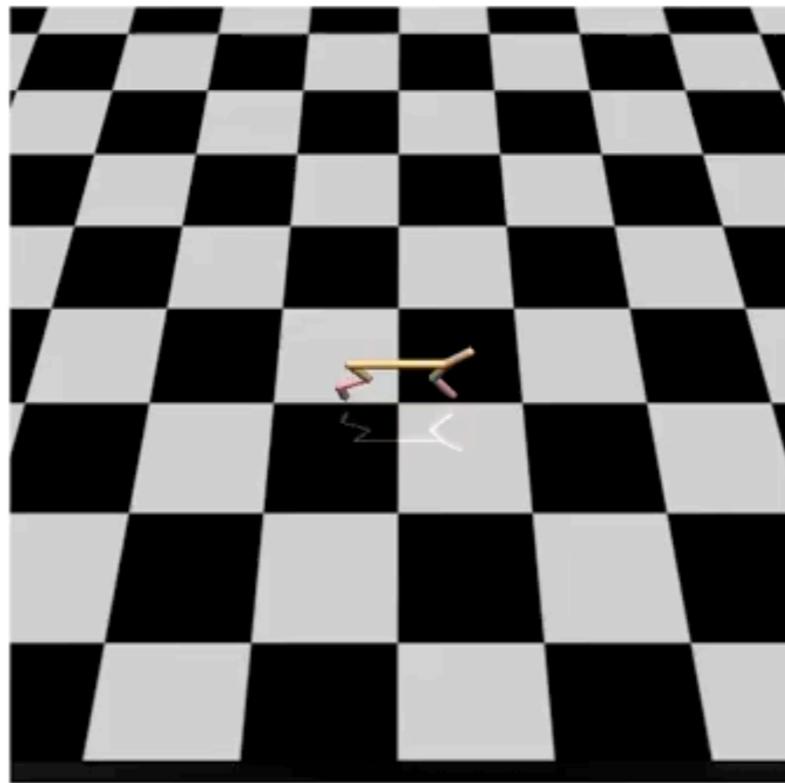
Build Models and Plan with Them

Probabilistic Ensembles with Trajectory Sampling (PETS)

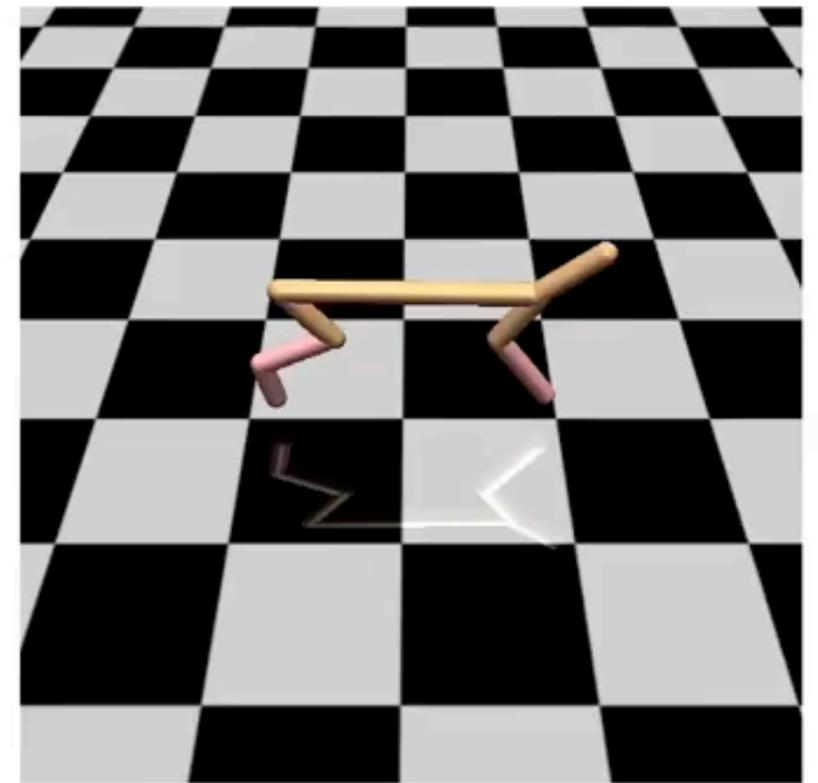
Trial 1



Our method: PE-TS



[Nagabandi 2018]
(D-E)



[Haarnoja 2018]
SAC

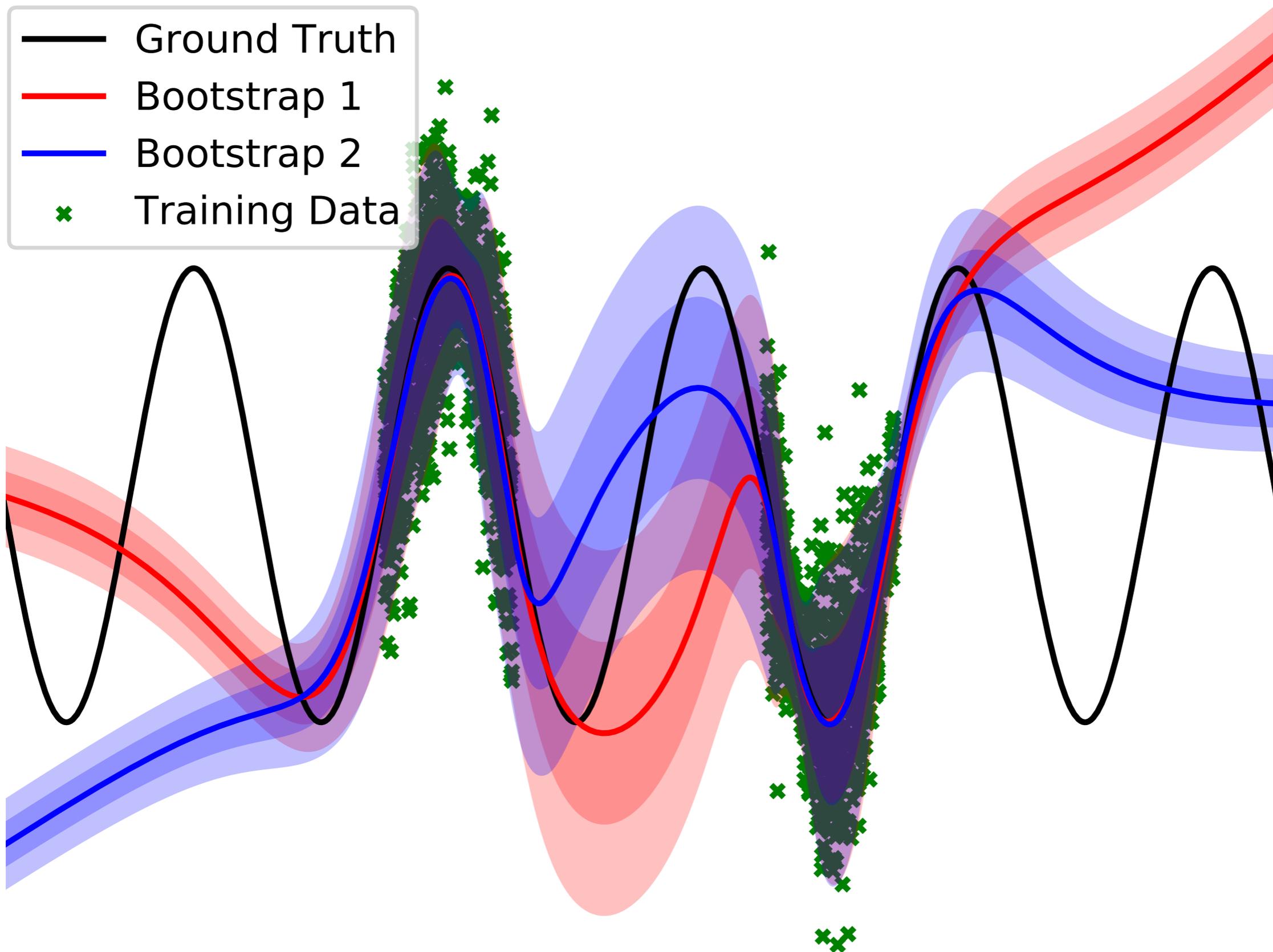
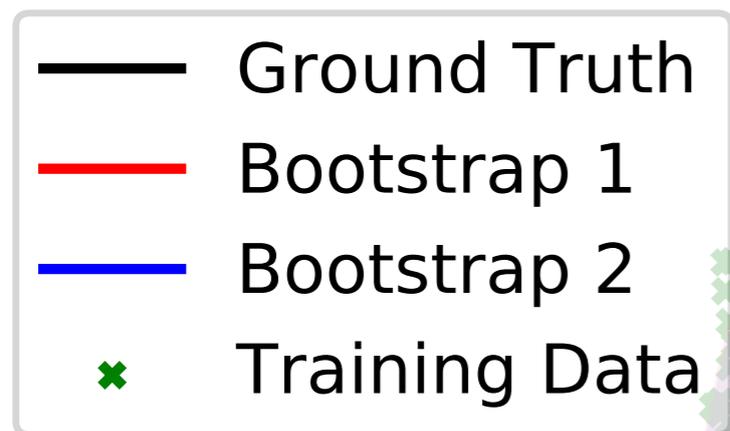
Uncertainties

- Aleatoric uncertainty (inherent stochasticity)
- Epistemic uncertainty (lack of data)

Table 1: Model uncertainties captured.

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

Uncertainties



Planning

- Model predictive control

$$\mathbf{a}_{t:t+T} \doteq \{\mathbf{a}_t, \dots, \mathbf{a}_{t+T}\}$$

$$\arg \max_{\mathbf{a}_{t:t+T}} \sum_{\tau=t}^{t+T} \mathbb{E}_{\tilde{f}} [r(\mathbf{s}_\tau, \mathbf{a}_\tau)].$$

- Cross entropy method

1. Choose an initial parameter vector $\hat{\mathbf{v}}_0$. Let $N^e = \lceil \rho N \rceil$. Set $t = 1$ (level counter).
2. Generate $\mathbf{X}_1, \dots, \mathbf{X}_N \sim_{\text{iid}} f(\cdot; \hat{\mathbf{v}}_{t-1})$. Calculate the performances $S(\mathbf{X}_i)$ for all i , and order them from smallest to largest: $S_{(1)} \leq \dots \leq S_{(N)}$. Let $\hat{\gamma}_t$ be the sample $(1 - \rho)$ -quantile of performances; that is, $\hat{\gamma}_t = S_{(N - N^e + 1)}$.
3. Use the **same** sample $\mathbf{X}_1, \dots, \mathbf{X}_N$ and solve the stochastic program

$$\max_{\mathbf{v}} \frac{1}{N} \sum_{k=1}^N I_{\{S(\mathbf{x}_k) \geq \hat{\gamma}_t\}} \ln f(\mathbf{X}_k; \mathbf{v}). \quad (9)$$

Denote the solution by $\hat{\mathbf{v}}_t$.

4. If some stopping criterion is met, stop; otherwise, set $t = t + 1$, and return to Step 2.

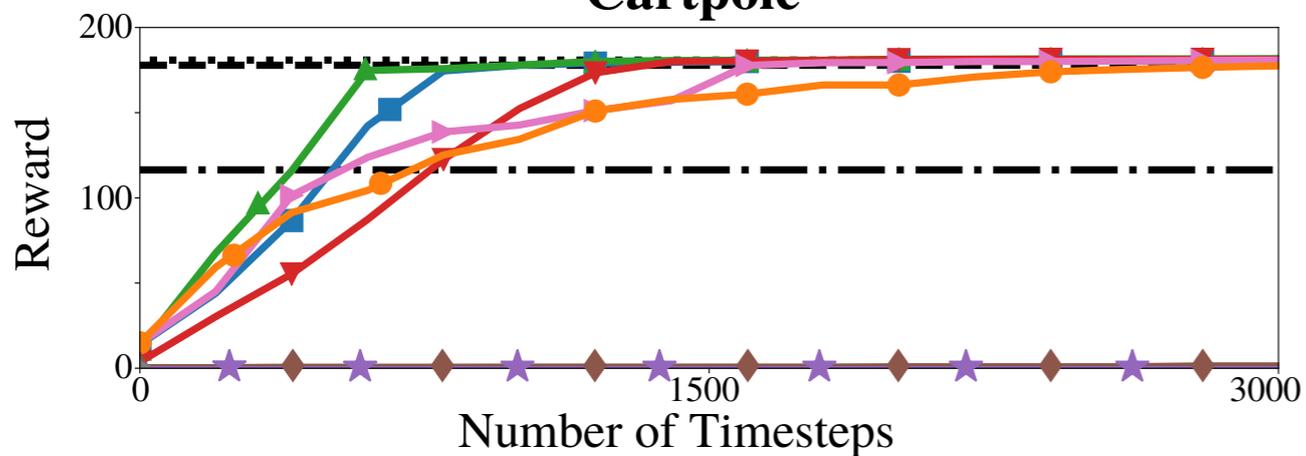
PETS

Algorithm 1 Our model-based MPC algorithm ‘*PETS*’:

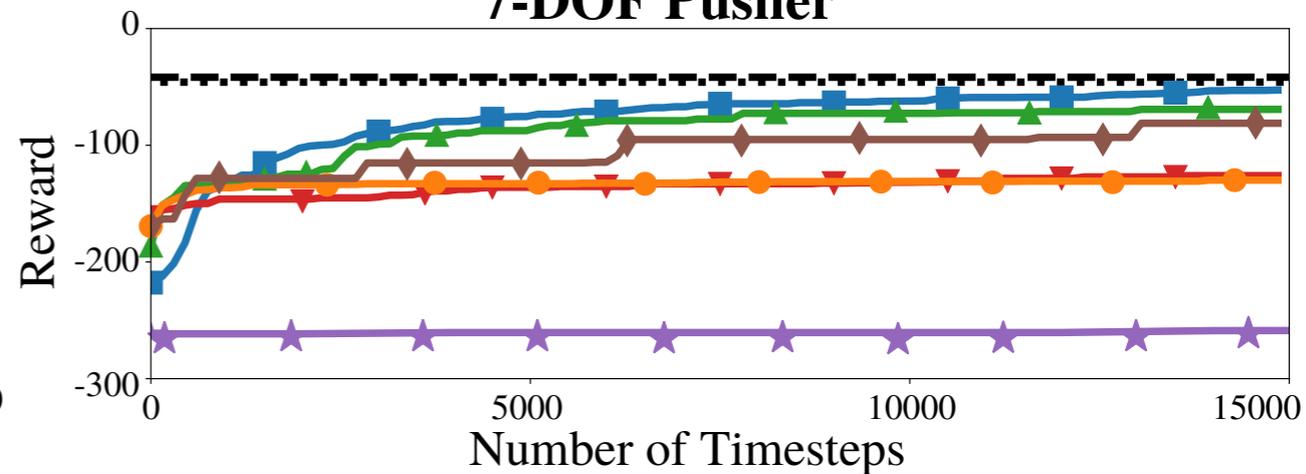
- 1: Initialize data \mathbb{D} with a random controller for one trial.
 - 2: **for** Trial $k = 1$ to K **do**
 - 3: Train a *PE* dynamics model \tilde{f} given \mathbb{D} .
 - 4: **for** Time $t = 0$ to TaskHorizon **do**
 - 5: **for** Actions sampled $\mathbf{a}_{t:t+T} \sim \text{CEM}(\cdot)$, 1 to $\tilde{N}\text{Samples}$ **do**
 - 6: Propagate state particles \mathbf{s}_τ^p using *TS* and $\tilde{f} | \{\mathbb{D}, \mathbf{a}_{t:t+T}\}$.
 - 7: Evaluate actions as $\sum_{\tau=t}^{t+T} \frac{1}{P} \sum_{p=1}^P r(\mathbf{s}_\tau^p, \mathbf{a}_\tau)$
 - 8: Update $\text{CEM}(\cdot)$ distribution.
 - 9: Execute first action \mathbf{a}_t^* (only) from optimal actions $\mathbf{a}_{t:t+T}^*$.
 - 10: Record outcome: $\mathbb{D} \leftarrow \mathbb{D} \cup \{\mathbf{s}_t, \mathbf{a}_t^*, \mathbf{s}_{t+1}\}$.
-

Experiments

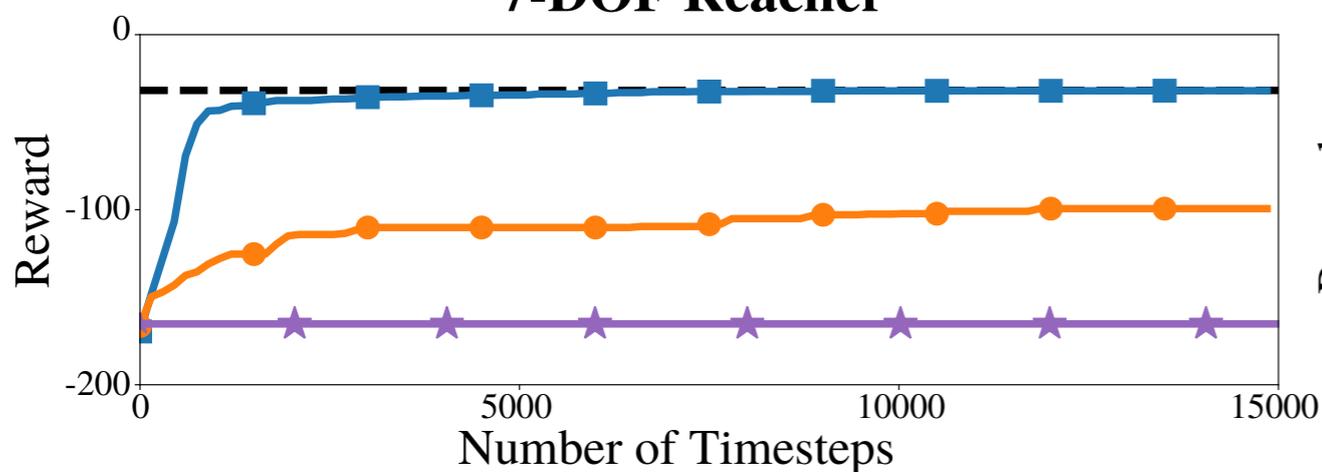
Cartpole



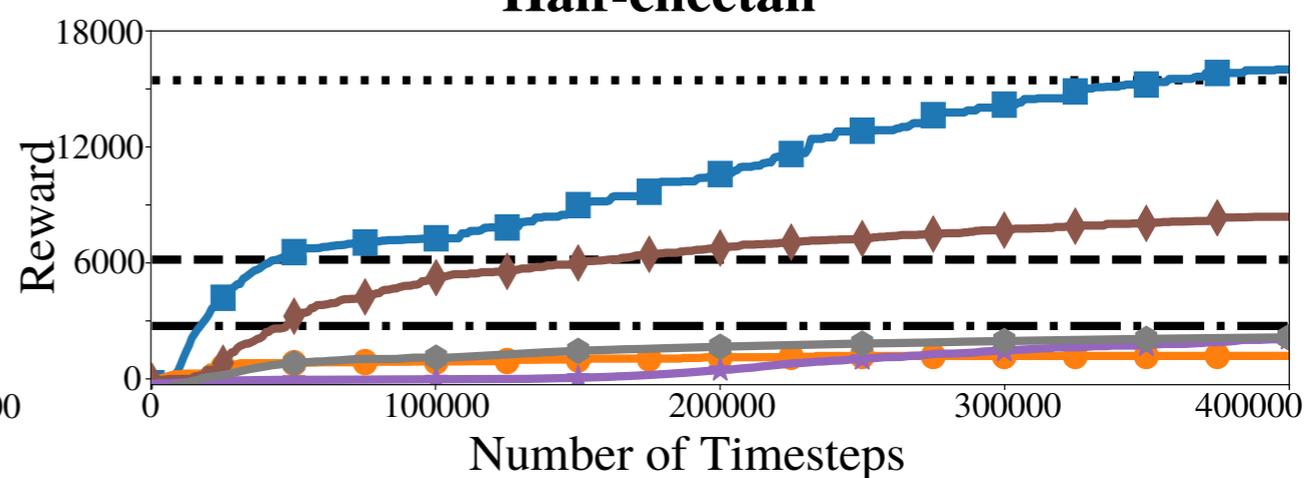
7-DOF Pusher



7-DOF Reacher



Half-cheetah



Experiments

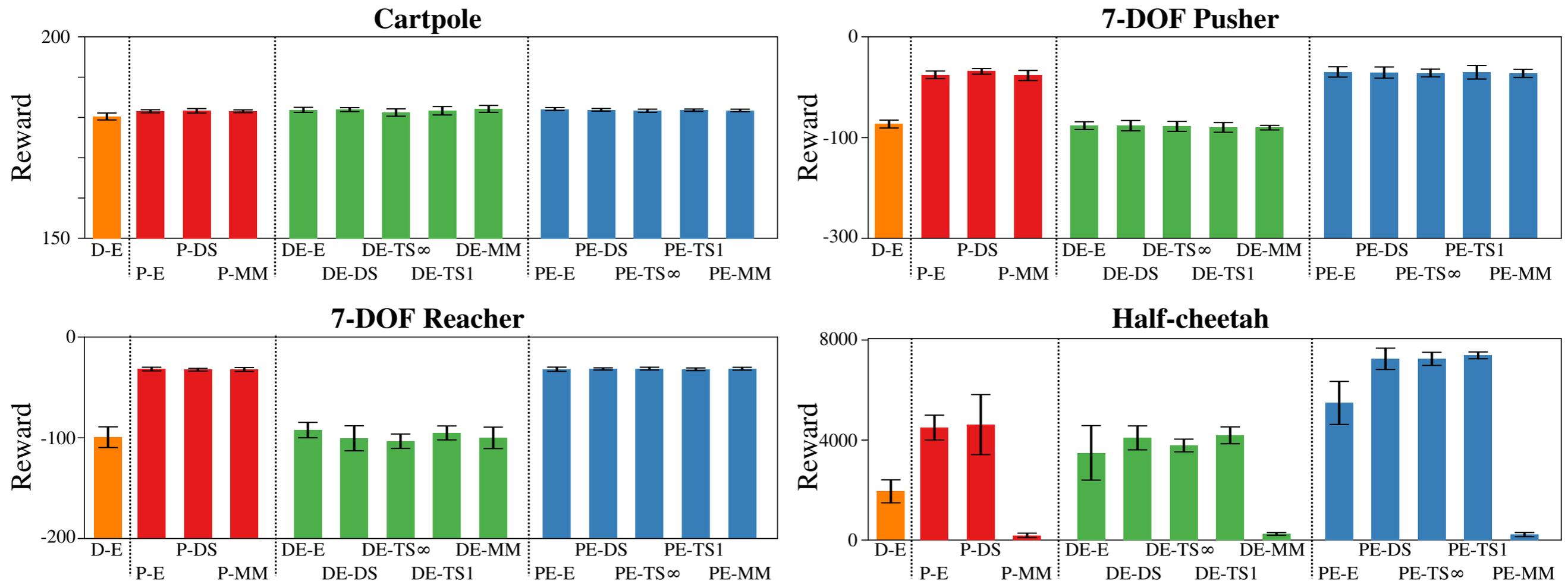


Figure 4: Final performance for different tasks, models, and uncertainty propagation techniques. The model choice seems to be more important than the technique used to propagate the state/action space. Among the models the ranking in terms of performance is: $PE > P > DE > D$. A linear model comparison can also be seen in Appendix A.10.

Thank you