

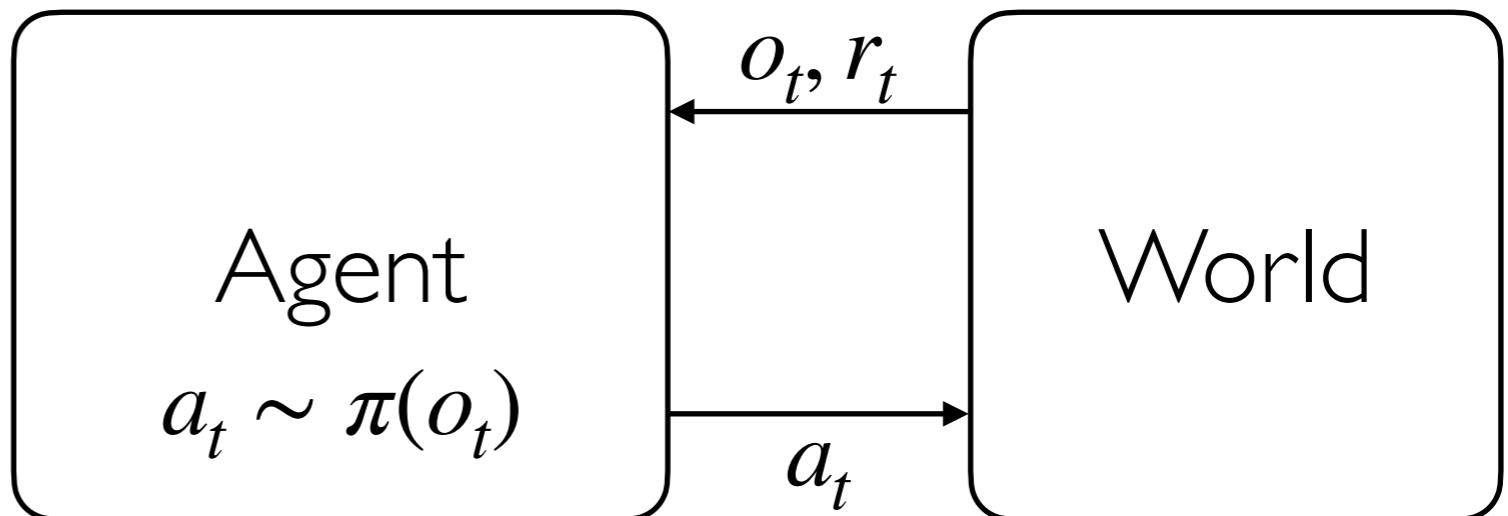
Learning Models for RL

Saurabh Gupta

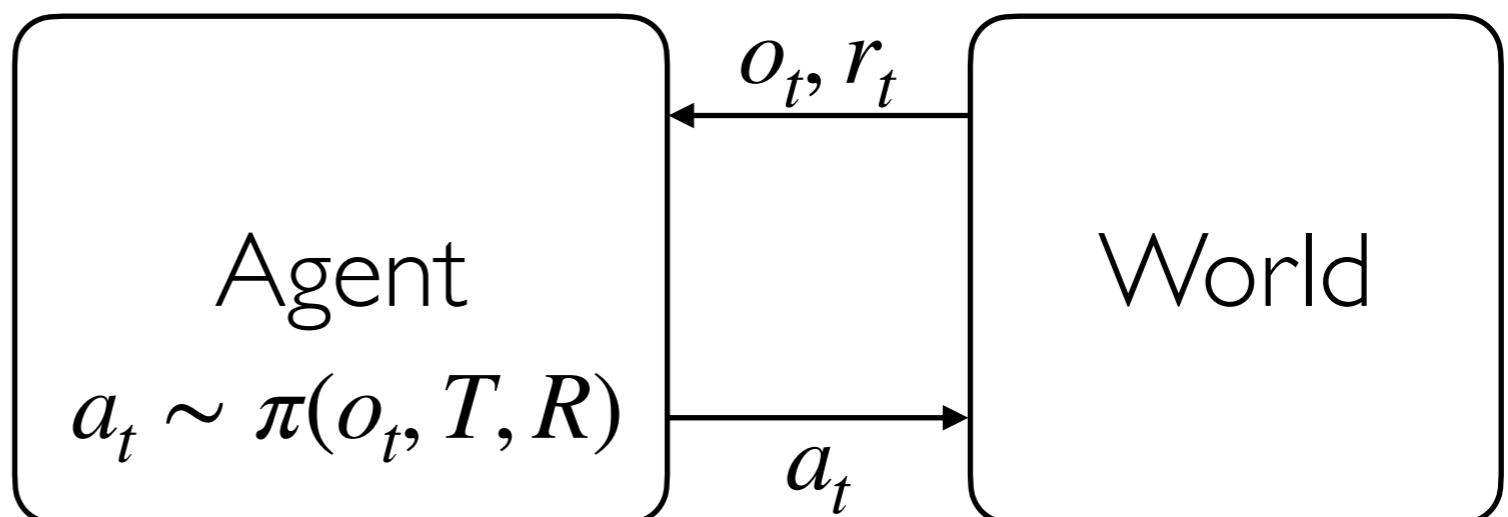
Solving MDPs

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

Most General Case



More Specific Case



Fully Observed System

$$o_t = s_t$$

Known Transition Function

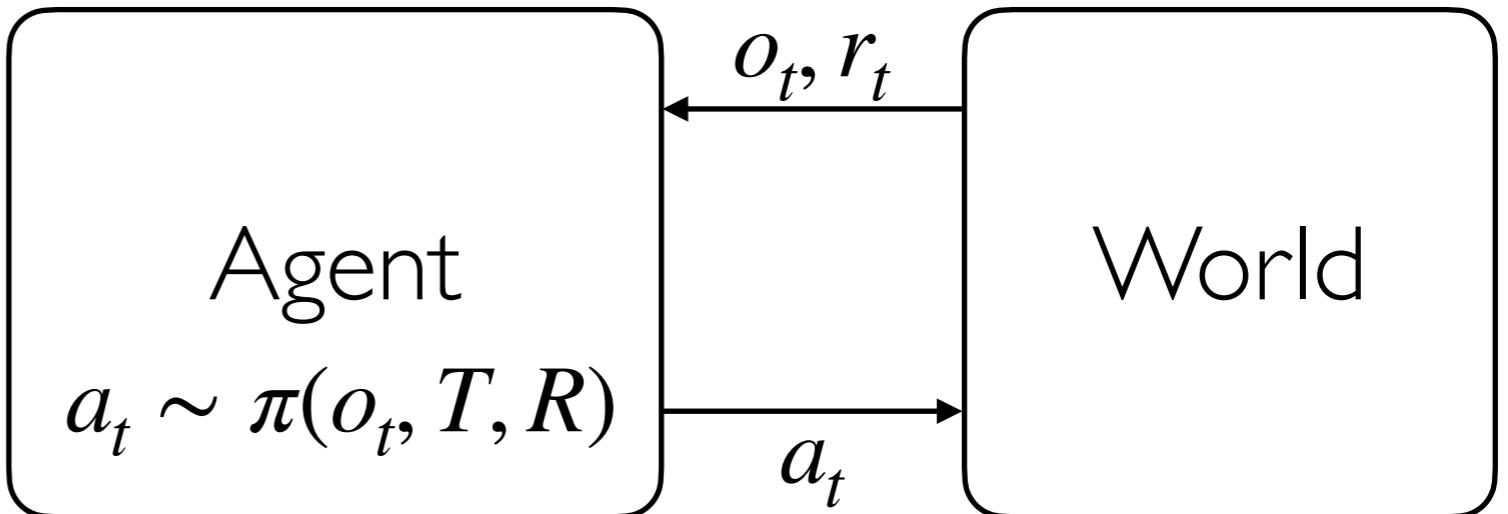
$$s_{t+1} \sim T(s_t, a_t)$$

Known Reward Function

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

Build Models and Plan with Them

Relax Assumptions For Model-Based RL



Fully Observed System

$$o_t = s_t$$

Known Transition Function

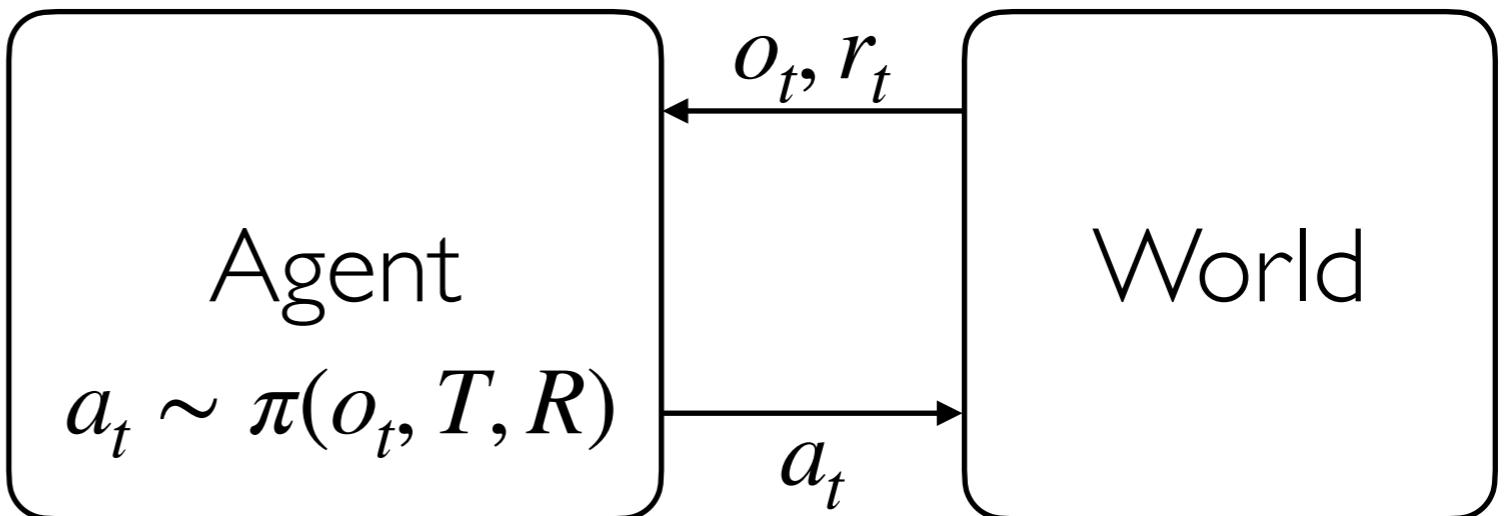
$$s_{t+1} \sim T(s_t, a_t)$$

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

I. 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

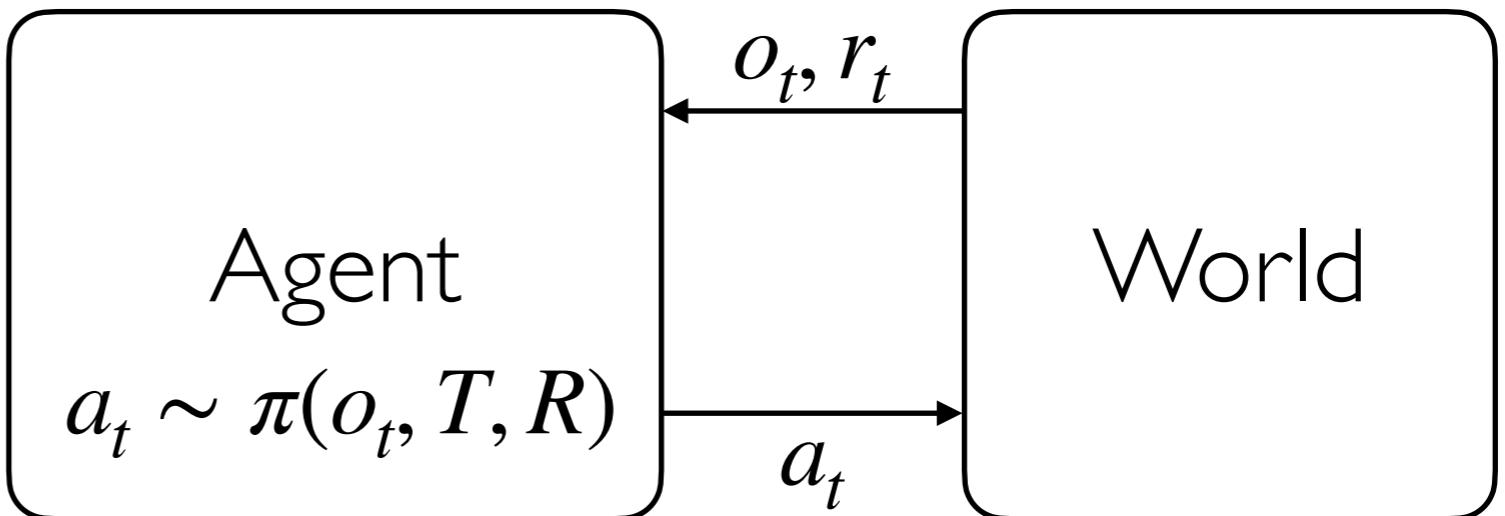
[PILCO] M. Deisenroth et al. PILCO:A Model-based and Data-Efficient Approach to Policy Search. ICML 2011

[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 $o_t = s_t$

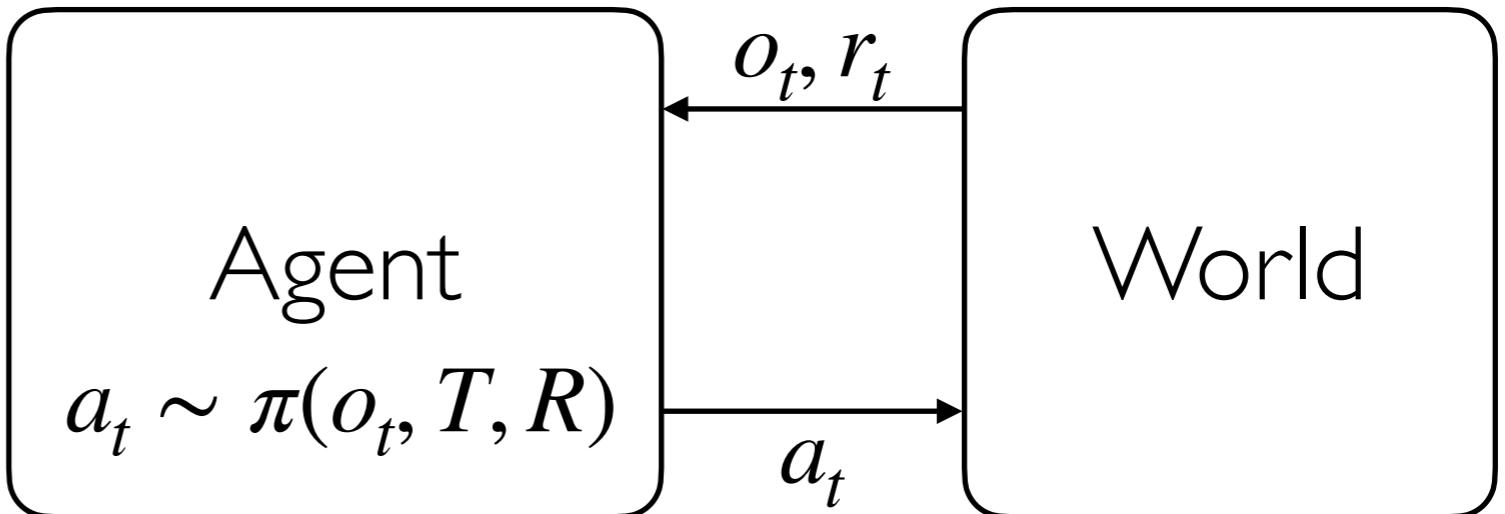
~~Known Transition Function~~ $s_{t+1} \sim T(s_t, a_t)$

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~~ $o_t \equiv s_t$

~~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) $o_{t+1} \sim \hat{T}(o_t, a_t)$

I. 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~~

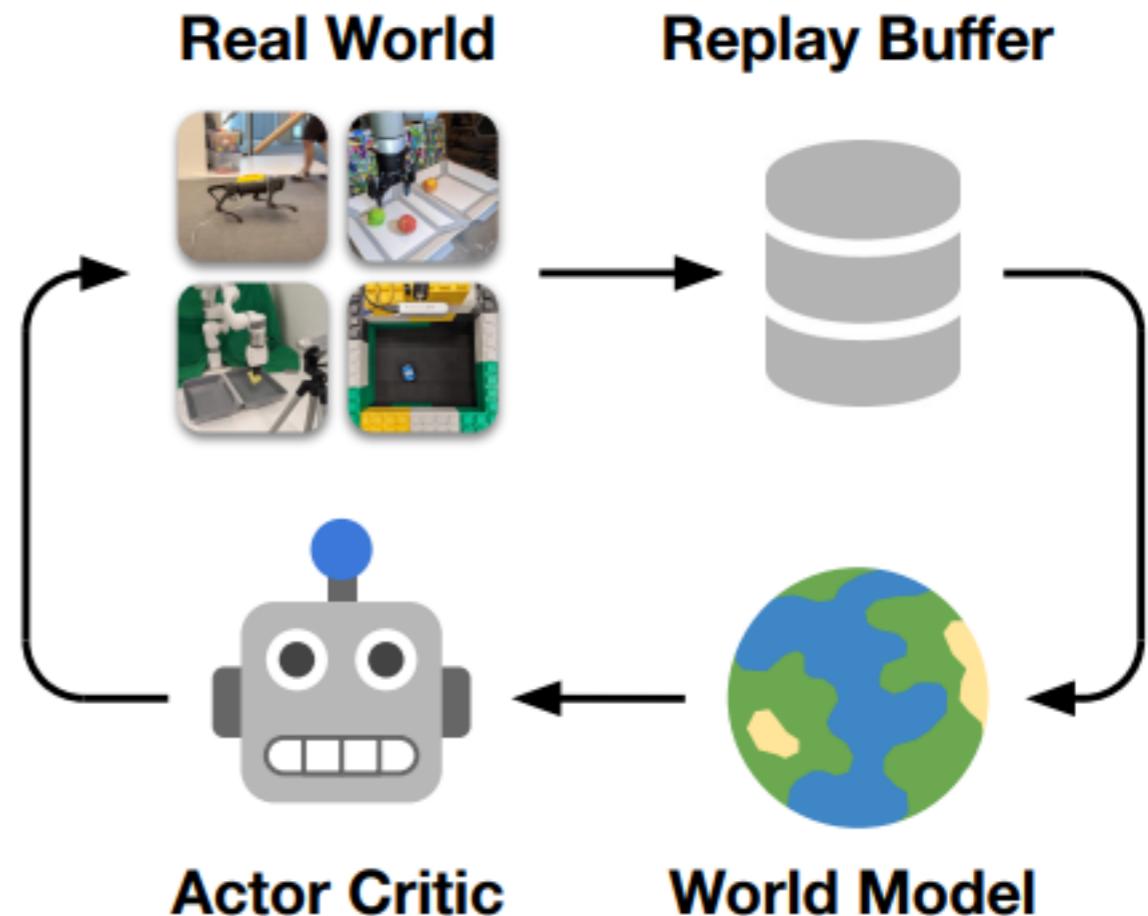
Neural Network

3. Plan using function \hat{T}

Build a Models and Learn a Policy within it

Examples:

- ME-TRPO
- World models
- Dreamers
 - Dreamer V1
 - Dreamer V2
 - Day Dreamer
 - ...



Build Models and Plan with Them

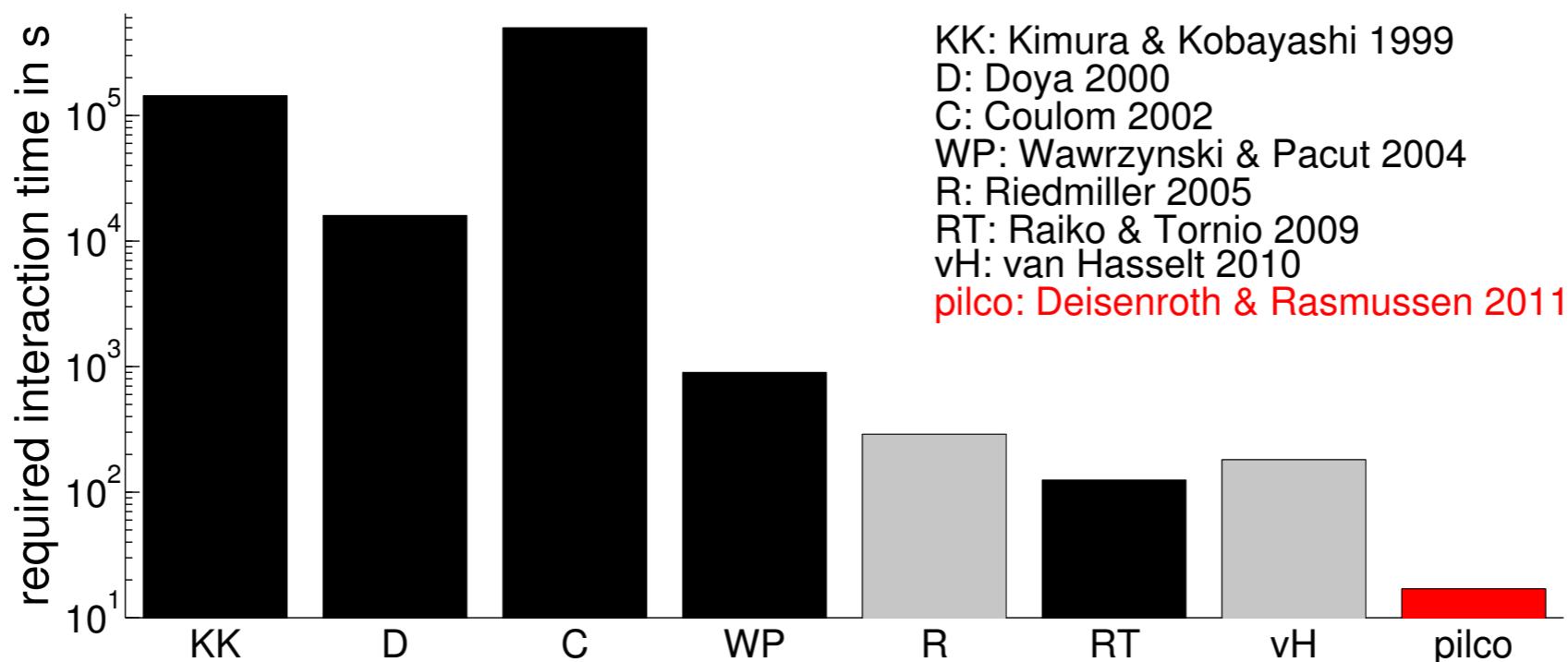
PILCO - Inverting a pendulum



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}

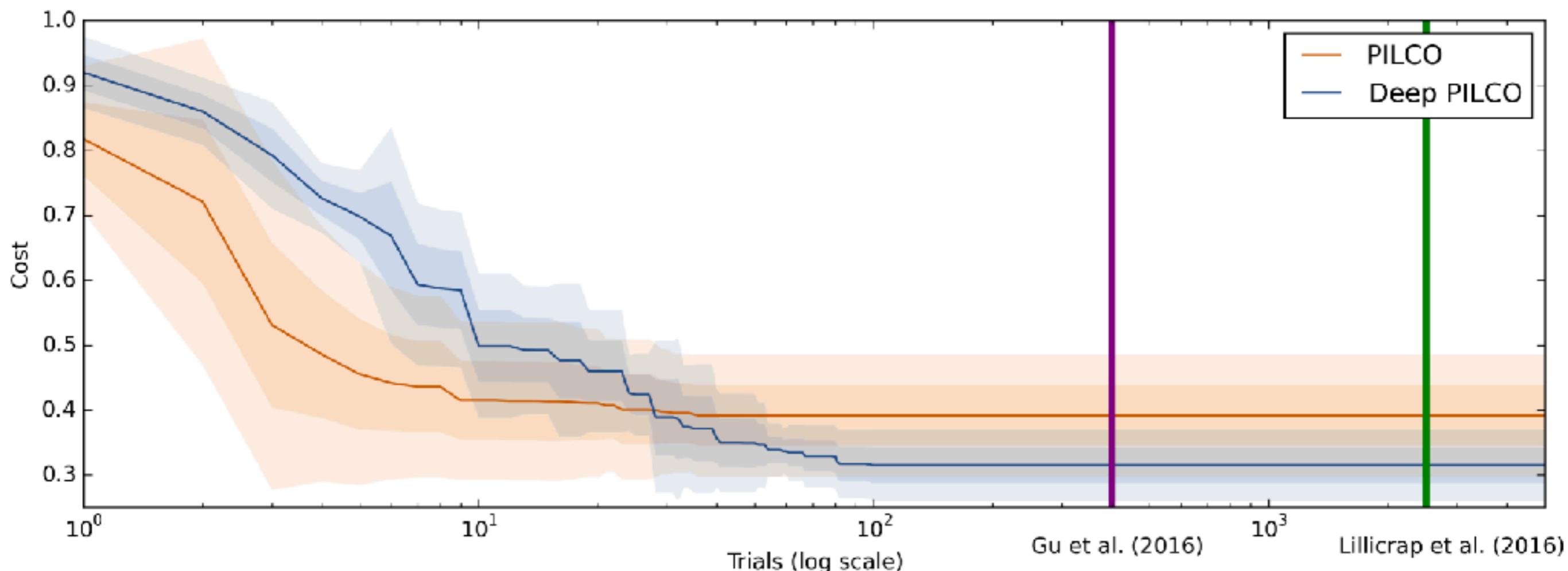


Build Models and Plan with Them

PILCO - Inverting a pendulum

Improving PILCO with Bayesian Neural Network Dynamics Models

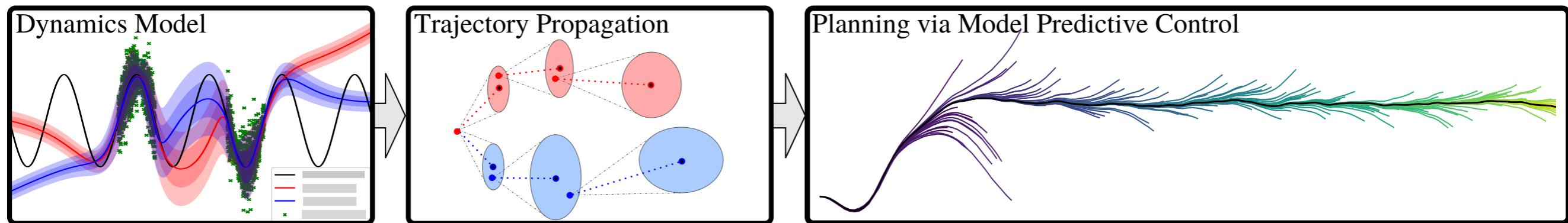
Yarin Gal and Rowan Thomas McAllister and Carl Edward Rasmussen¹



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

PETS

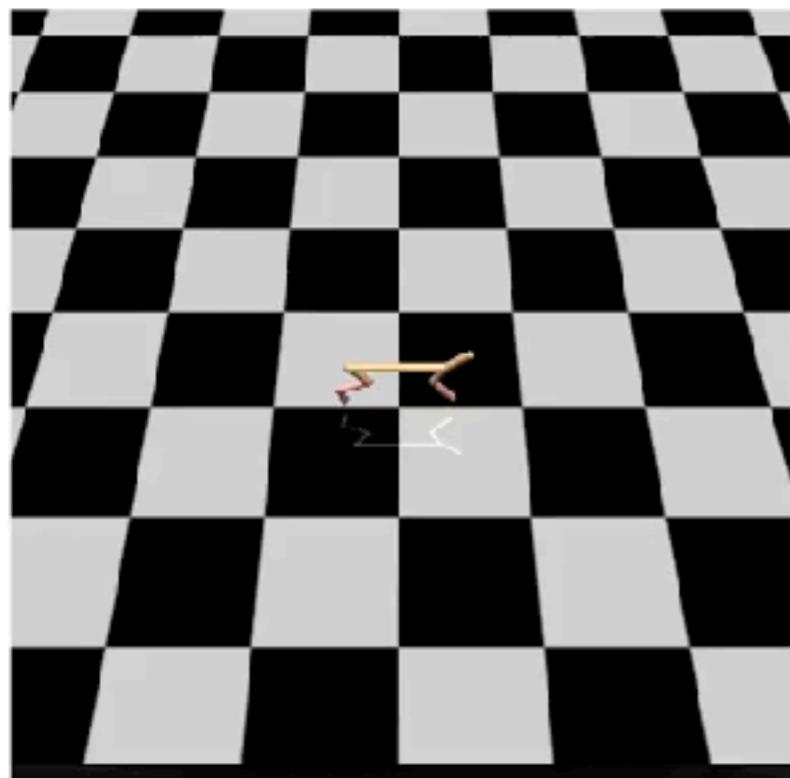
- Employ Neural Network based probabilistic models for Model-based RL
- Use an ensemble of probabilistic models
- Act at test time using CEM method



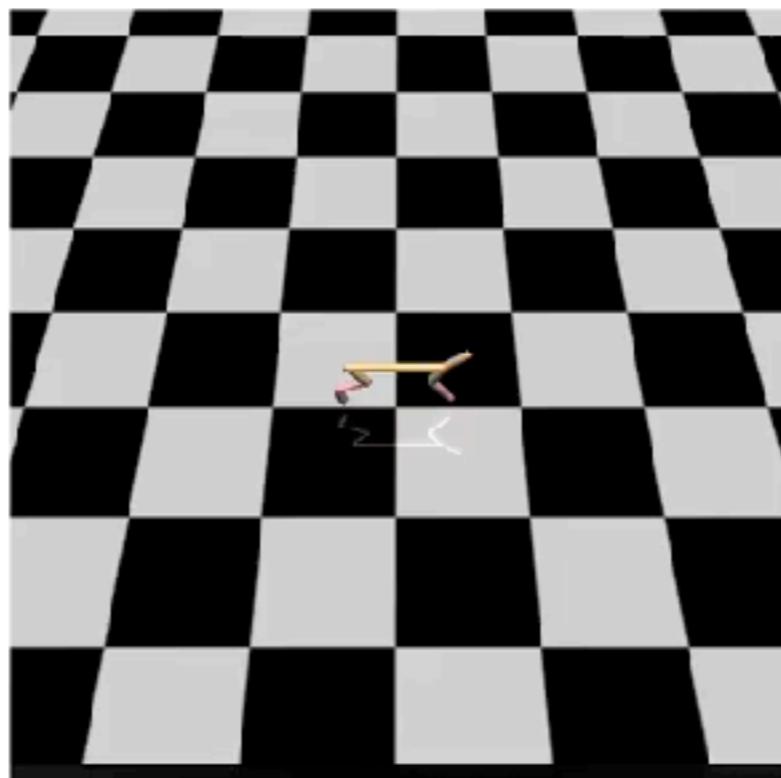
Build Models and Plan with Them

Probabilistic Ensembles with Trajectory Sampling (PETS)

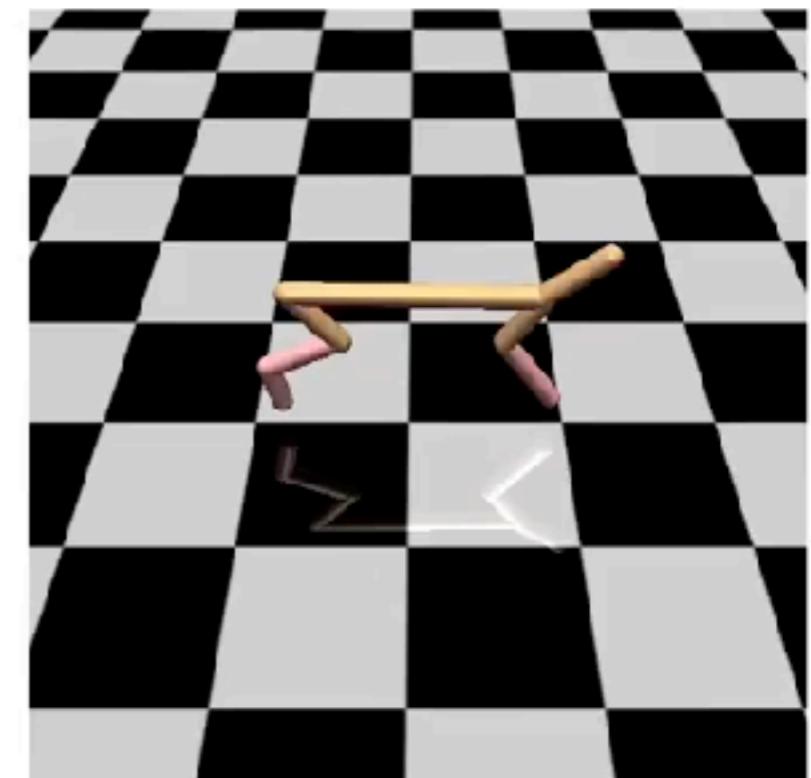
Trial 1



Our method: PE-TS



[Nagabandi 2018]
(D-E)



[Haarnoja 2018]
SAC

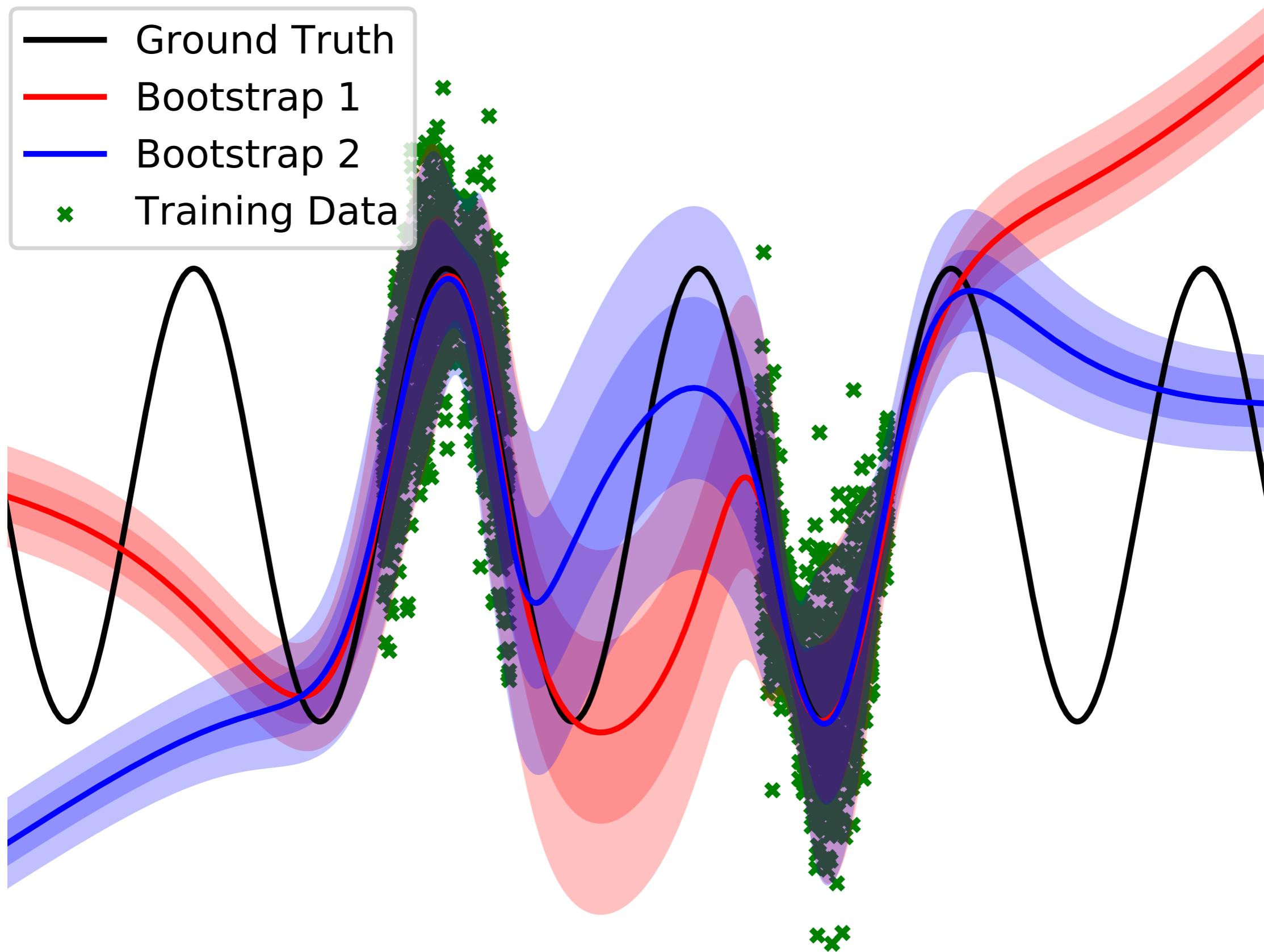
Uncertainties

- 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 $N^e = \lceil \varrho N \rceil$. 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 $(1 - \varrho)$ -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

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

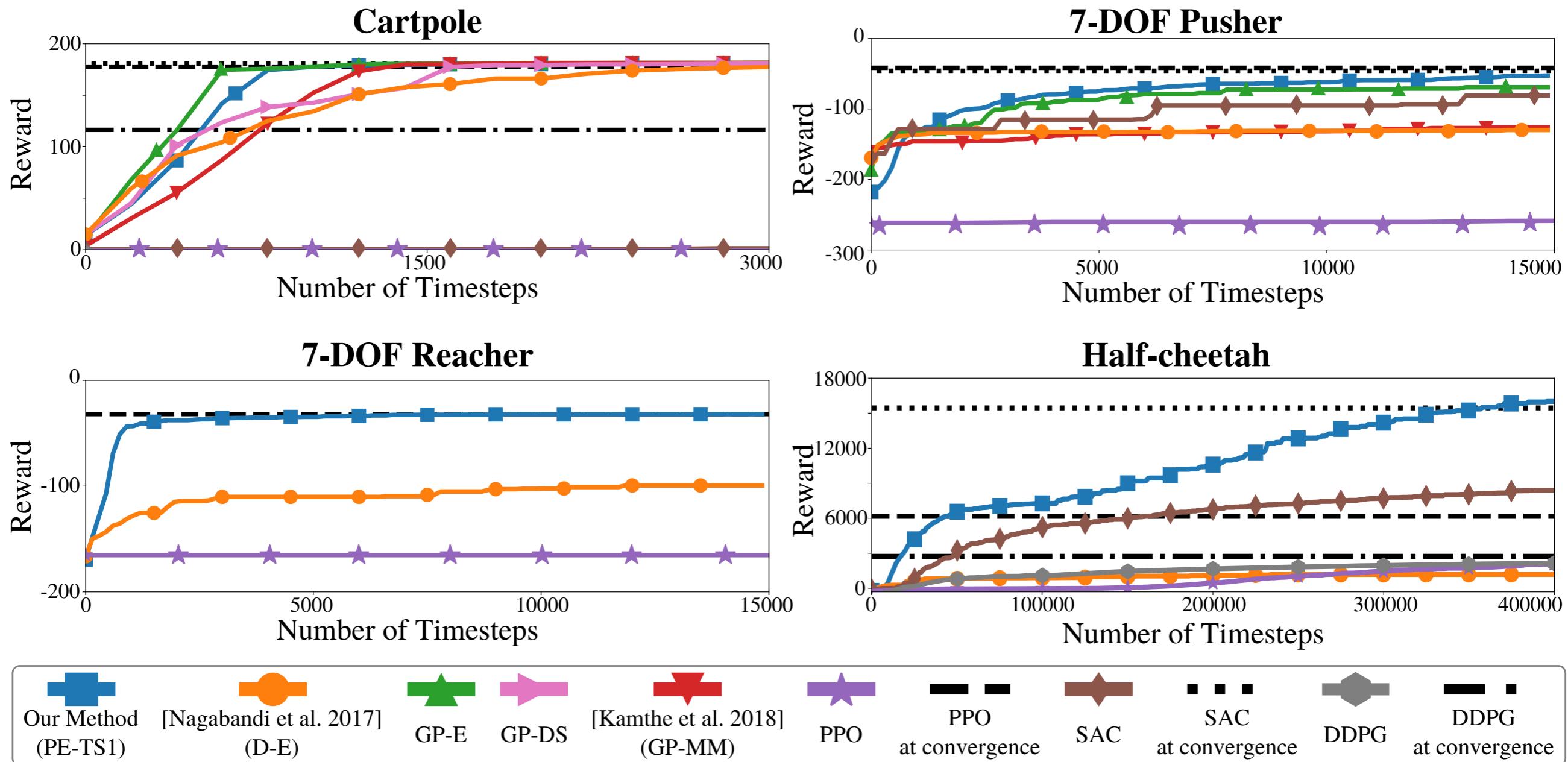
Denote the solution by $\hat{\mathbf{v}}_t$.

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

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 NSamples **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



Experiments

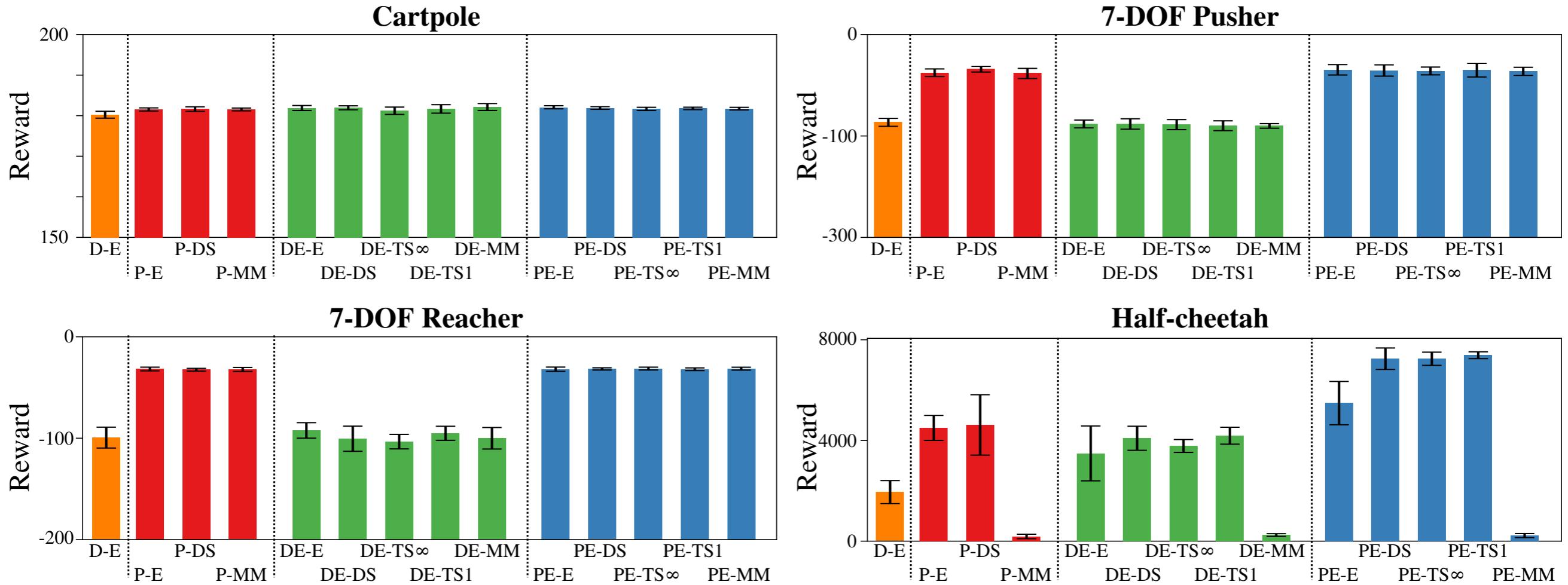


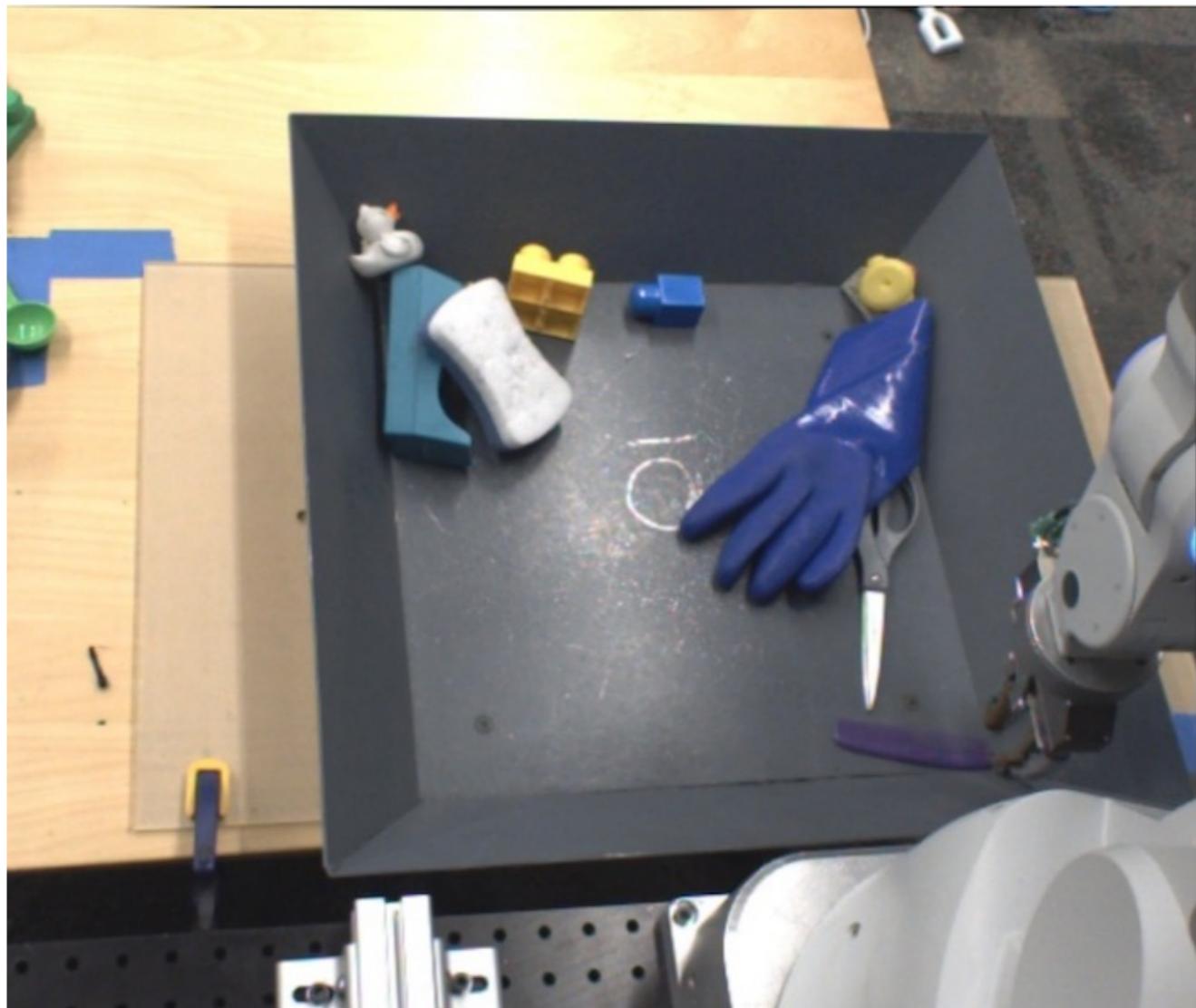
Figure 4: Final performance for different tasks, models, and uncertainty propagation techniques. The model choice seems to be more important than the technique used to propagate the state/action space. Among the models the ranking in terms of performance is: $PE > P > DE > D$. A linear model comparison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

Deep Visual Foresight for Planning Robot Motion

Chelsea Finn^{1,2} and Sergey Levine^{1,2}

- Build a action conditioned predictive model
 - in pixel space
 - using largely unsupervised interaction
- Use predictive model for planning

Data collection



- Action space?
 - Target location of gripper

Data collection

Data Collection

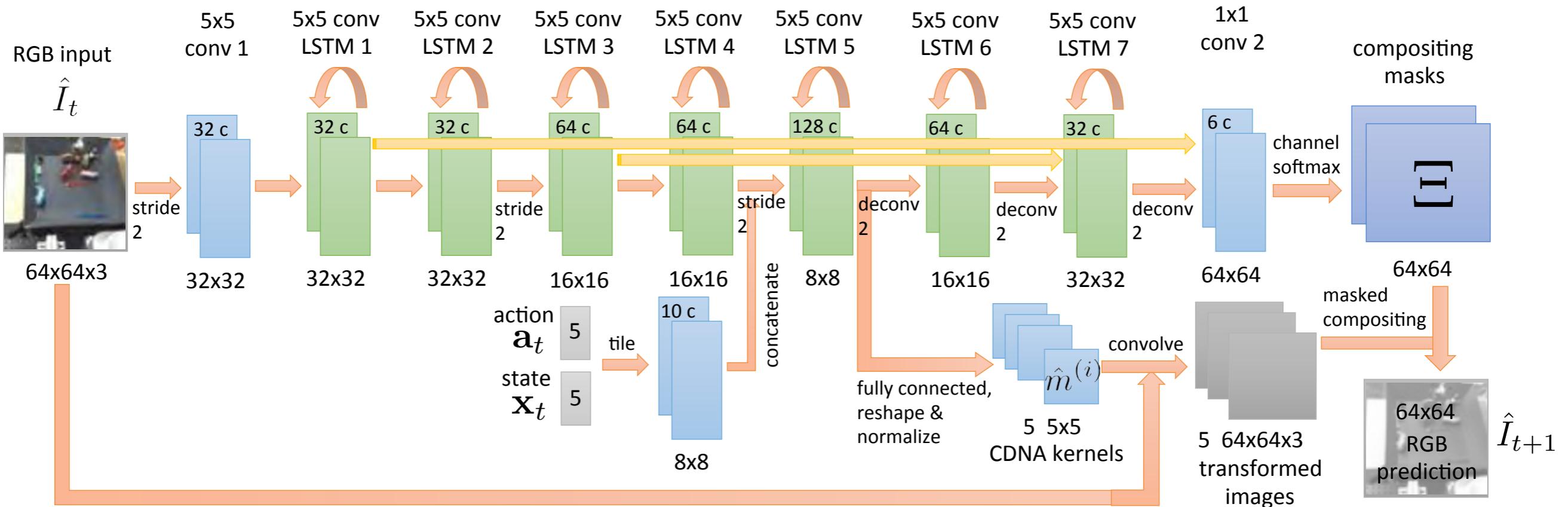


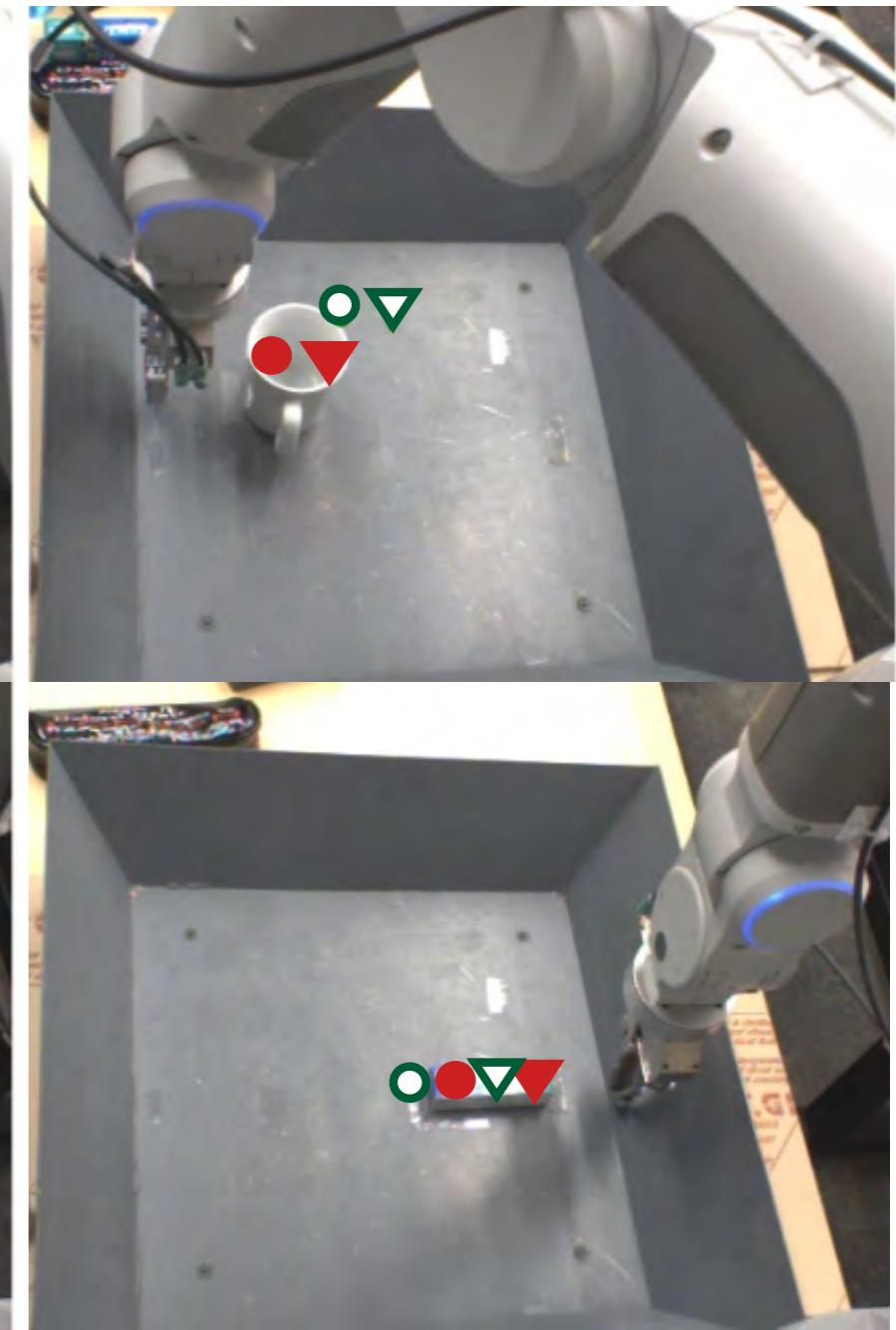
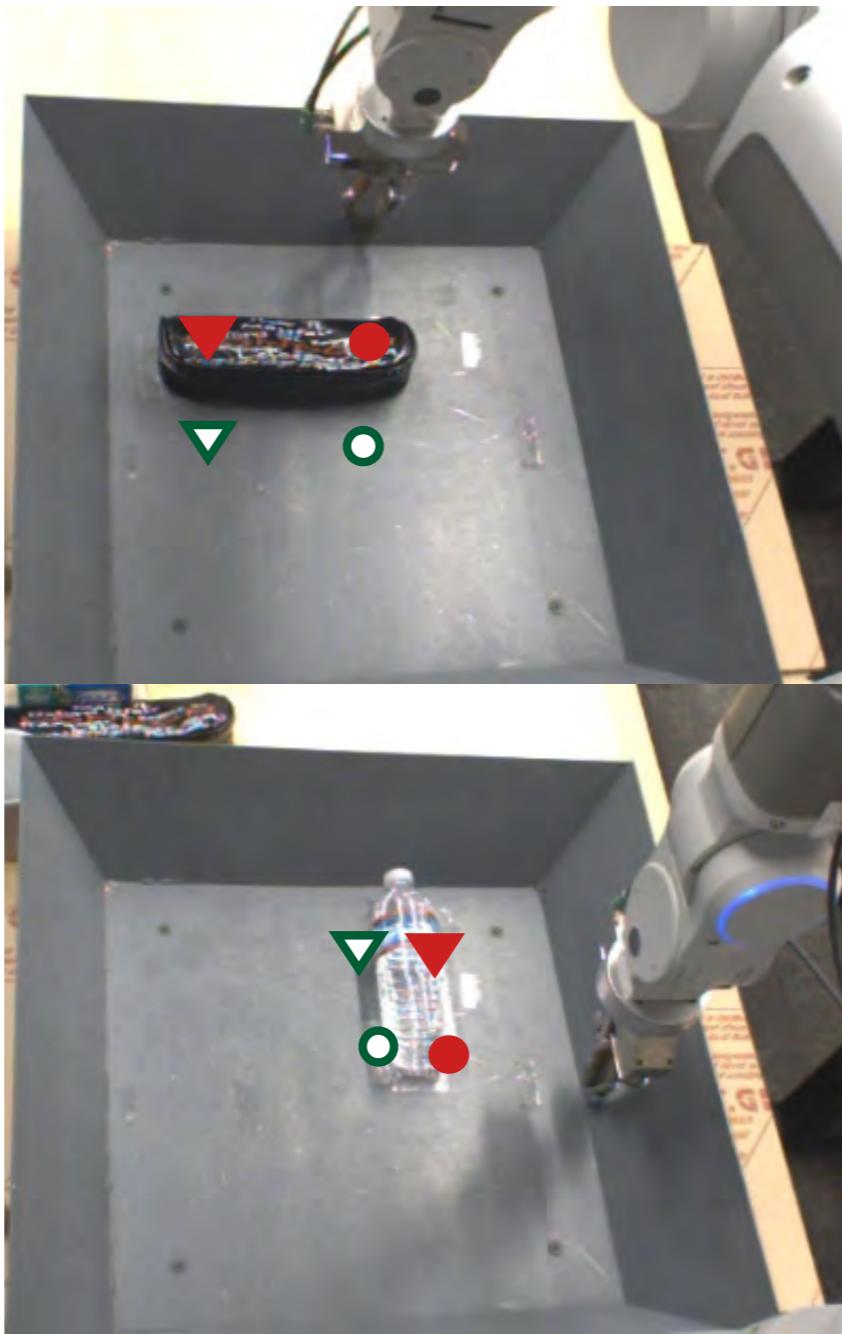
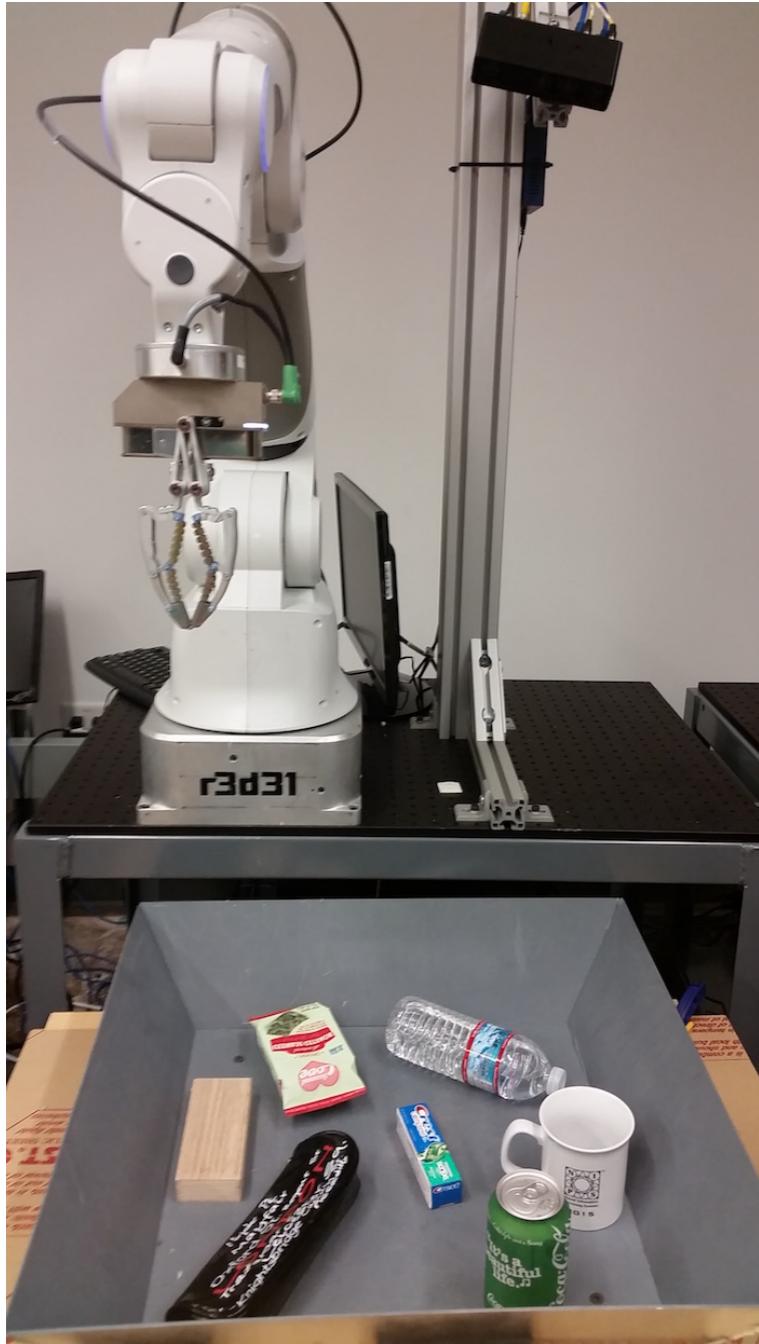
Fig. 2. Our video prediction model predicts stochastic pixel flow transformation from the current frame to the next frame, which allows it to generate predictions for subsequent images conditioned on a sequence of future actions. Predicted stochastic flow is parameterized by a set of normalized convolution filters that give rise to an independent Gaussian distribution over future images. Internally, our model predicts multiple stochastic flow channels, in order to process multiple separate moving objects. These channels are composited using learned object masks. The only supervision to the model consists of video, action, and state sequences, with no explicit supervision over flows or masks. The flows and masks therefore are an emergent property of the model, which we will exploit to perform visual MPC using high-level user commands. Further details about the model may be found in prior work [3].

Forward Models with Images



- Extremely blurry predictions, why?

Task Specification



Test time planning

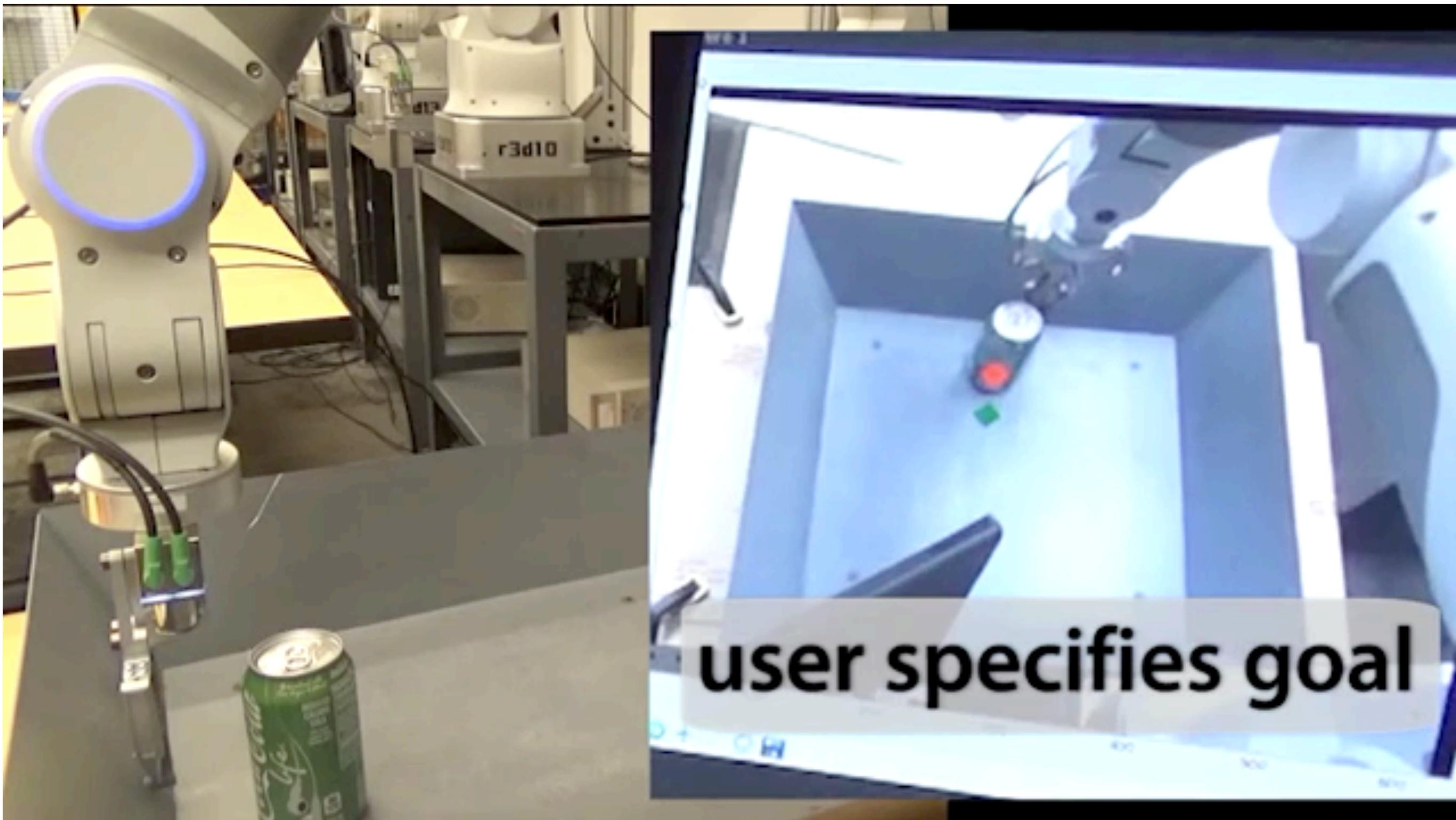
Algorithm 1 Visual MPC with Deep Predictive Models

```
0: inputs: predictive model  $\mathcal{M}$ , designated pixel  $\{d_0\}$ , pixel goal  
position  $\{g\}$   
1: for  $t = 1 \dots T$  do  
2:   Initialize  $Q_1$  with uniform distribution.  
3:   for  $j = 1 \dots J_t$  do  
4:     Sample  $M$  action sequences  $\{\mathbf{a}_{t:t+H-1}^{(m)}\}$  from  $Q_j$ .  
5:     Use model  $\mathcal{M}$  to compute the distributions over future  
pixel locations  $\mathbf{P}_{t+H-1}^{(m)}$  using Equation 1  
6:     Fit multivariate Gaussian distribution  $Q_{j+1}$  to  $K$  samples  
with highest probability of success  $\mathbf{P}_{t+H-1}^{(m)}$   
7:   end for  
8:   Execute action  $\mathbf{a}_t^*$  with highest probability of success  
9:   Observe new image  $I_{t+1}$ .  
10:  Set next designated pixel location  $d_{t+1}$  using optical flow  
computed from image observations  $I_{t:t+1}$ , and  $d_t$ .  
11: end for
```

Results

method	mean pixel distance
initial pixel position	5.10 ± 2.25
1) random actions	4.05 ± 1.75
2) move end-effector to goal	3.79 ± 2.66
3) move end-effector along vector (with replanning)	3.19 ± 1.68
visual MPC (ours)	2.52 ± 1.06

Results



Extensions

<https://bair.berkeley.edu/blog/2018/11/30/visual-rl/>

Thank you