

COM SCI 188

# Intro to Robotics

## Lecture 12

Yuchen Cui  
Winter 2026

# Agenda

- Announcements
- Recap: Reinforcement Learning
- TODAY: Imitation Learning
  - Dynamic Movement Primitives
  - Behavioral Cloning
  - Inverse Reinforcement Learning

# Announcements

- Coding Assignment 3 is out, due 2/23
- Midterm next Thursday 2/19 (in-class)



**Reinforcement Learning**  
**First trial...**

# Deep Reinforcement Learning for Robotics: A Survey of Real-World Successes

Chen Tang<sup>1,\*</sup>, Ben Abbatematteo<sup>1,\*</sup>,  
Jiaheng Hu<sup>1,\*</sup>, Rohan Chandra<sup>2</sup>,  
Roberto Martín-Martín<sup>1</sup>, Peter Stone<sup>1,3</sup>

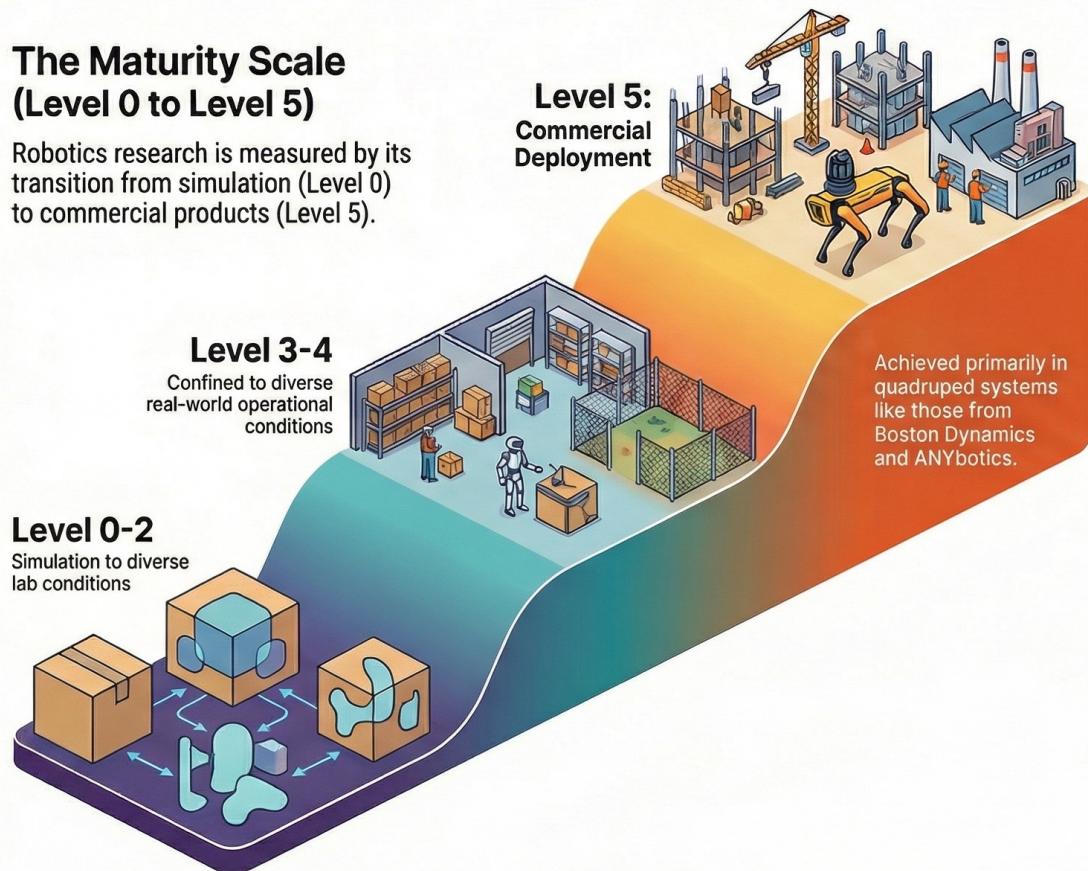
<https://arxiv.org/pdf/2408.03539>

# The State of Deep Reinforcement Learning (DRL) in Real-World Robotics

A survey of DRL successes in robotics, contrasting high-success areas like quadruped locomotion with emerging fields like human-robot interaction, and outlining the maturity from simulation to commercial deployment.

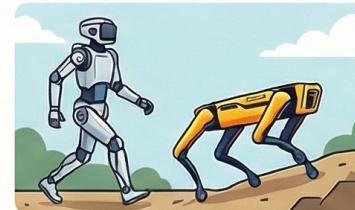
## The Maturity Scale (Level 0 to Level 5)

Robotics research is measured by its transition from simulation (Level 0) to commercial products (Level 5).



## Competency Maturity & Future Hurdles

### Locomotion vs. Interaction Gaps



**Mature Competency: Locomotion**  
Locomotion is highly mature.



**Early Stage: Interaction & Coordination**  
Human-Robot Interaction and Multi-Robot coordination remain in early stages.

### Success in Restricted Domains

DRL excels when tasks are enumerable a priori, like grasping and assembly.



### Three Major Open Challenges



**Improving Sample Efficiency**  
Learning faster from fewer examples.



**Enabling Safe Real-World Learning**  
Preventing damage during training.



**Mastering Long-Horizon Tasks**  
Planning and executing multi-step sequences.

# Practical Challenges of RL for Robotics

- Sample Efficiency
- Safety
- Reset
- Reward Specification



Specifying reward for RL is hard...



**Reward hacking:** RL agent learns to exploit loopholes or unintended behaviors in its reward function to achieve high rewards without actually accomplishing the intended task

# Imitation Learning

# Why learn from demonstrations?

- Natural and expressive
- No expert knowledge required
- Valuable human intuition
- Program new tasks as-needed

Human babies imitate

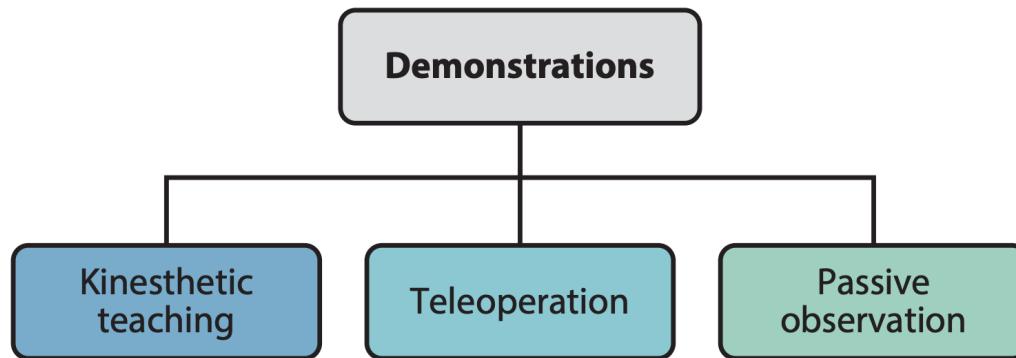


# What to imitate?

Demonstration(s) -> Autonomous Behavior

- ***Dynamic Movement Primitives (DMP)***: replay the **motion**
- ***Behavior Cloning (BC)***: supervised learning of expert **policy**
- ***Inverse Reinforcement Learning (IRL)***: inferring the underlying **intent**

# Types of Demonstrations



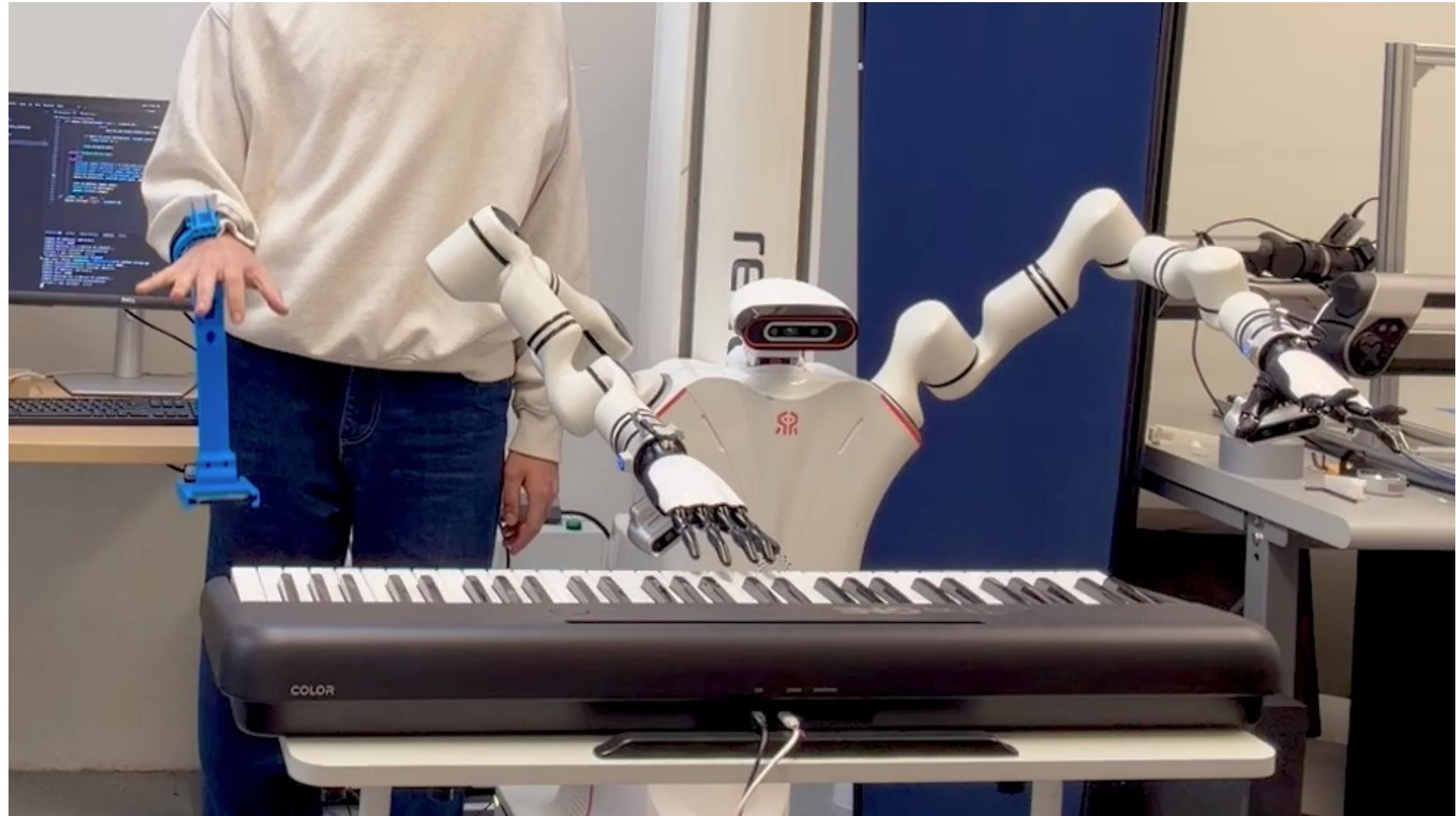
Ravichandar, Harish, et al. "Recent advances in robot learning from demonstration."  
*Annual review of control, robotics, and autonomous systems* 3.1 (2020): 297-330.

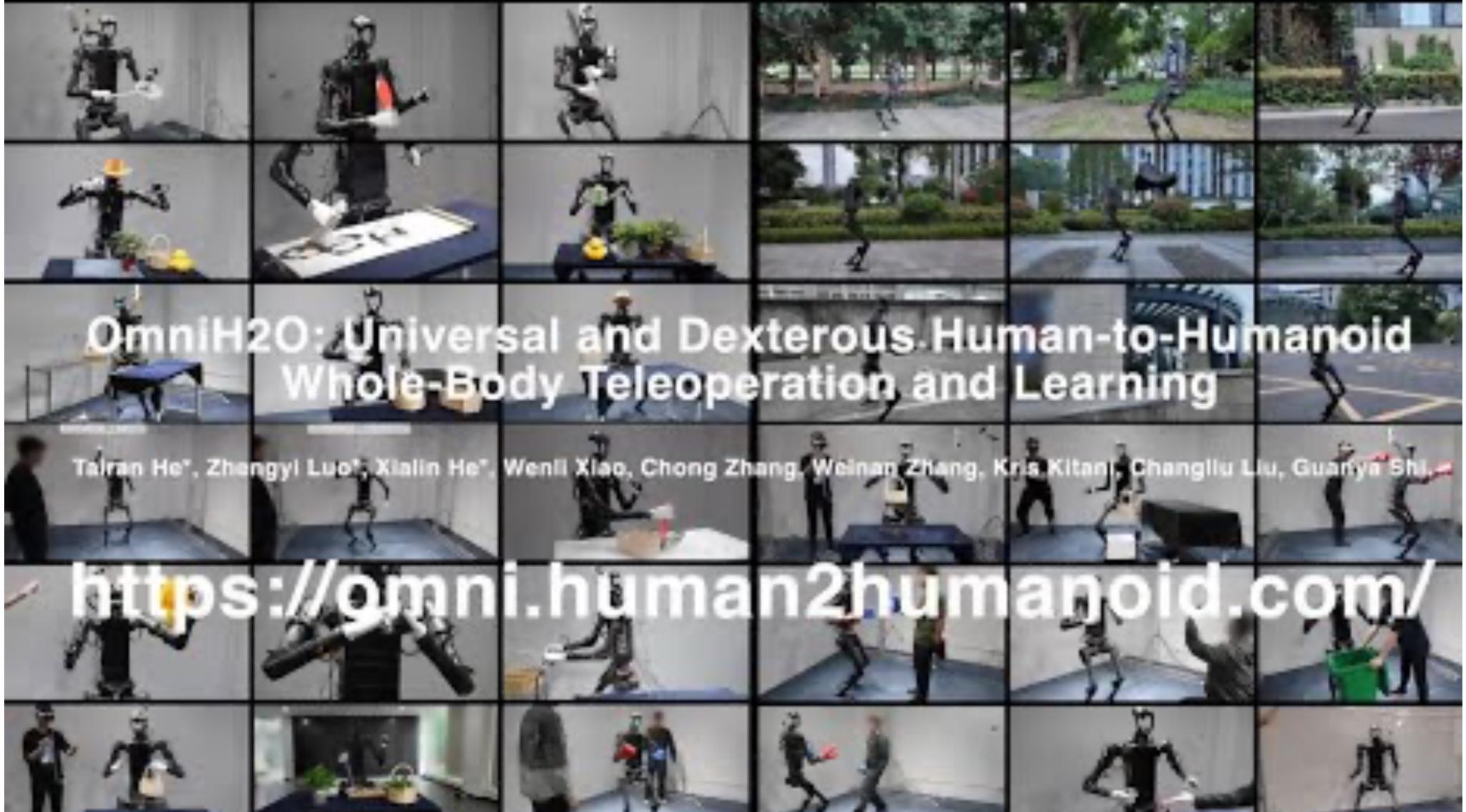


ntional Center







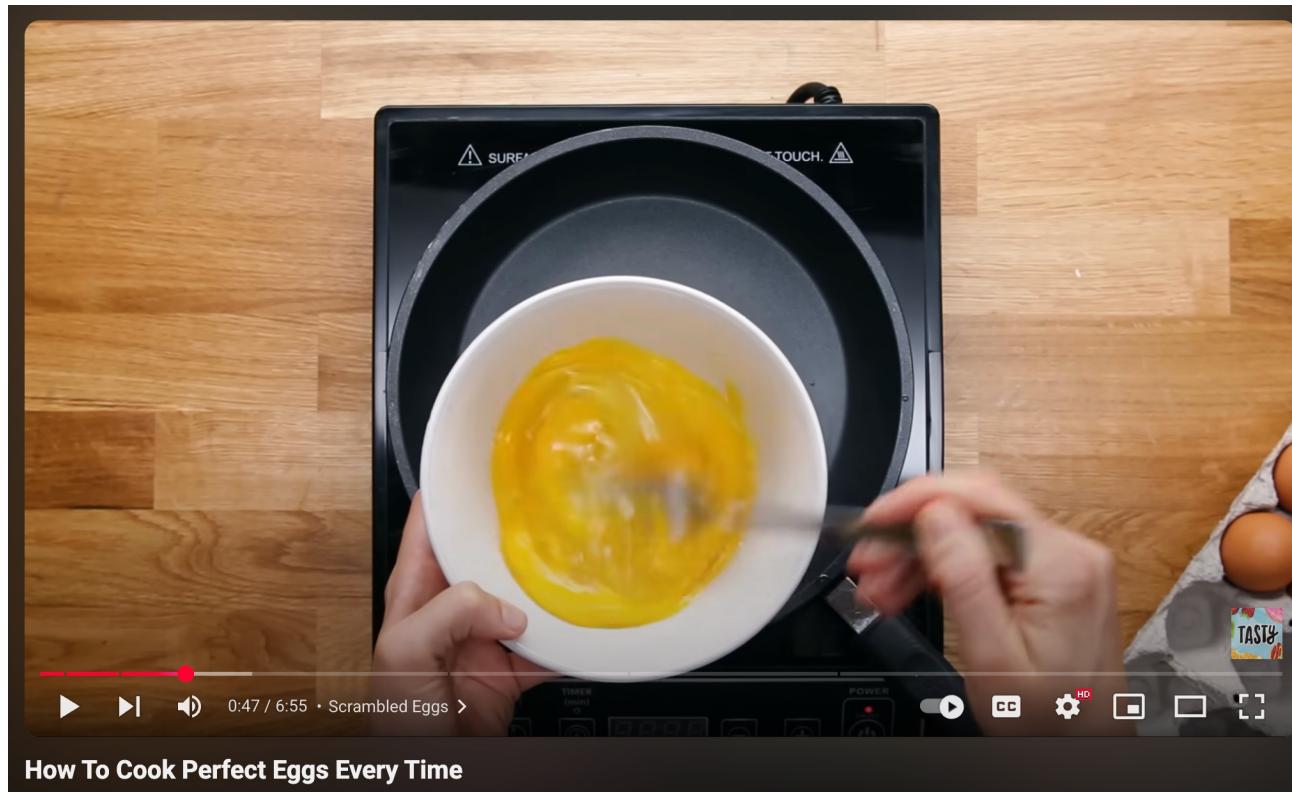


## OmniH2O: Universal and Dexterous Human-to-Humanoid Whole-Body Teleoperation and Learning

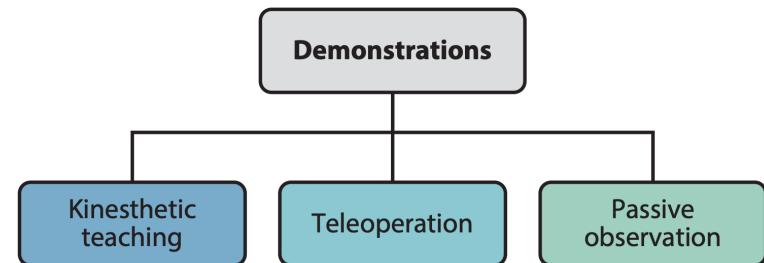
Taiyan He\*, Zhengyi Luo\*, Xialin He\*, Wenli Xiao, Chong Zhang, Weinan Zhang, Kris Kitani, Changlu Liu, Guanya Shi, Ming Tang

<https://omni.human2humanoid.com/>

# Learning from Observations



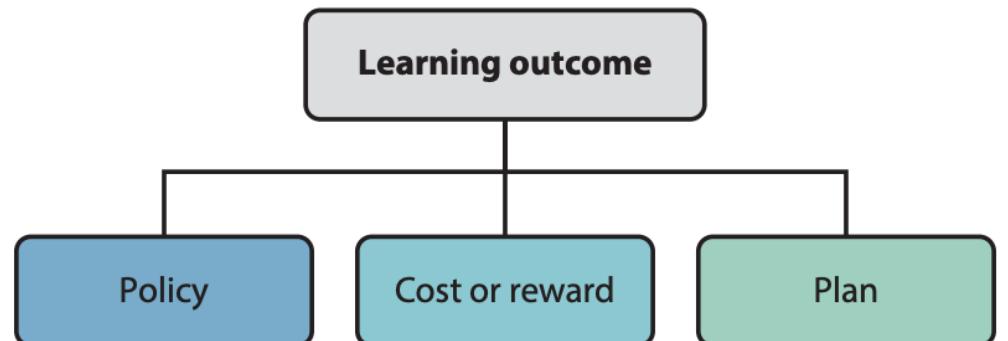
# Types of Demonstrations



Demonstration	Ease of demonstration	High DOFs	Ease of mapping
Kinesthetic teaching	✓		✓
Teleoperation		✓	✓
Passive observation	✓	✓	

Ravichandar, Harish, et al. "Recent advances in robot learning from demonstration." *Annual review of control, robotics, and autonomous systems* 3.1 (2020): 297-330.

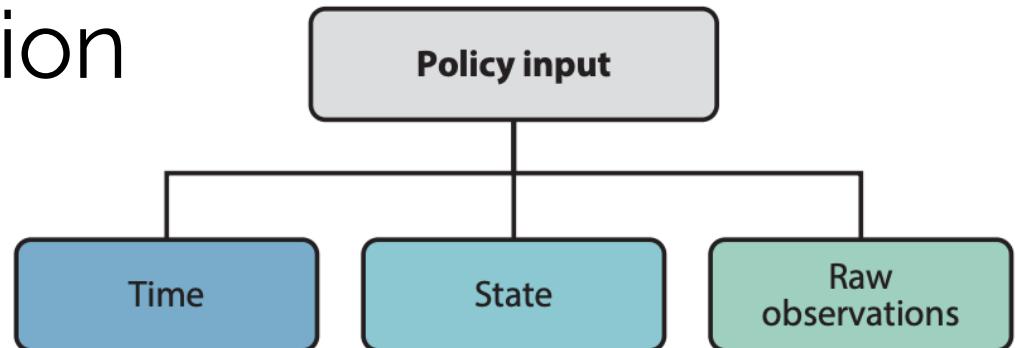
# Learning Outcomes



Learning outcome	Low-level control	Action space continuity	Compact representation	Long-horizon planning	Multistep tasks
Policy	✓	✓	✓		
Cost or reward	✓	✓		✓	
Plan			✓	✓	✓

Ravichandar, Harish, et al. "Recent advances in robot learning from demonstration." *Annual review of control, robotics, and autonomous systems* 3.1 (2020): 297-330.

# Policy Parameterization



Policy input	Ease of design	Performance guarantees	Robustness to perturbations	Task variety	Algorithmic efficiency
Time	✓	✓			✓
State		✓	✓		✓
Raw observations	✓		✓	✓	

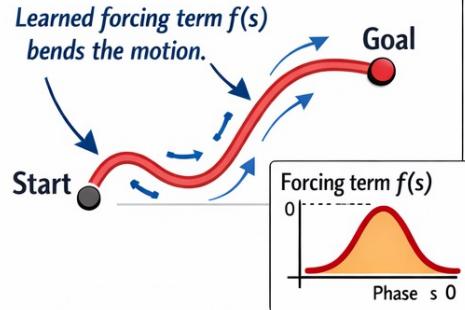
Ravichandar, Harish, et al. "Recent advances in robot learning from demonstration." *Annual review of control, robotics, and autonomous systems* 3.1 (2020): 297-330.

# Dynamic Movement Primitives

# Dynamic Movement Primitives

## learnable dynamical systems for generating movements

stable, spring-like motion models augmented with a learned shaping term that lets you reproduce and flexibly adapt demonstrated movements to new goals and speeds

Basic Spring-Damper System	DMP with Learned Forcing Term	Flexible Goals & Speed
 <p>Pure spring–damper: always goes straight to the goal.</p>	 <p>Learned forcing term <math>f(s)</math> bends the motion.</p> <p>Forcing term <math>f(s)</math></p> <p>Phase <math>s = 0</math></p> <p>Spring–damper + learned forcing term → Same goal, shaped motion.</p>	 <p>Start</p> <p>Goal A</p> <p>Goal B</p> <p>Fast</p> <p>Same learned shape, flexible goal &amp; speed.</p>

# Dynamic Movement Primitives

$$\begin{aligned}\tau \dot{v} &= K(\underline{g} - x) - Dv + (\underline{g} - \underline{x}_0) f \\ \tau \dot{x} &= v\end{aligned}$$

goal                          goal    initial state  
 K: spring constant  
 D: damping term

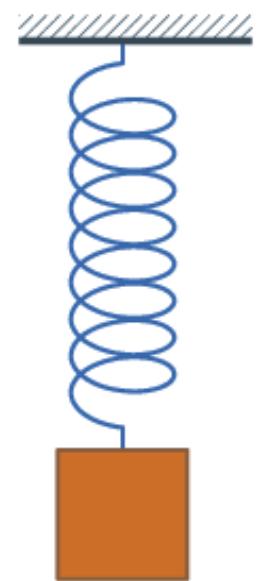
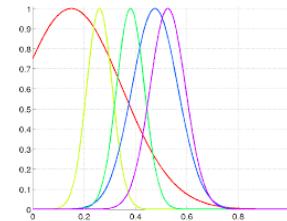
Non-linear force function:

$$f(s) = \frac{\sum_i w_i \psi_i(s)s}{\sum_i \psi_i(s)} \quad \psi_i(s) = \exp(-h_i(s - c_i)^2)$$

Gaussian basis functions

canonical system:  $\tau \dot{s} = -\alpha s$

s: phase variable



Pastor, Peter, et al. "Learning and generalization of motor skills by learning from demonstration." *2009 IEEE international conference on robotics and automation*. IEEE, 2009.

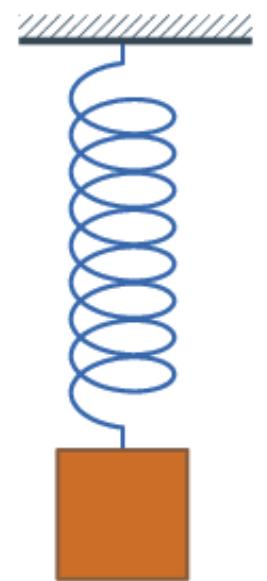
# Dynamic Movement Primitives

Learning:

$$f_{\text{target}}(s) = \frac{-K(g - x) + Dv + \tau\dot{v}}{g - x_0} \quad f(s) = \frac{\sum_i w_i \psi_i(s)s}{\sum_i \psi_i(s)}$$

$$J = \sum_s (f_{\text{target}}(s) - f(s))^2$$

Linear regression



Pastor, Peter, et al. "Learning and generalization of motor skills by learning from demonstration." *2009 IEEE international conference on robotics and automation*. IEEE, 2009.

# Characteristics of DMPs

- Convergence to the goal  $g$  is guaranteed (for bounded weights) since  $f(s)$  vanishes at the end of a movement.
- The weights  $w_i$  can be learned to generate any desired *smooth* trajectory.
- The equations are spatial and temporal invariant, i.e., movements are self-similar for a change in goal, start point, and temporal scaling without a need to change the weights  $w_i$ .
- The formulation generates movements which are robust against perturbation due to the inherent attractor dynamics of the equations.

Pastor, Peter, et al. "Learning and generalization of motor skills by learning from demonstration." *2009 IEEE international conference on robotics and automation*. IEEE, 2009.

# Summary

- DMP enable learning “movement styles” while enabling generalization to new movement targets
- DMP is a purely kinematic account
- DMP addresses timing, but account of coordination is limited
- DMP for different tasks and their combination...?

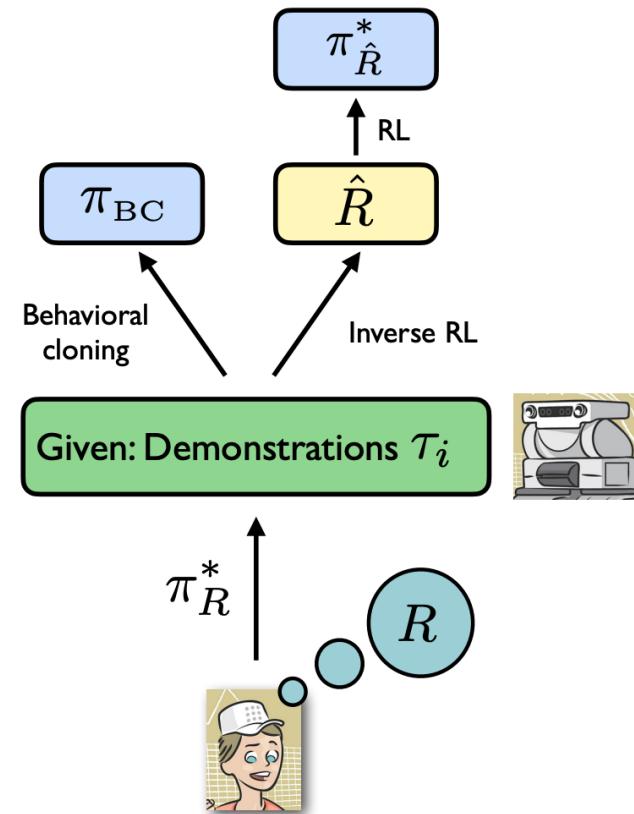
credit: Gregor Schöner



© BBC/Rosie Thomas

How to imitate complex behavior?

# Imitation Learning



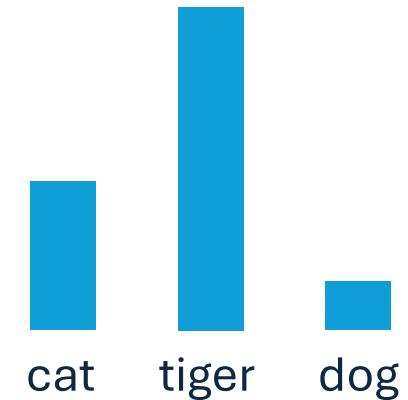
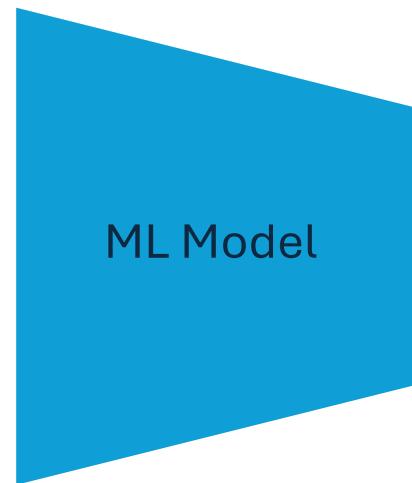
slide credit: Scott Niekum

Berkeley CS285: [https://www.youtube.com/watch?v=tbLaFtYpWWU&list=PL\\_iWQOsF6TfVYGEGiAOMaOzzv41Jfm\\_Ps&index=4](https://www.youtube.com/watch?v=tbLaFtYpWWU&list=PL_iWQOsF6TfVYGEGiAOMaOzzv41Jfm_Ps&index=4)

# Supervised Learning

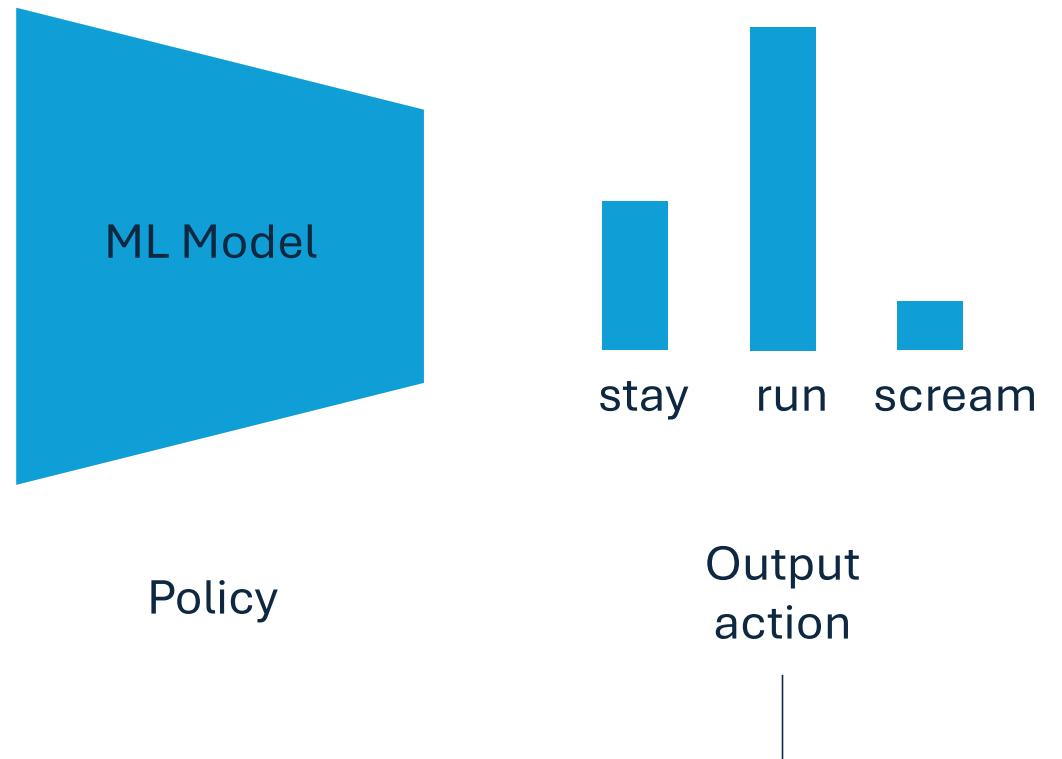


Input  
image



Output  
class label

# Supervised Learning / Behavior Cloning



**NOT i.i.d. -> independent and identically distributed**

# The i.i.d. Assumption

“The training and testing data are **independent** and **identically** distributed.”

# The i.i.d. Assumption

**Training**



**Testing**



# Input Data Distribution

“Poodle”



“Chihuahua”



“Shar-Pei”

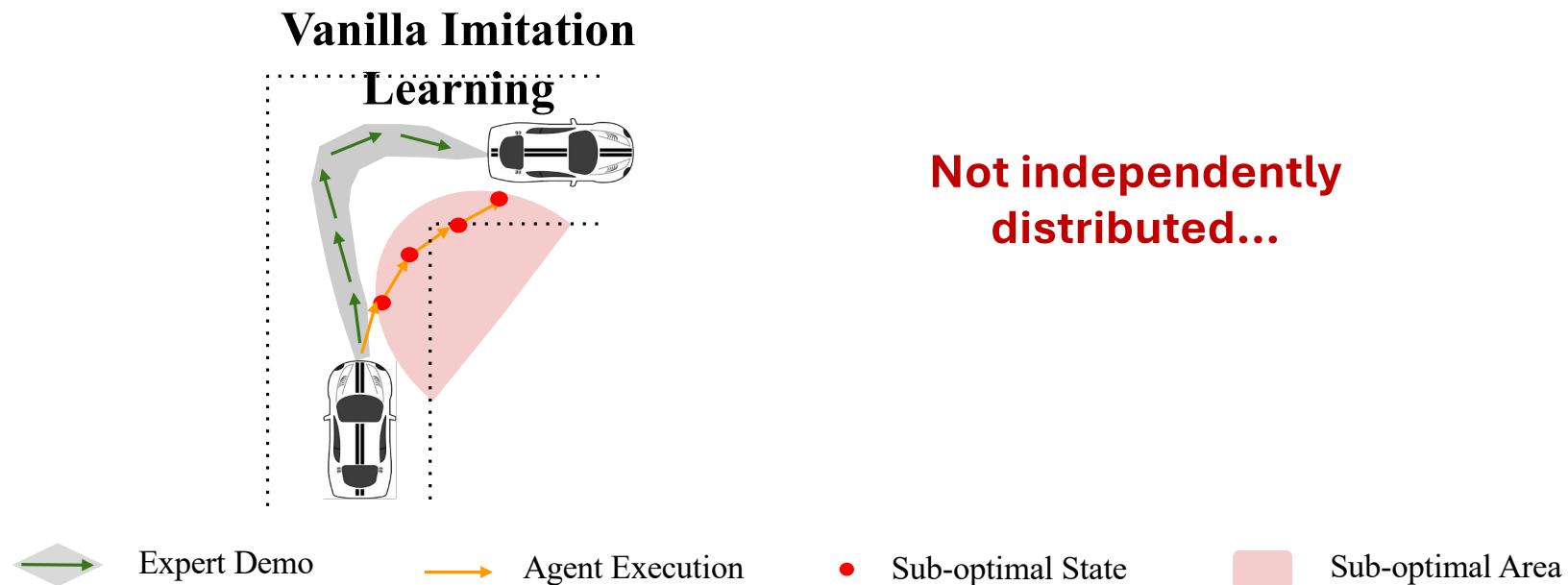


**Not identically distributed...**

**Robustness** refers to the ability of a system, model, or method to maintain performance or produce reliable results **despite variations, noise, errors, or adversarial conditions** in the input or environment.

# Input Data Distribution

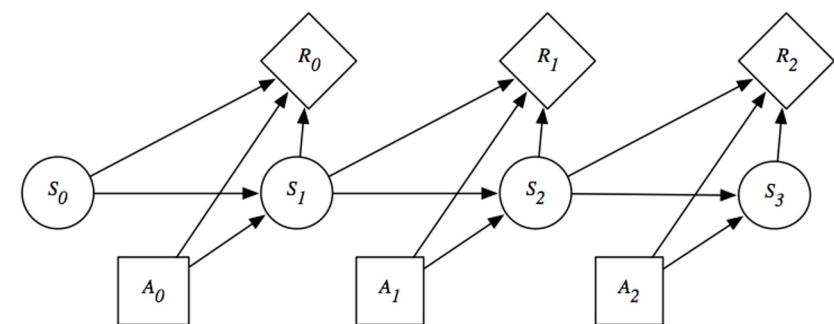
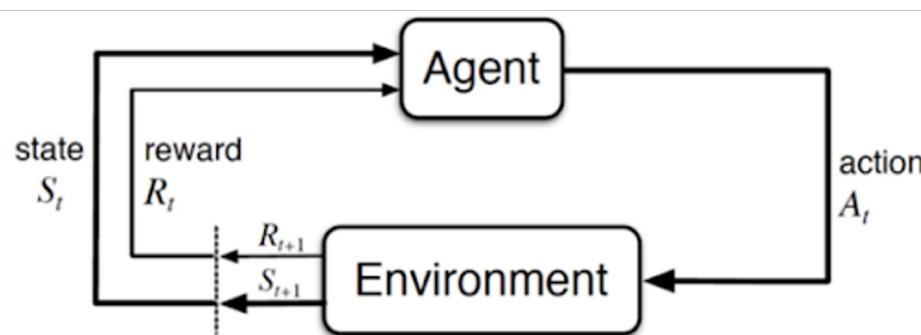
End-to-End Control Tasks



Berkeley CS285: [https://www.youtube.com/watch?v=tbLaFtYpWWU&list=PL\\_iWQOsF6TfVYGEGiAOMaOzzv41Jfm\\_Ps&index=4](https://www.youtube.com/watch?v=tbLaFtYpWWU&list=PL_iWQOsF6TfVYGEGiAOMaOzzv41Jfm_Ps&index=4)

# Formalizing Sequential Decision Making

Markov Decision Process  $\langle S, A, P, R \rangle$



Berkeley CS285: [https://www.youtube.com/watch?v=tbLaFtYpWWU&list=PL\\_iWQOsE6TfVYGEGiAOMaOzzv41Jfm\\_Ps&index=4](https://www.youtube.com/watch?v=tbLaFtYpWWU&list=PL_iWQOsE6TfVYGEGiAOMaOzzv41Jfm_Ps&index=4)

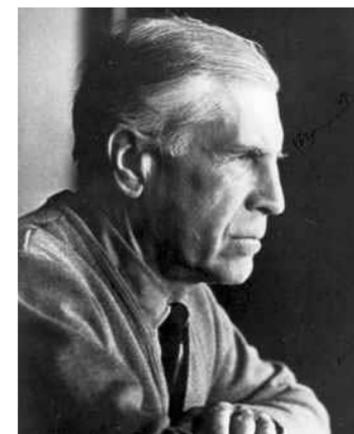
$s_t$  – state  
 $a_t$  – action



Richard Bellman

American applied mathematician  
introduced **dynamic programming** in 1953

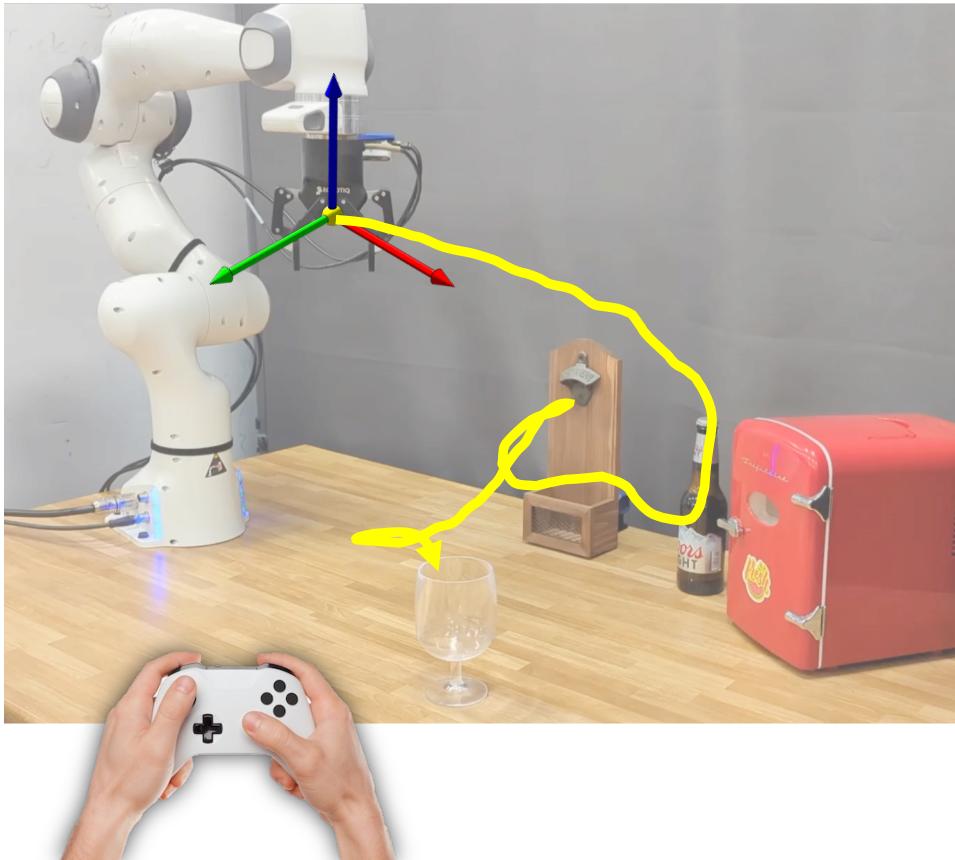
$x_t$  – state  
 $u_t$  – action      управление



Lev Pontryagin

Soviet mathematician  
algebraic topology, differential topology and **optimal control**

# Behavioral Cloning



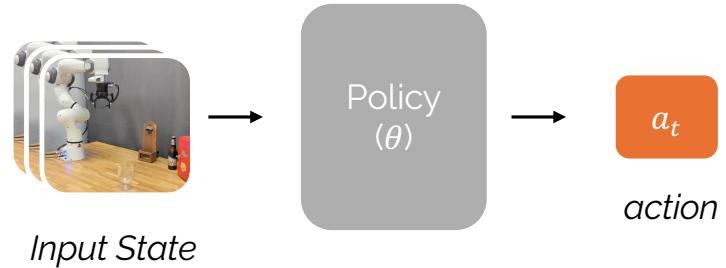
$s = (\text{Third-person View}, \text{Wrist-camera View}, \text{proprio.})$

$a = (\Delta x, \Delta y, \Delta z, \Delta rx, \Delta ry, \Delta rz, \text{close\_gripper})$

Gripper Velocity

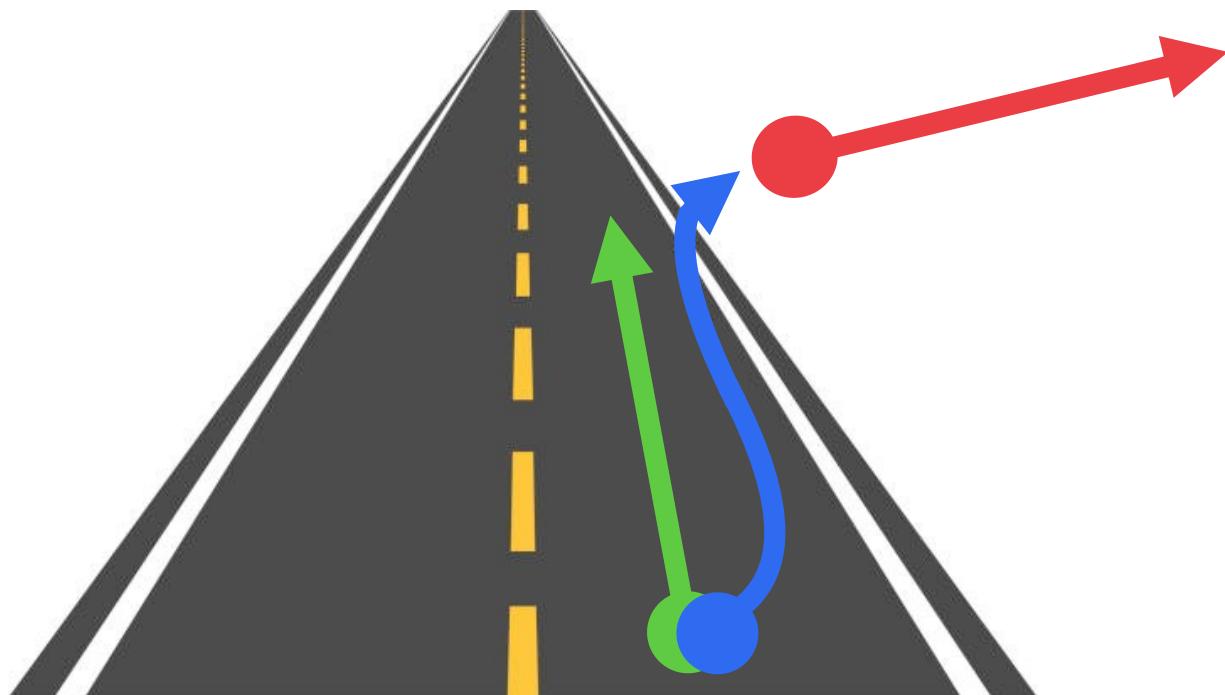
Grasping?

$$D = \{(s, a)\}$$



$$\mathcal{L}(\theta) = -\mathbb{E}_{(s,a) \sim D} [\log \pi_\theta(a|s)]$$

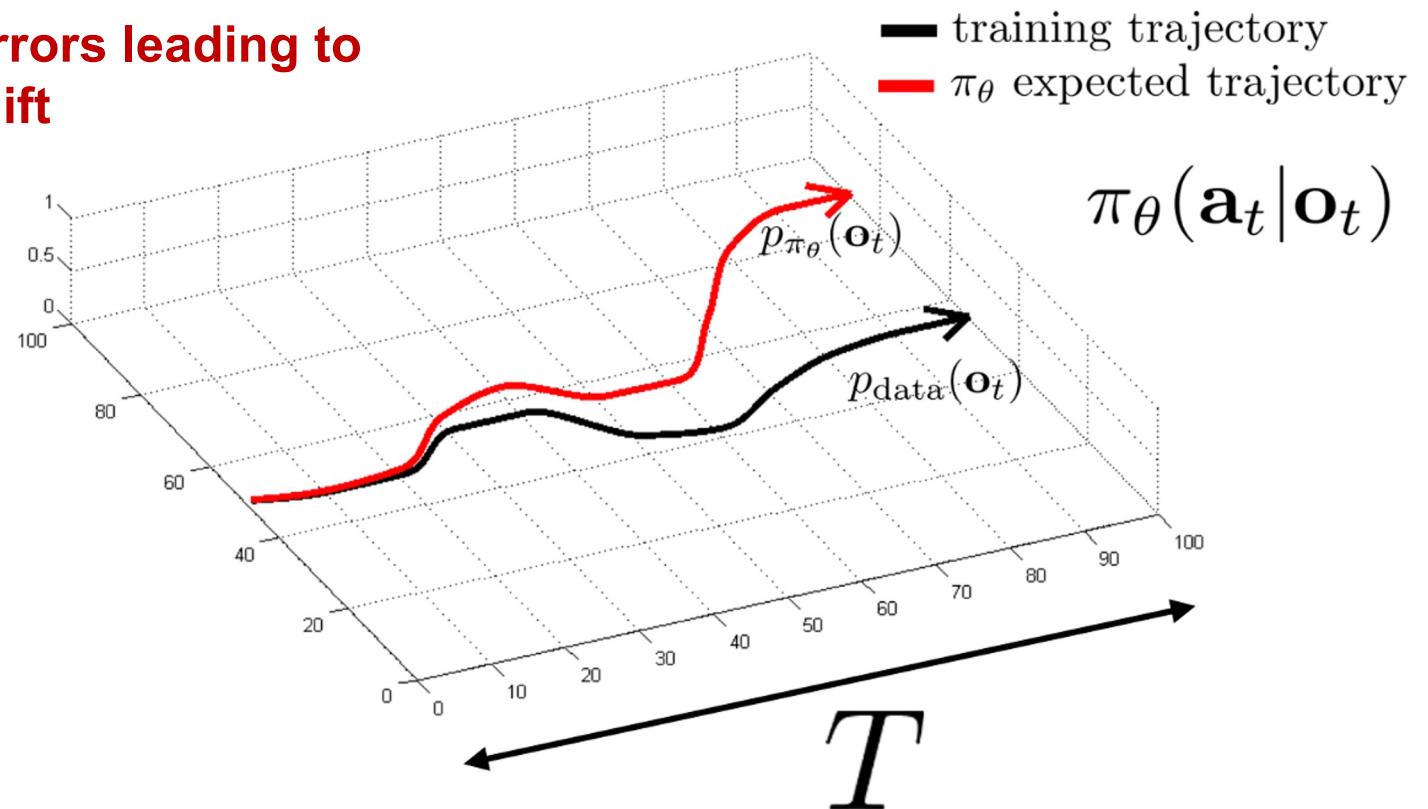
# Downsides of behavioral cloning



slide credit: Scott Niekum

# Challenges in imitation learning #1

Compounding Errors leading to  
Distributional Shift



# Quadratic Regret

**Regret** (in decision theory) measures the difference between the reward (or outcome) one actually received and the best possible reward one *could have received* if one had made the optimal choice

$$\hat{\pi}_{sup} = \arg \min_{\pi \in \Pi} \mathbb{E}_{s \sim d_{\pi^*}} [\ell(s, \pi)]$$

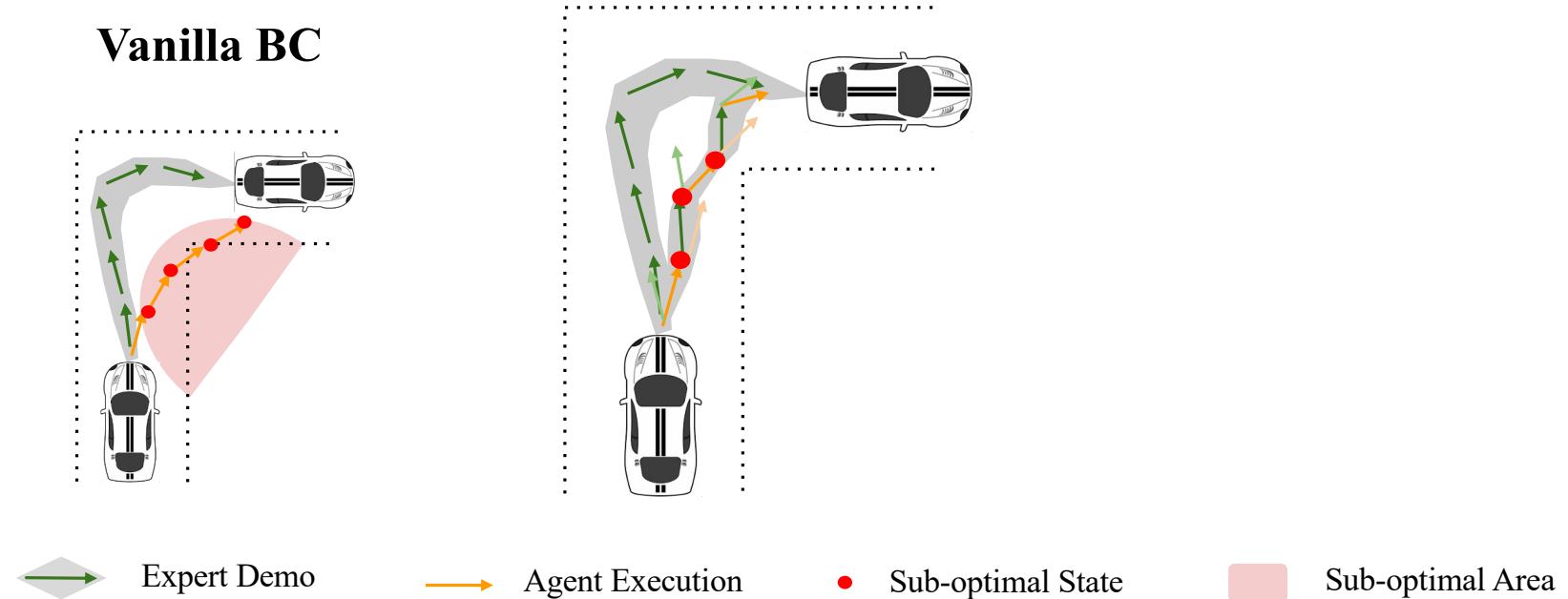
Assuming  $\ell(s, \pi)$  is the 0-1 loss (or upper bound on the 0-1 loss) implies the following performance guarantee with respect to any task cost function  $C$  bounded in  $[0, 1]$ :

**Theorem 2.1.** (Ross and Bagnell, 2010) Let  $\mathbb{E}_{s \sim d_{\pi^*}} [\ell(s, \pi)] = \epsilon$ , then  $J(\pi) \leq J(\pi^*) + T^2\epsilon$ .

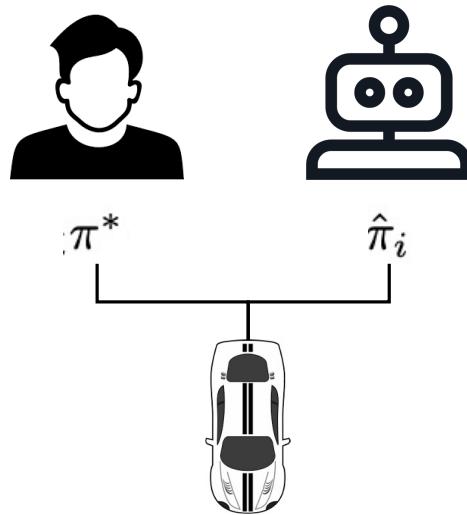
Compare to typical supervised learning loss that grows as:  $O(\epsilon T)$

# DAgger: Dataset Aggregation

End-to-End Control Tasks



# Dagger



```
Initialize  $\mathcal{D} \leftarrow \emptyset$ .  
Initialize  $\hat{\pi}_1$  to any policy in  $\Pi$ .  
for  $i = 1$  to  $N$  do  
    Let  $\pi_i = \beta_i \pi^* + (1 - \beta_i) \hat{\pi}_i$ .  
    Sample  $T$ -step trajectories using  $\pi_i$ .  
    Get dataset  $\mathcal{D}_i = \{(s, \pi^*(s))\}$  of visited states by  $\pi_i$   
        and actions given by expert.  
    Aggregate datasets:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_i$ .  
    Train classifier  $\hat{\pi}_{i+1}$  on  $\mathcal{D}$ .  
end for  
Return best  $\hat{\pi}_i$  on validation.
```

**Algorithm 3.1:** DAGGER Algorithm.

**Key idea:** keep collecting demonstration data that is on-distribution for current policy,  
and reduce dependence on expert over time

slide credit: Scott Niekum

Ross, Stéphane, Geoffrey Gordon, and Drew Bagnell. "A reduction of imitation learning and structured prediction to no-regret online learning." *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011.

# DAgger

**Theorem 2.2.** Let  $\pi$  be such that  $\mathbb{E}_{s \sim d_\pi}[\ell(s, \pi)] = \epsilon$ , and  $Q_{T-t+1}^{\pi^*}(s, a) - Q_{T-t+1}^{\pi^*}(s, \pi^*) \leq u$  for all action  $a$ ,  $t \in \{1, 2, \dots, T\}$ ,  $d_\pi^t(s) > 0$ , then  $J(\pi) \leq J(\pi^*) + uT\epsilon$ .

If difference between optimal t-step Q and any other action is u  
(e.g. the worst single action regret):

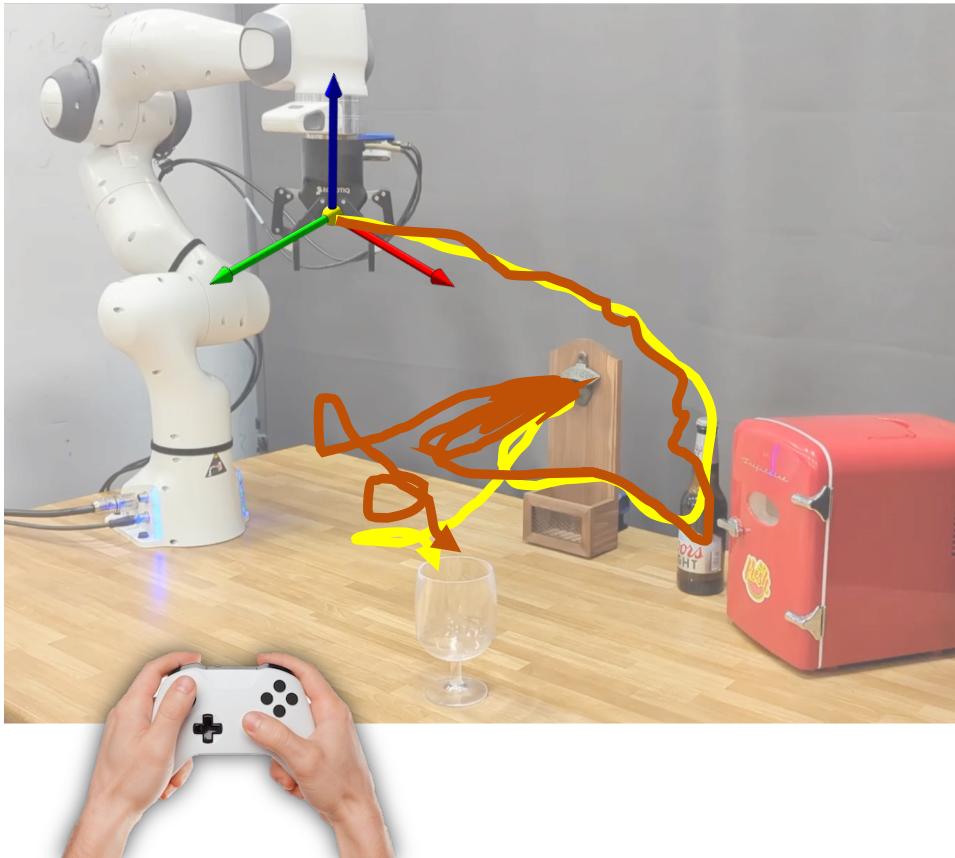
The end cost is (no) worse than optimal plus number of mistakes  
times u, the worst possible regret of each mistake

## What are some other ways to improve robustness?

**Robustness** refers to the ability of a system, model, or method to maintain performance or produce reliable results **despite variations, noise, errors, or adversarial conditions** in the input or environment.

$$J(\pi) \leq J(\pi^*) + \boxed{T^2} \epsilon.$$

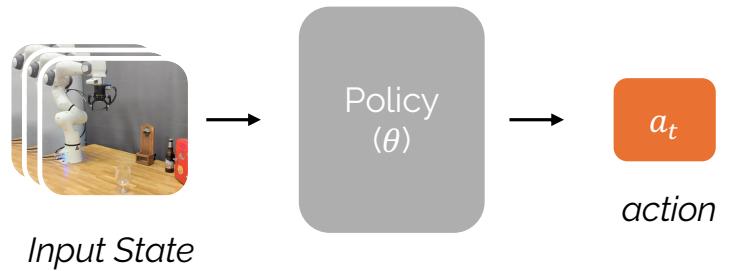
# Behavioral Cloning



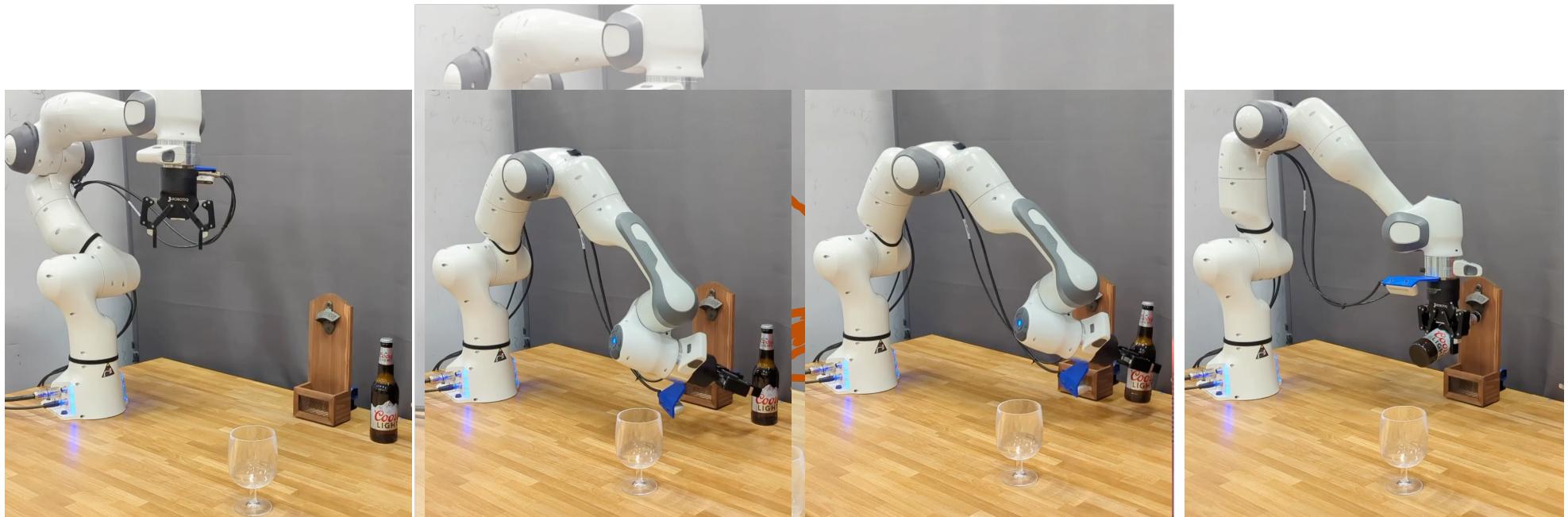
$s = (\text{Third-person View}, \text{Wrist-camera View}, \text{proprio.})$

$a = (\Delta x, \Delta y, \Delta z, \Delta rx, \Delta ry, \Delta rz, \text{close\_gripper})$   
Gripper Velocity      Grasping?

$$D = \{(s, a)\}$$



$$\mathcal{L}(\theta) = -\mathbb{E}_{(s,a) \sim D} [\log \pi_\theta(a|s)]$$



*reach bottle*

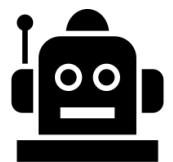
*pick up bottle*

*move bottle to opener*

*open bottle*



$a = (\Delta x, \Delta y, \Delta z, \Delta rx, \Delta ry, \Delta rz, gripper)$   
 $a_0 a_1 a_2 \dots$       ...  $a_k \dots$       ...  $aT$



**Unintended low-level motions** constitute noise in demonstrations.

The demonstrator's **high-level actions** are optimal!

Free-space Reaching



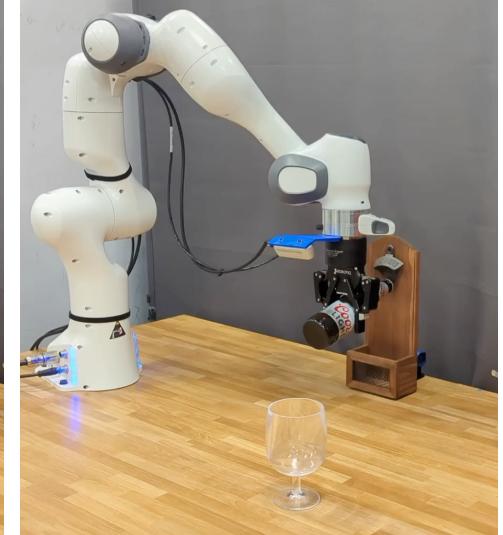
Contact-Rich Interaction



Free-space Reaching



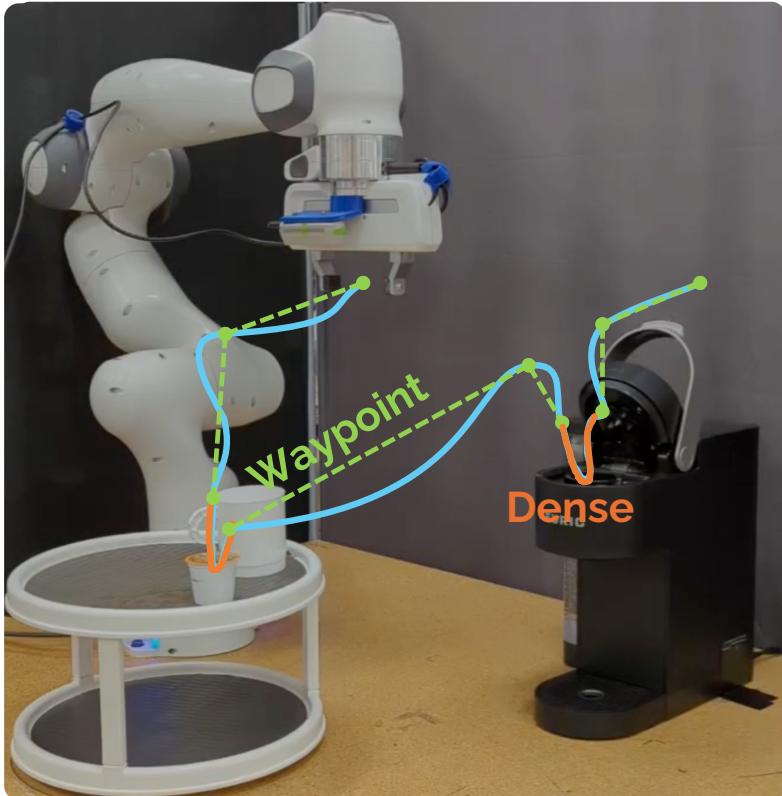
Contact-Rich Interaction



The demonstrator's **high-level actions** are optimal!

These actions can be categorized into ***two general modes***.

# HYDRA: Hybrid Robot Actions for Imitation Learning in Manipulation



$$D = \{(s, a, w, m)\}$$

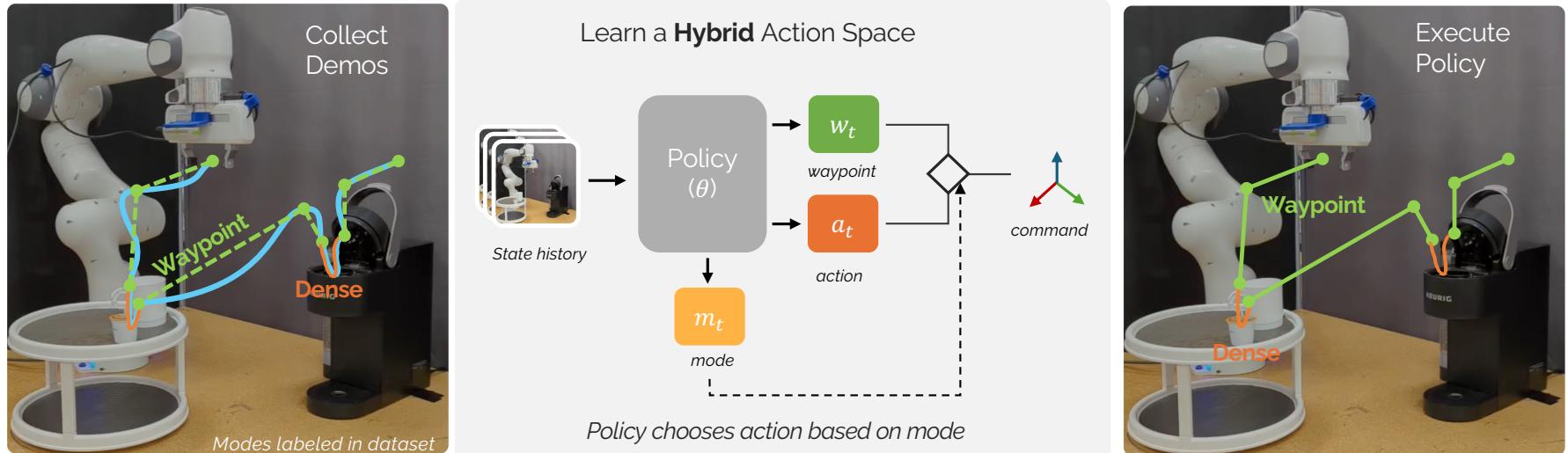
$$s = (RGB_{external}, RGB_{wrist}, proprio.)$$

$$a = (\Delta x, \Delta y, \Delta z, \Delta rx, \Delta ry, \Delta rz, close\_gripper)$$

$$w = (x, y, z, rx, ry, rz, gripper\_state)$$

$$m = 0 \text{ if waypoint, else } 1$$

# HYDRA: Hybrid Robot Actions for Imitation Learning in Manipulation



$$\mathcal{L}_a(\theta) = -\mathbb{E}_{(s,a,w,m) \sim D} [(1 - \alpha) \log \pi_\theta^A(a|s) + \alpha \log \pi_\theta^W(w|s)]$$

$$\mathcal{L}_m(\theta) = -\mathbb{E}_{(s,a,w,m) \sim D} [(m \log \pi_\theta^M(m=1|s) + (1 - m) \log \pi_\theta^M(m=0|s))]$$



## Make Coffee



■ BC-RNN

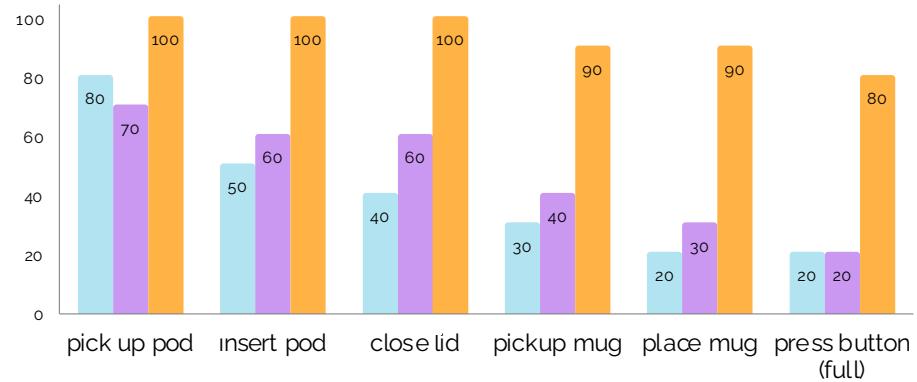
Mandlekar et al. What Matters in Learning from Offline Human Demonstrations for Robot Manipulation. CoRL'21

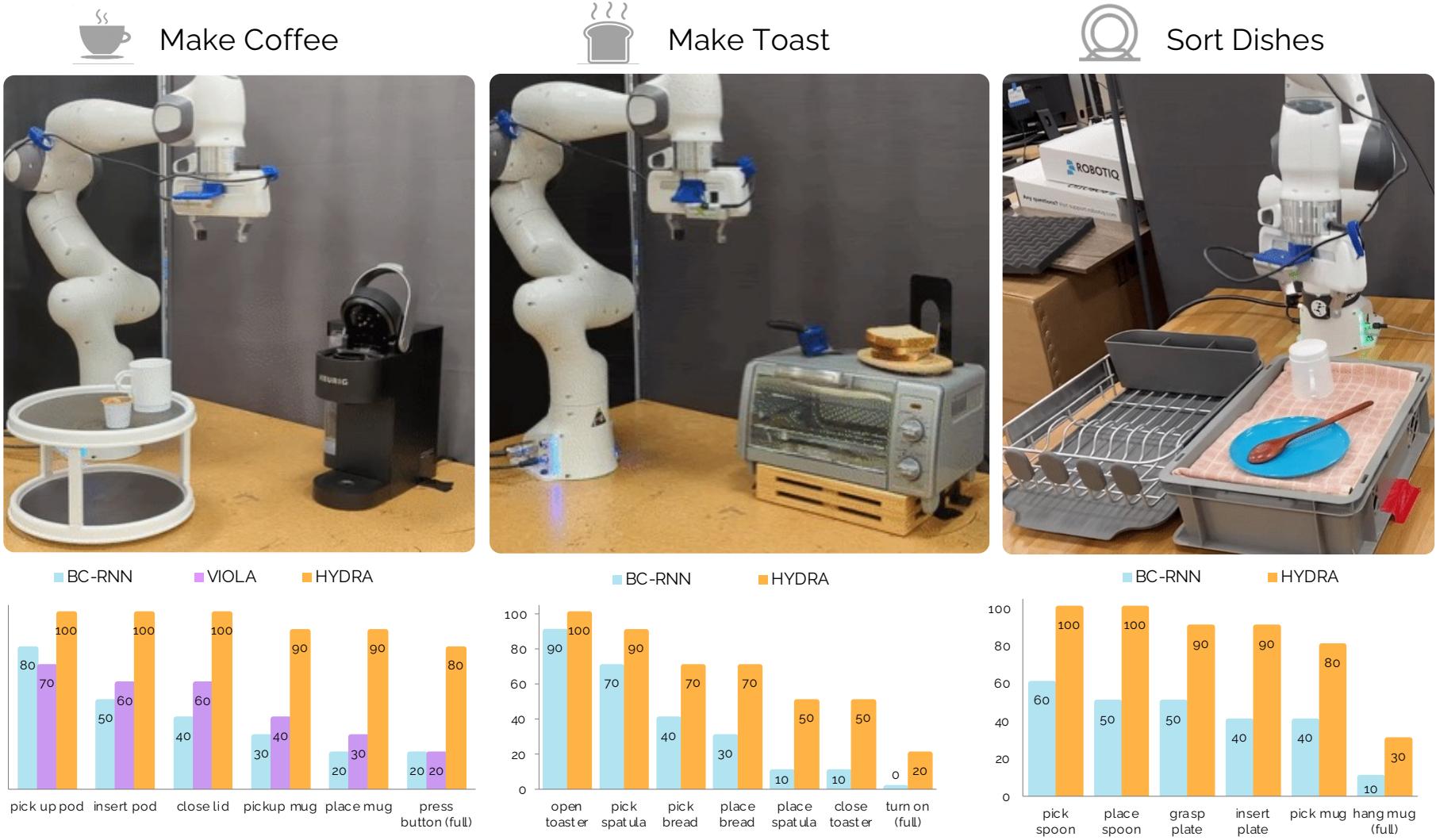
■ VIOLA

Zhu et al. "VIOLA: Object-Centric Imitation Learning for Vision-Based Robot Manipulation." CoRL'22

■ HYDRA

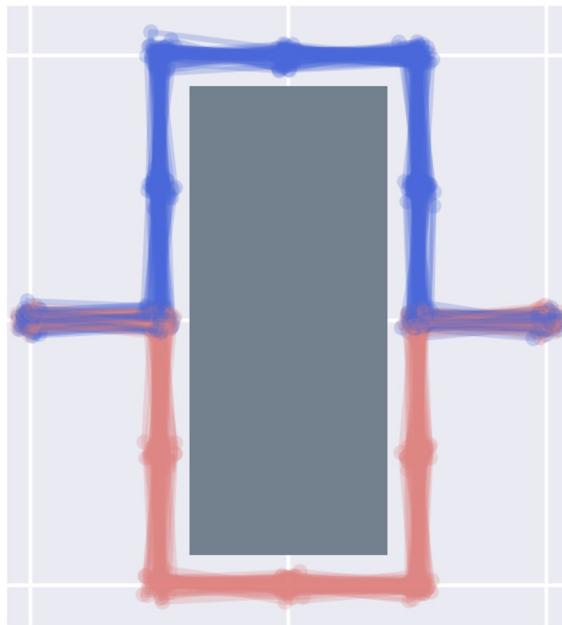
OURS



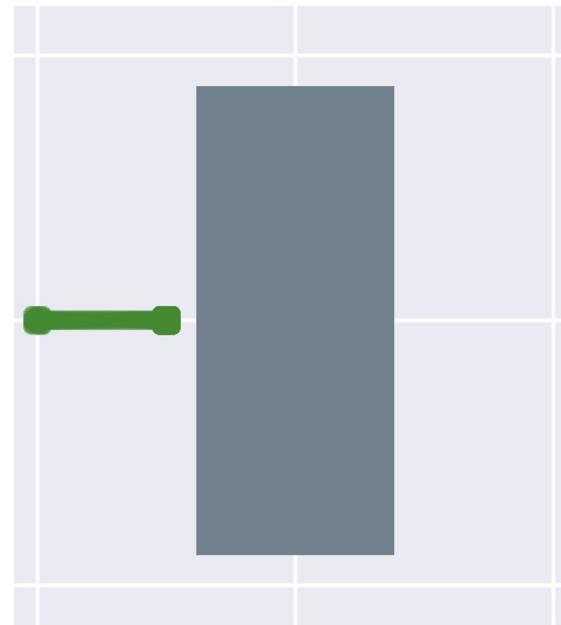


# Challenges in imitation learning #2

Dataset



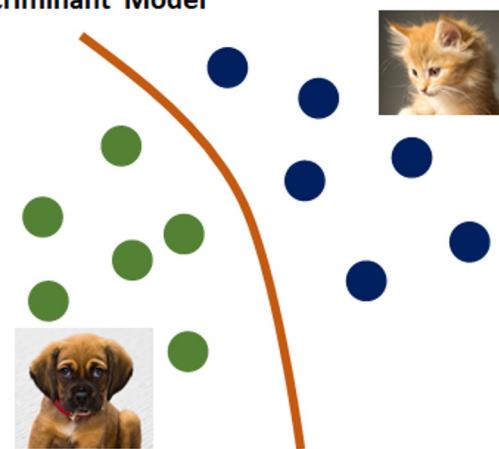
Unimodal BC



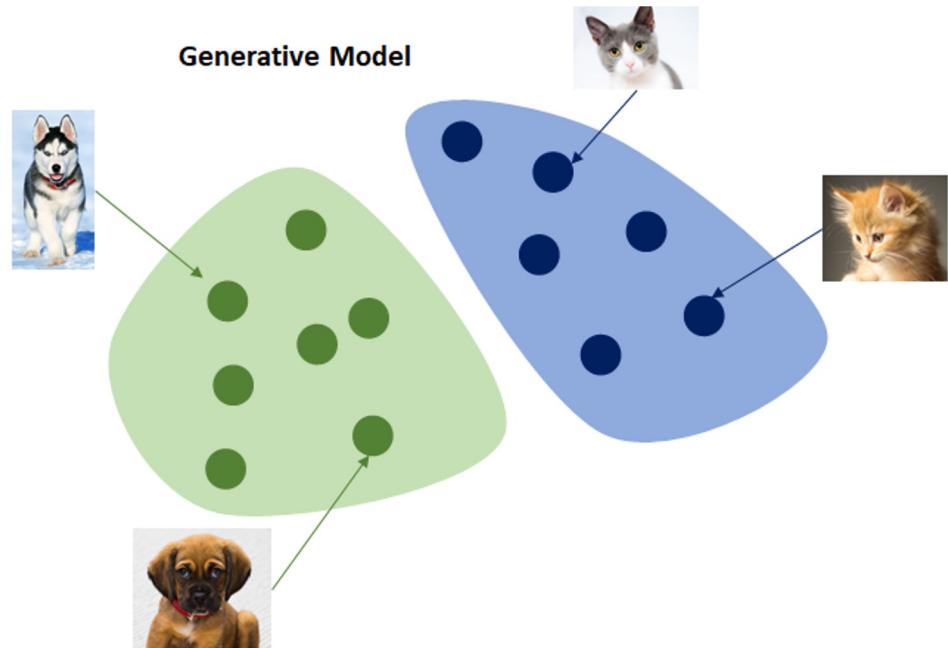
Multimodal demonstrations

# Generative Modeling

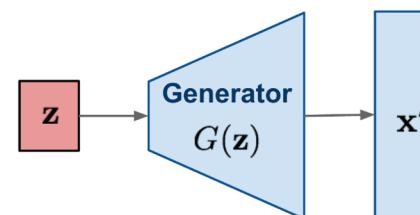
Discriminant Model



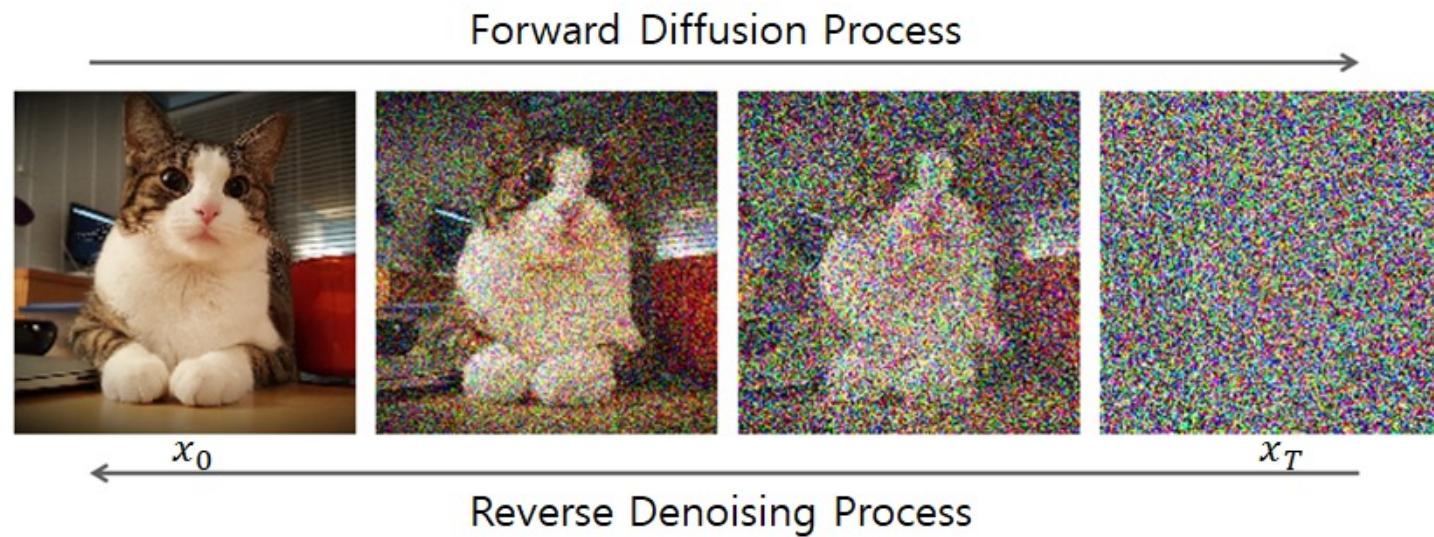
Generative Model



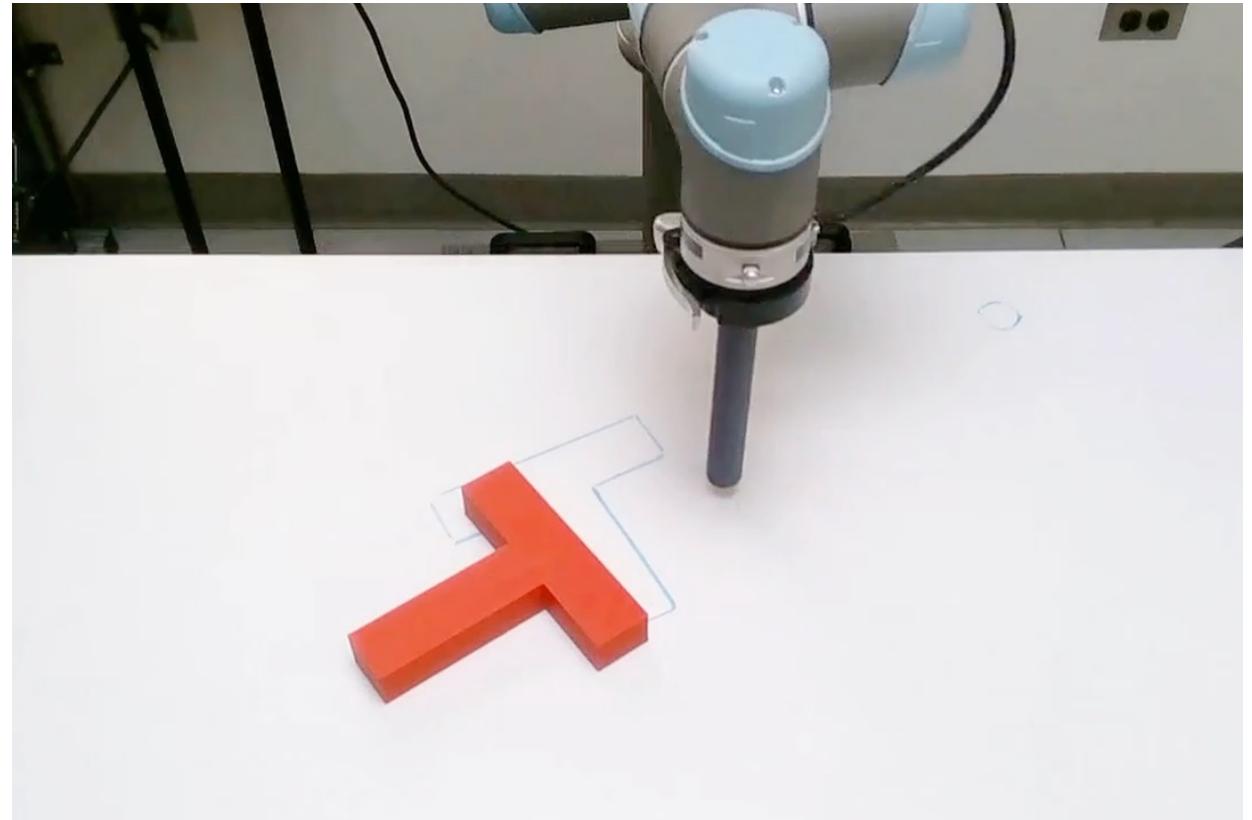
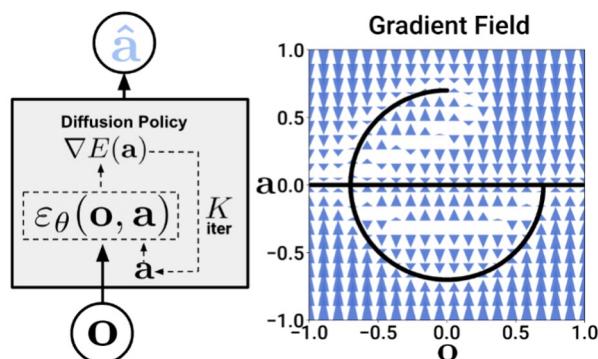
# Generative Models



# Diffusion



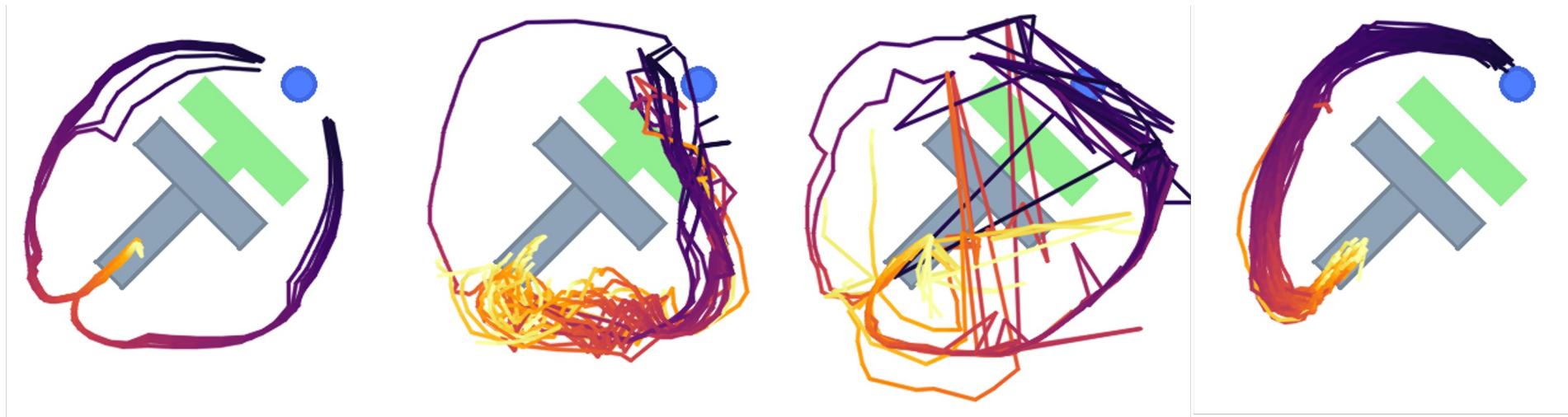
# Diffusion Policy



<https://diffusion-policy.cs.columbia.edu/>

Chi, Cheng, et al. "Diffusion policy: Visuomotor policy learning via action diffusion." *The International Journal of Robotics Research* (2023): 02783649241273668.

# Diffusion Policy



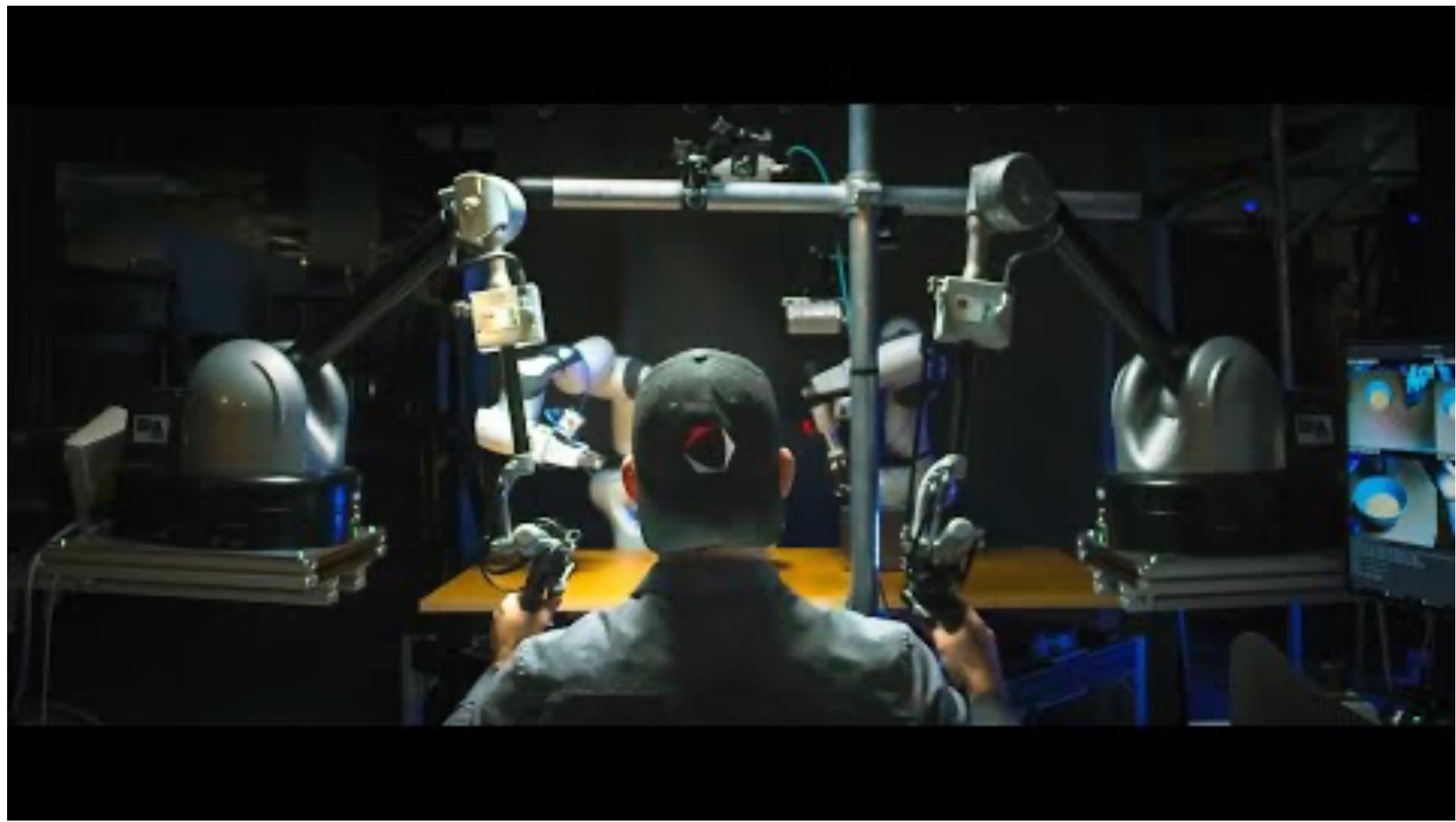
Diffusion Policy

LSTM-GMM

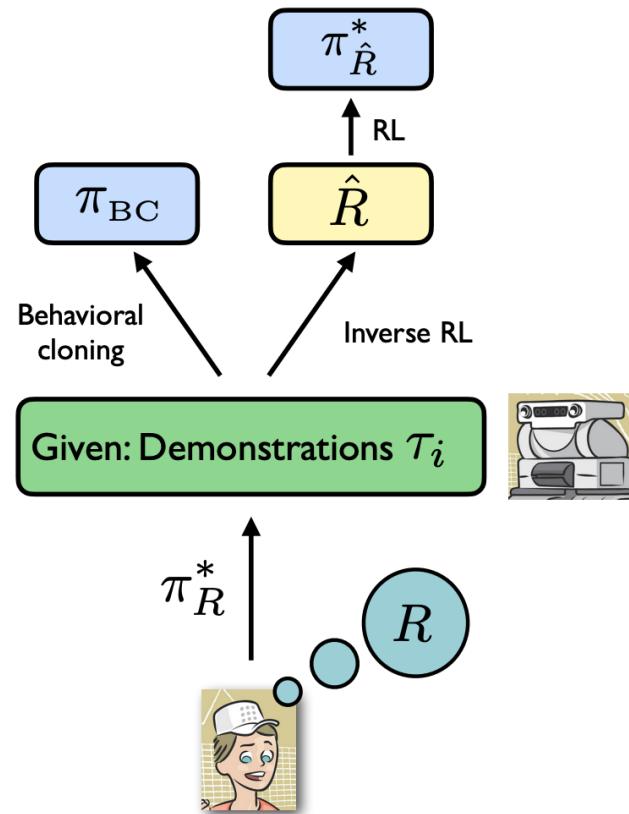
BET

IBC

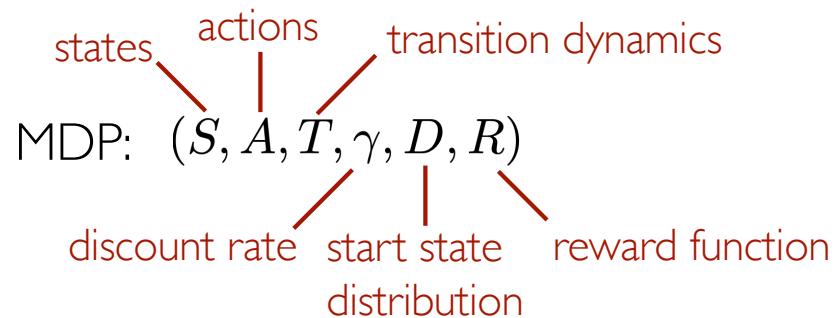
Chi, Cheng, et al. "Diffusion policy: Visuomotor policy learning via action diffusion." *The International Journal of Robotics Research* (2023): 02783649241273668.



## Imitation Learning



# Inverse reinforcement learning



Policy:  $\pi(s, a) \rightarrow [0, 1]$

Value function:  $V^\pi(s_0) = \sum_{t=0}^{\infty} \gamma^t R(s_t)$

What if we have an **MDP/R**?

# Inverse reinforcement learning

1. Collect user demonstration  $(s_0, a_0), (s_1, a_1), \dots, (s_n, a_n)$   
and assume it is sampled from the expert's policy,  $\pi^E$
2. Explain expert demos by finding  $R^*$  such that:

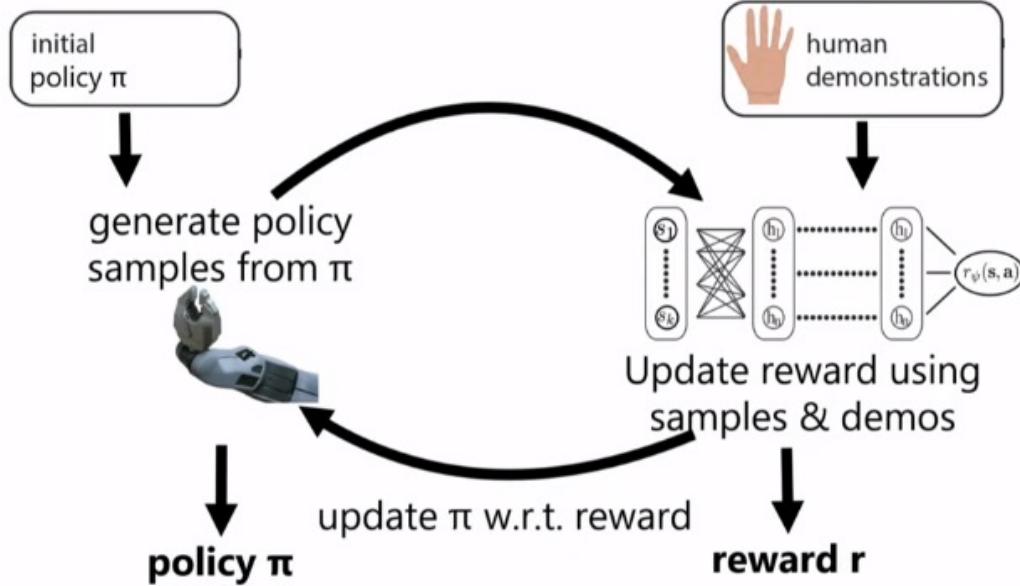
$$E[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi^E] \geq E[\sum_{t=0}^{\infty} \gamma^t R^*(s_t) | \pi] \quad \forall \pi$$

$$E_{s_0 \sim D}[V^{\pi^E}(s_0)] \geq E_{s_0 \sim D}[V^\pi(s_0)] \quad \forall \pi$$

slide credit: Scott Niekum

Abbeel, Pieter, and Andrew Y. Ng. "Apprenticeship learning via inverse reinforcement learning." *Proceedings of the twenty-first international conference on Machine learning*. 2004.

# Guided Cost Learning



Assume we don't know the dynamics, but we can sample, like in standard RL

---

## Algorithm 2 Nonlinear IOC with stochastic gradients

---

- 1: **for** iteration  $k = 1$  to  $K$  **do**
  - 2:   Sample demonstration batch  $\hat{\mathcal{D}}_{\text{demo}} \subset \mathcal{D}_{\text{demo}}$
  - 3:   Sample background batch  $\hat{\mathcal{D}}_{\text{samp}} \subset \mathcal{D}_{\text{samp}}$
  - 4:   Append demonstration batch to background batch:  
     $\hat{\mathcal{D}}_{\text{samp}} \leftarrow \hat{\mathcal{D}}_{\text{demo}} \cup \hat{\mathcal{D}}_{\text{samp}}$
  - 5:   Estimate  $\frac{d\mathcal{L}_{\text{IOC}}}{d\theta}(\theta)$  using  $\hat{\mathcal{D}}_{\text{demo}}$  and  $\hat{\mathcal{D}}_{\text{samp}}$
  - 6:   Update parameters  $\theta$  using gradient  $\frac{d\mathcal{L}_{\text{IOC}}}{d\theta}(\theta)$
  - 7: **end for**
  - 8: **return** optimized cost parameters  $\theta$
- 

slide credit: Sergey Levine

Finn, Chelsea, Sergey Levine, and Pieter Abbeel. "Guided cost learning: Deep inverse optimal control via policy optimization." *International conference on machine learning*. PMLR, 2016.

autonomous execution  
1x real-time

PR2

goal

our method  
100 samples from  $q(u_t | x_t)$



**That's it for today!**

**Questions?**