

Reinforcement Learning for Robot Control

Jens Kober

March 24, 2025

Workshop on Themes across Control and Reinforcement Learning

Robot Learning

Traditional Robot Programming

The image displays two side-by-side screenshots of a traditional robot programming environment. The left screenshot shows a code editor window with the following program code:

```
1 DEF DL_Cube_Pick_Basic()  
2 ;Programmer: Zabri@KUKA  
3  
4 INI  
5  
6 ; ROBOT : KR 16  
7  
8 $base=Base_data[1]  
9 $tool=Tool_data[1]  
10 PTP XNONE  
11 PTP XP4 C_DIS  
12 LIN XP5 C_DIS  
13 LIN XP2 C_DIS  
14 LIN XP1  
15 $OUT[2]= TRUE  
16 WAIT SEC 0.2  
17 LIN XP3 C_DIS  
18 LIN XP5 C_DIS  
19 PTP XP9 C_DIS  
20 LIN XP7 C_DIS  
21  
22  
23 END
```

The right screenshot shows a 3D simulation of a KUKA KR 16 robot arm in a virtual environment. The interface includes a toolbar with various navigation and tool options, a 'Workcell Diagram' on the left, and a 'KR 16 Properties' panel at the bottom. The robot is shown in a 3D view, and the simulation is running.

At the bottom of the image, there is a watermark: "This video may not be reproduced or disclosed to third parties without the express permission".

What's the problem?



CC BY Steve Jurvetson



CC BY-SA Mixabest



CC BY-NC-SA Brandon Heyer

What's the problem?

- Uncertainties and variations
 - Objects
 - Environment
 - Tasks
 - Human behavior
- Occurring in all robot domains!

(Reinforcement) Learning in Robotics

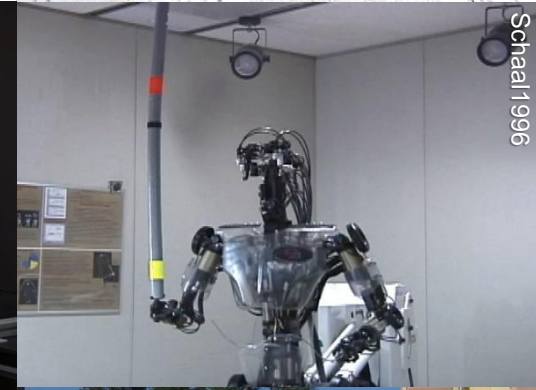
- Safe with real robots
- Fast learning
 - Sample efficient – “small” data
 - Incorporate prior knowledge
 - Few open parameters
- Real time computations



Google



Mahadevan1992



Schaal1996

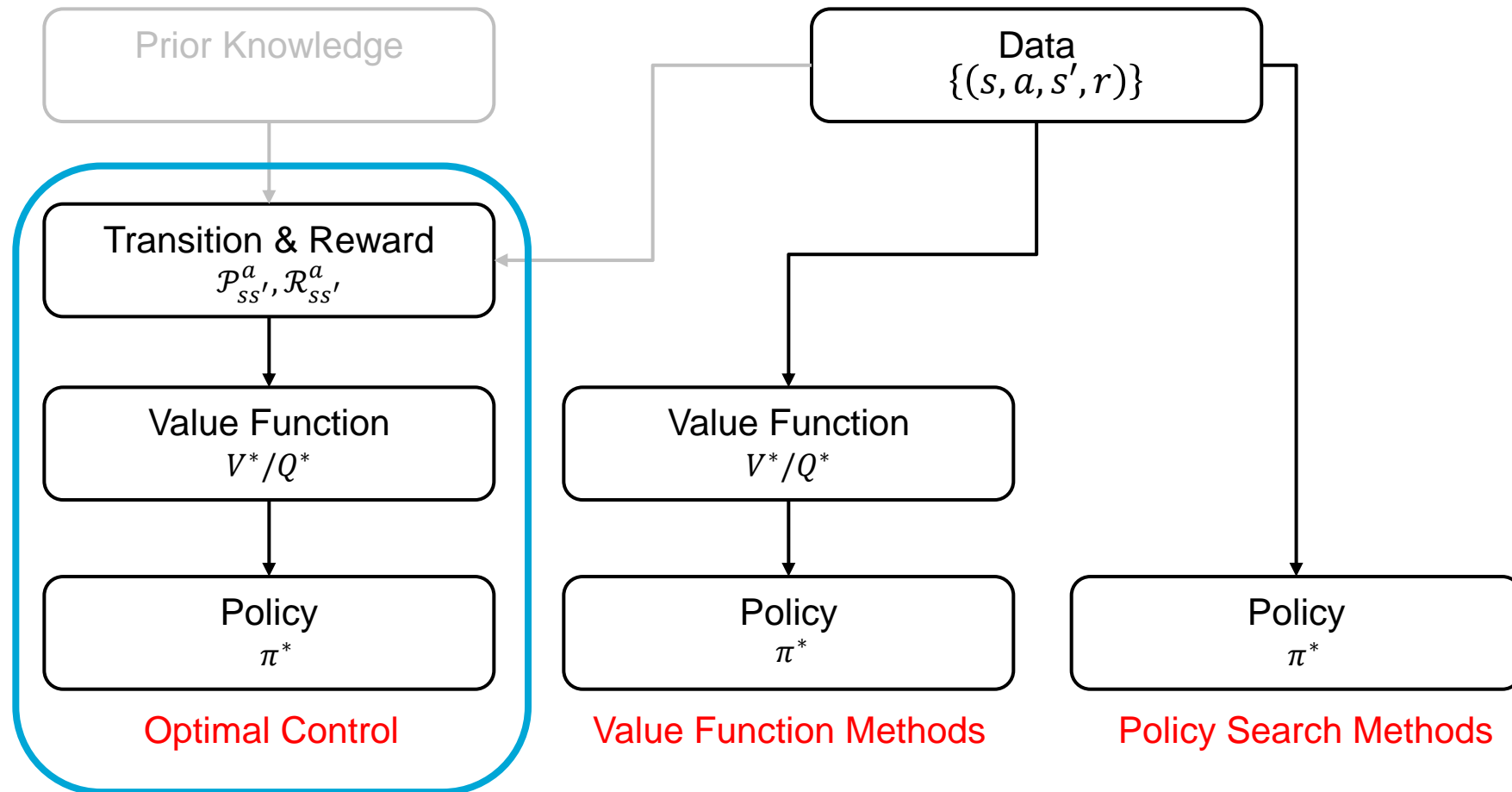


Rotmann2007



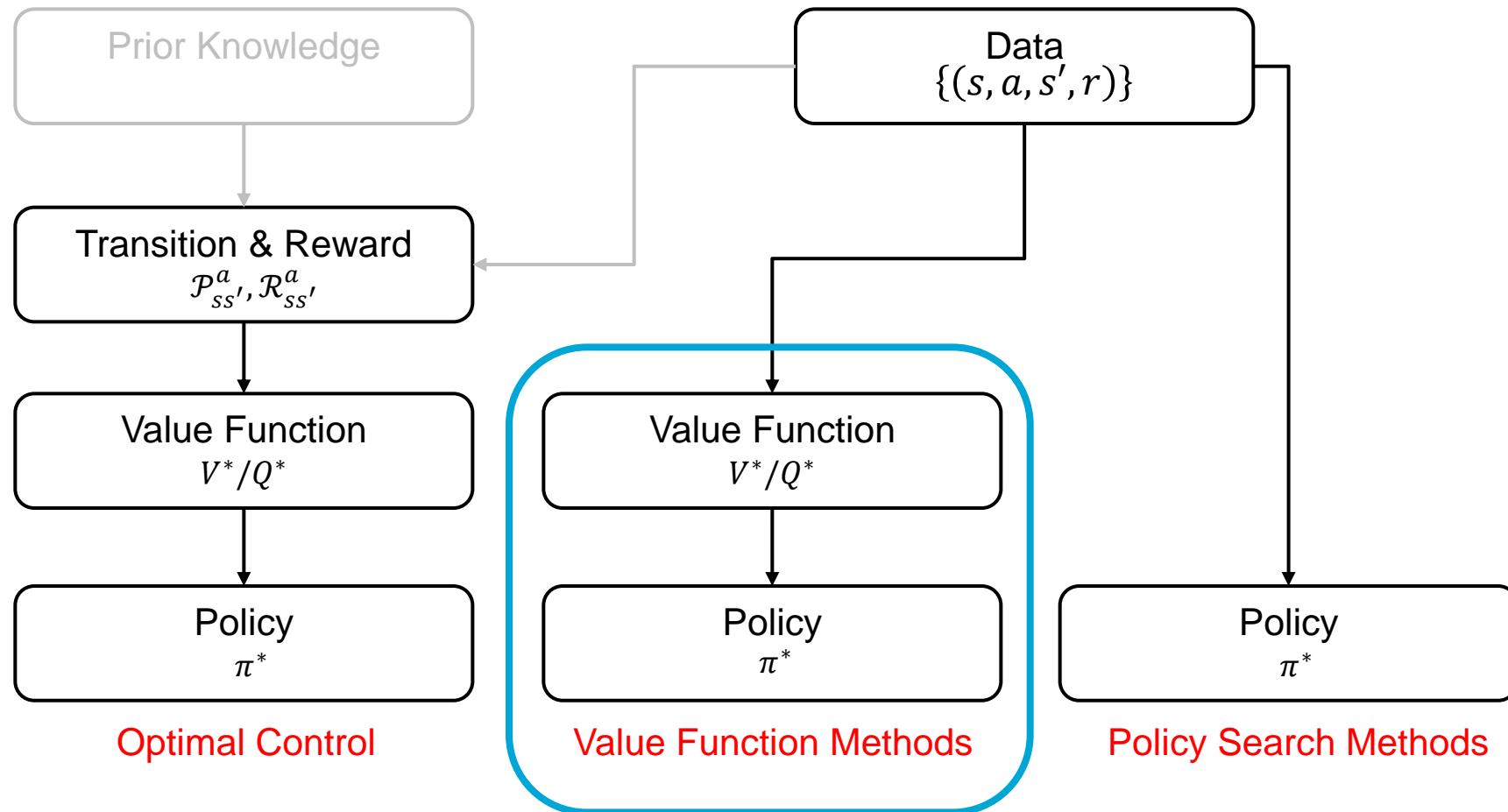
Kormushev2010

Types of RL Algorithms



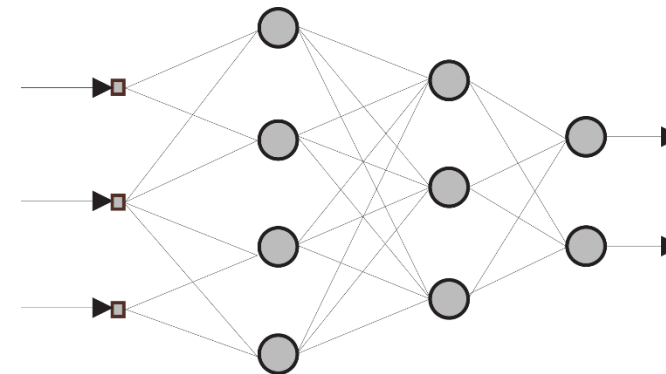
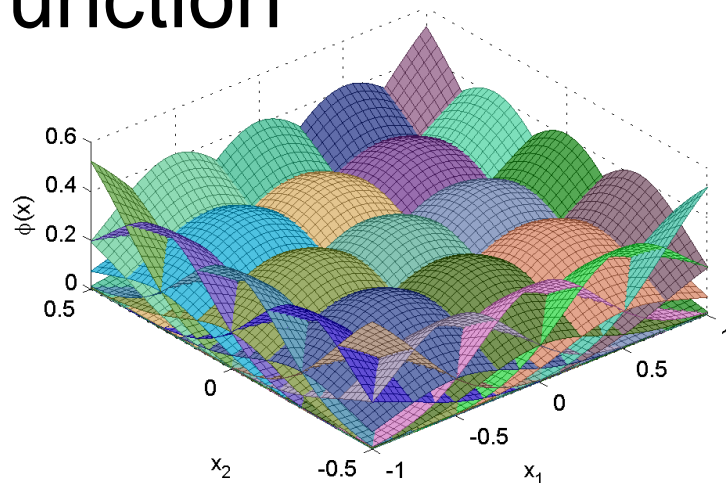
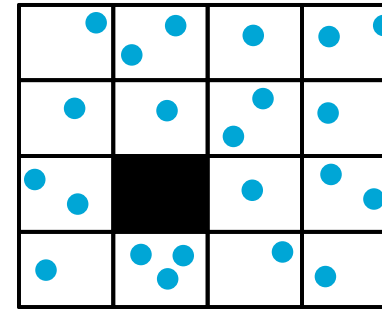


Types of RL Algorithms



Dealing with Continuous States/Actions

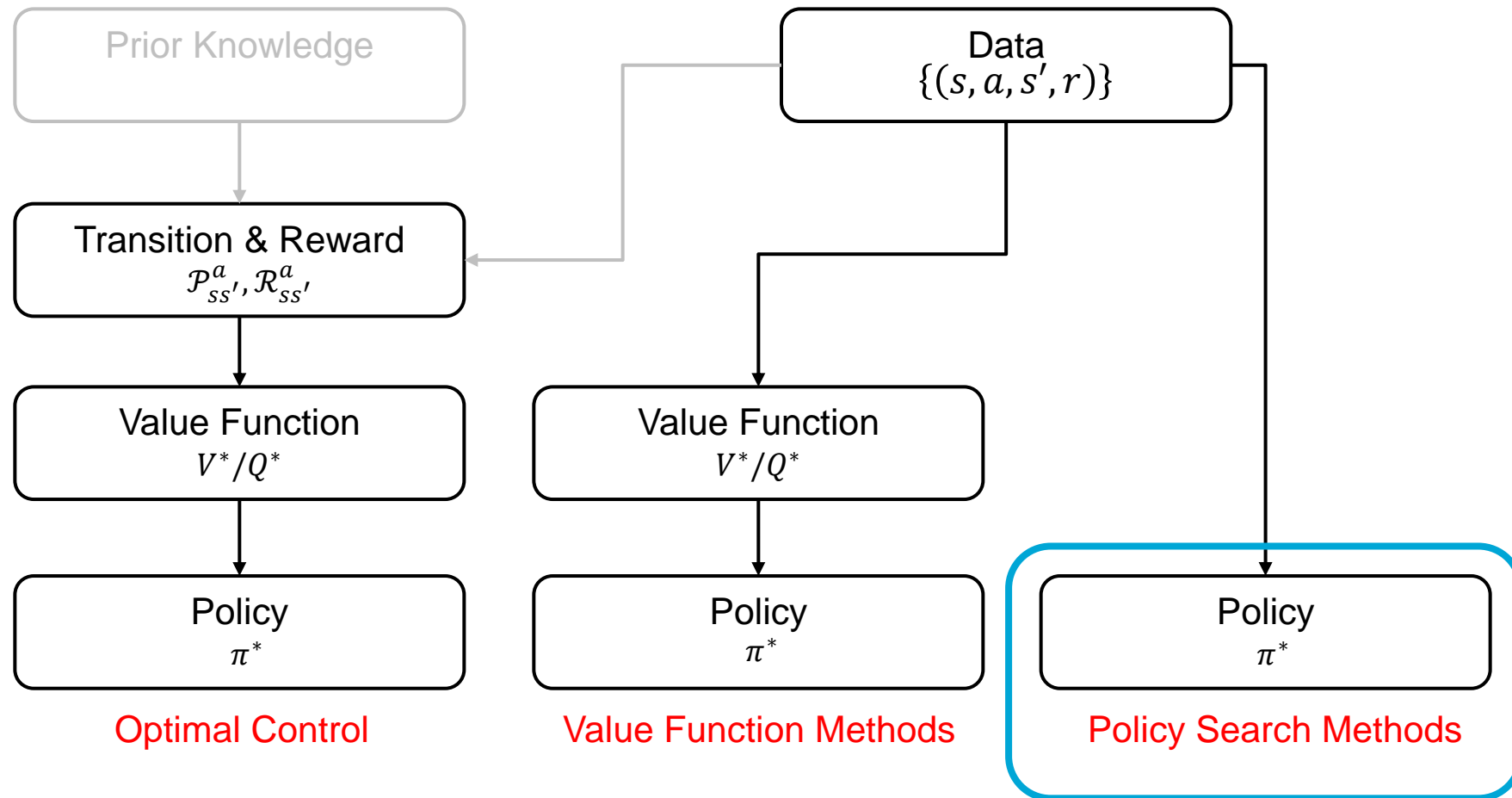
- Discretization
- Function approximation
 - Value function
 - Policy



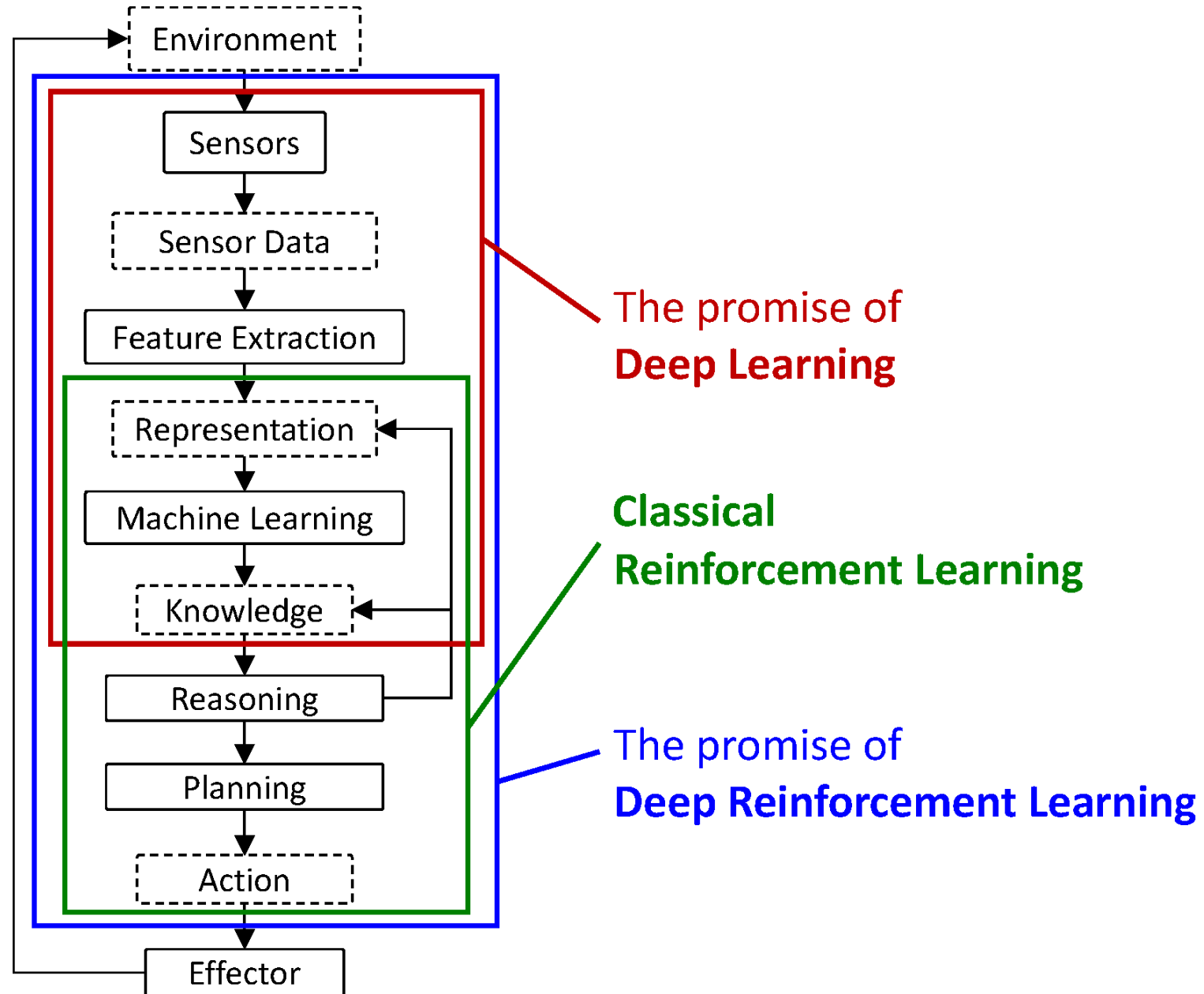


successful Neuro Dribbling at
RoboCup 2007 Atlanta, USA

Types of RL Algorithms

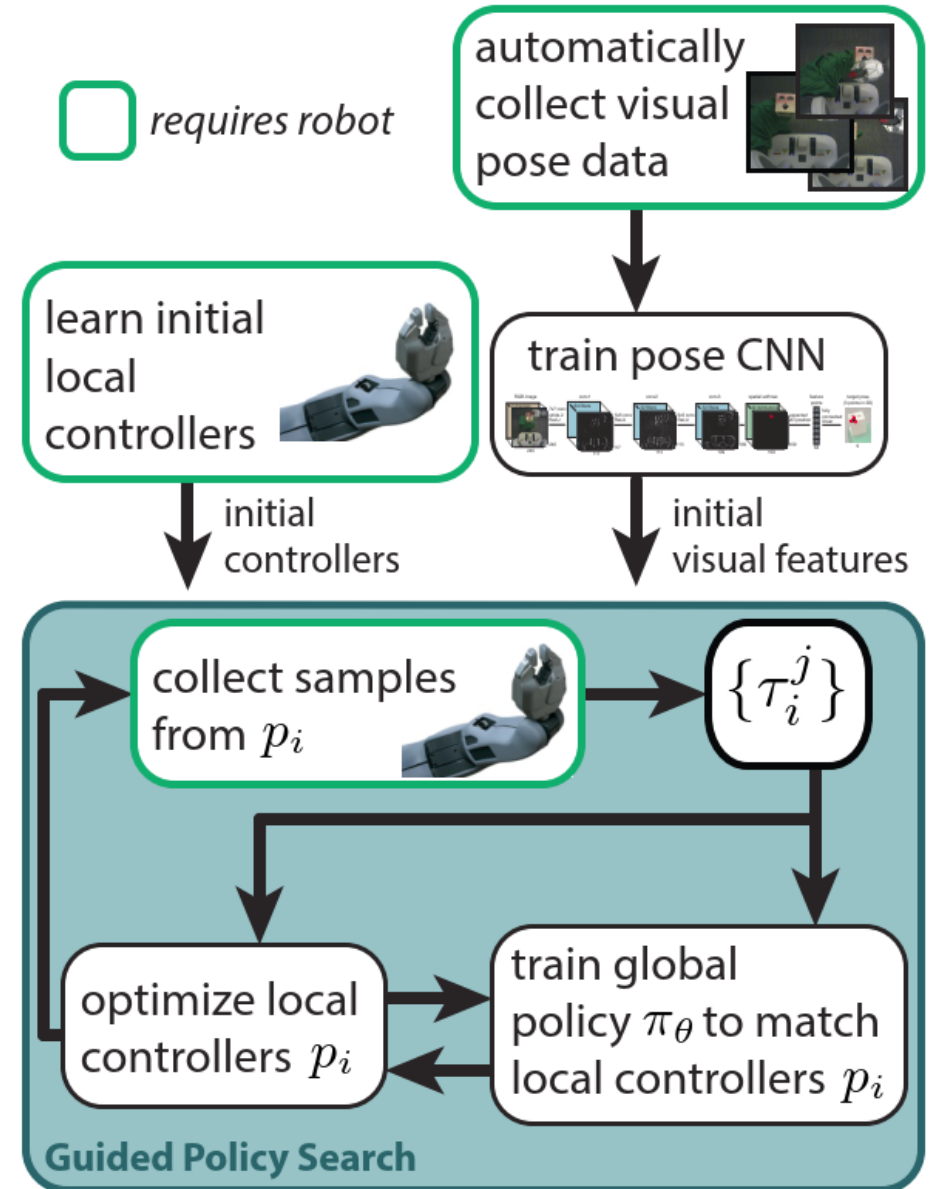


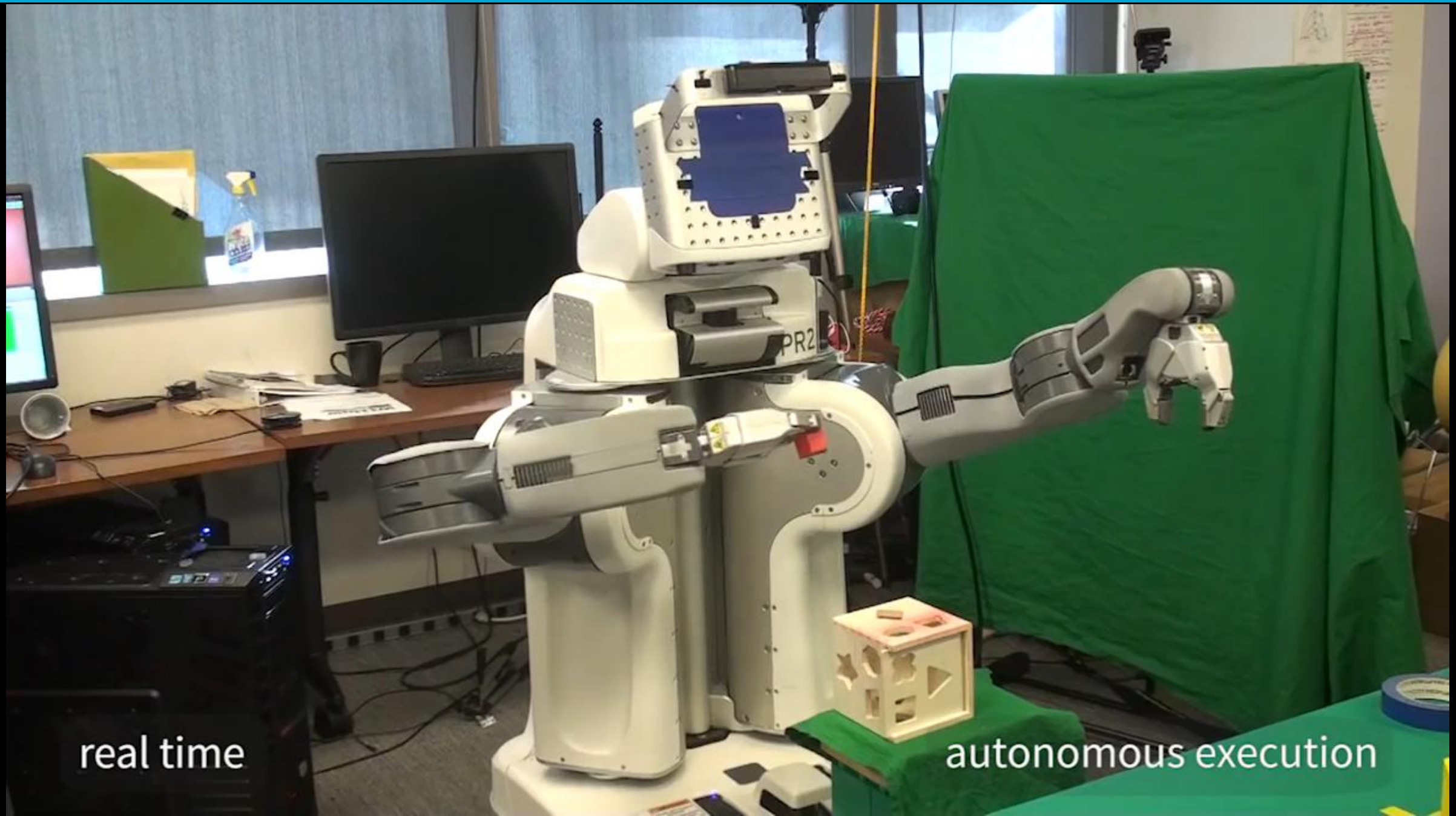




Guided Policy Search

- Real world
- Instrumented training





real time

autonomous execution

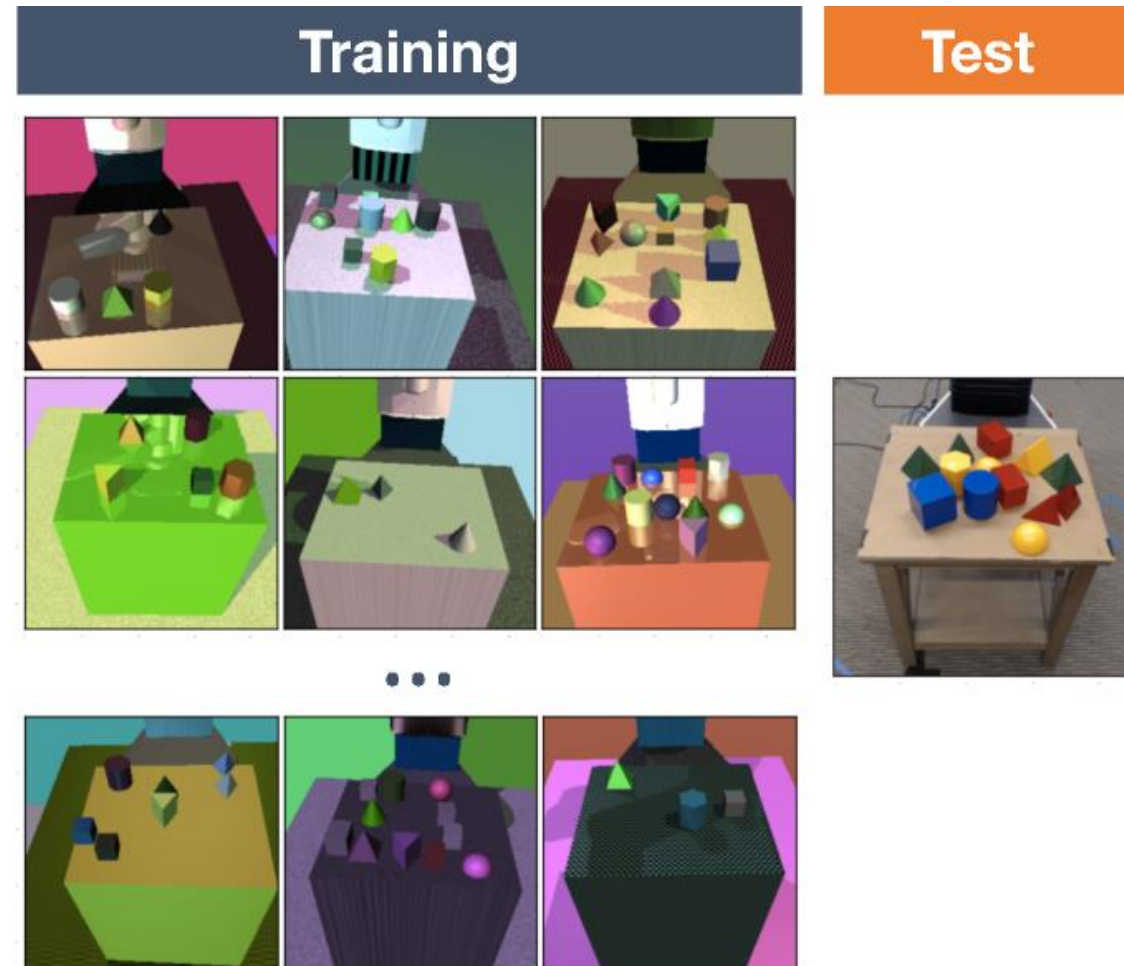
Sim2Real

Sim2Real

- Learn in Simulation, Transfer to Real-World
- Advantages:
 - Faster than real-time training
 - Simulated learning is inherently safe
 - Simulators can reset to any arbitrary state
 - Evaluating different task specifications is easier in simulation

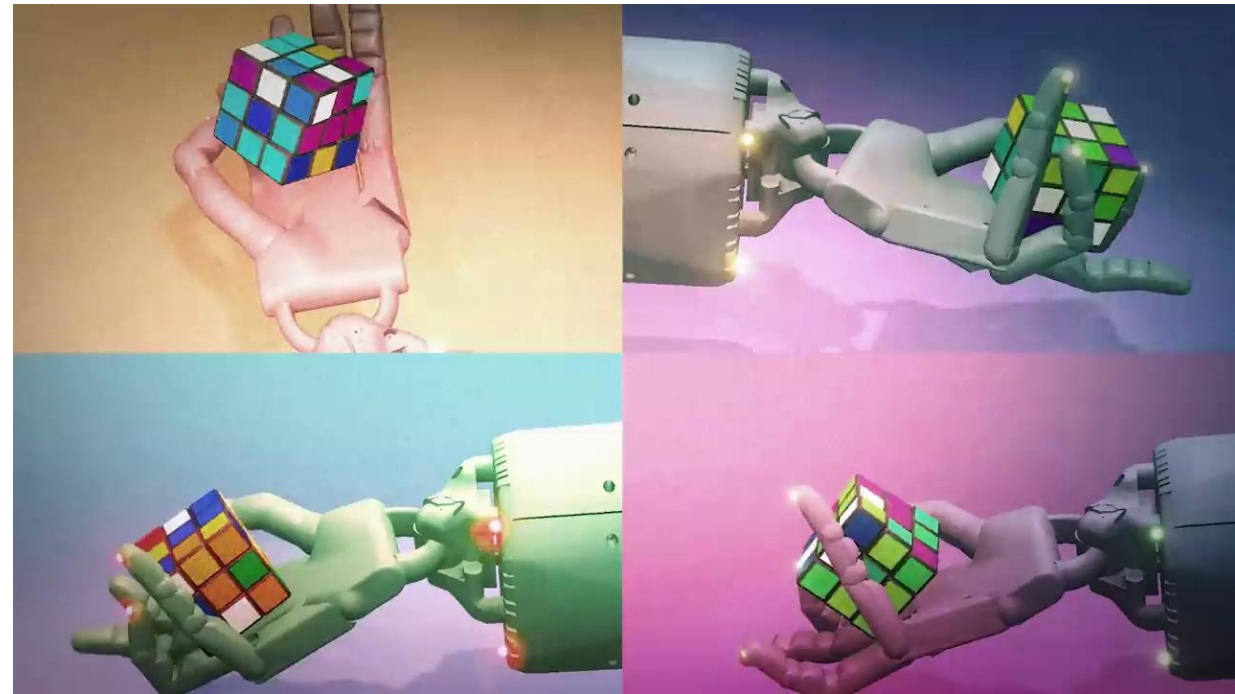
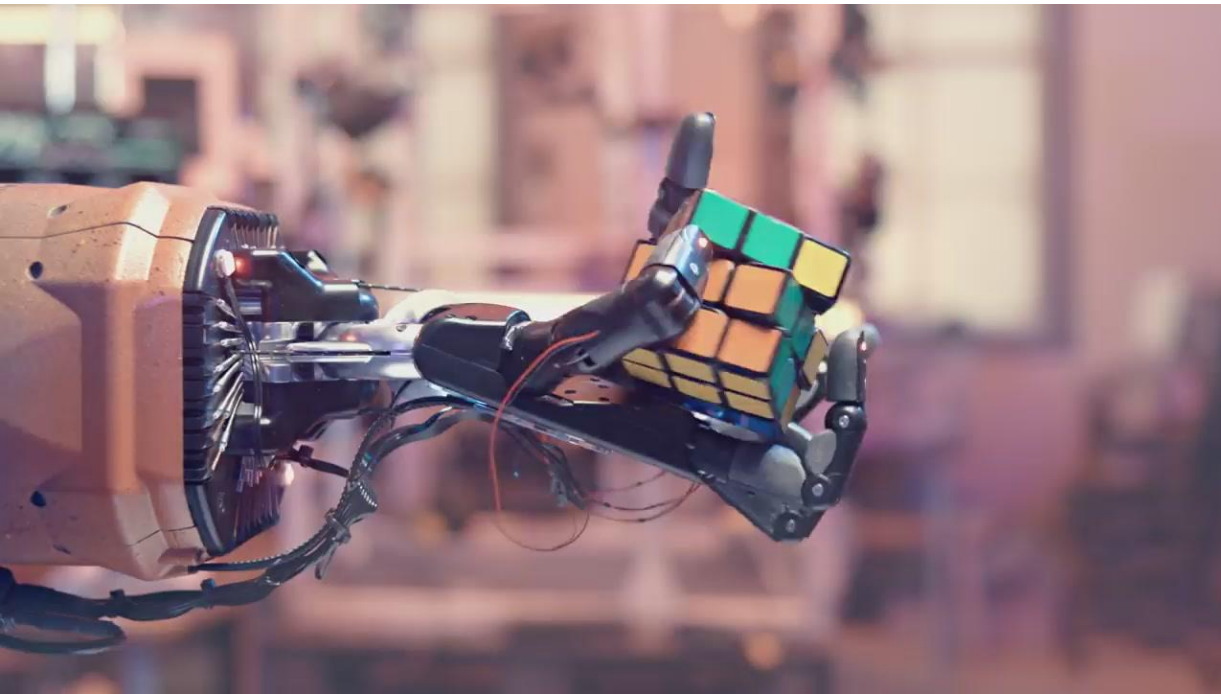
Training on Simulated Data

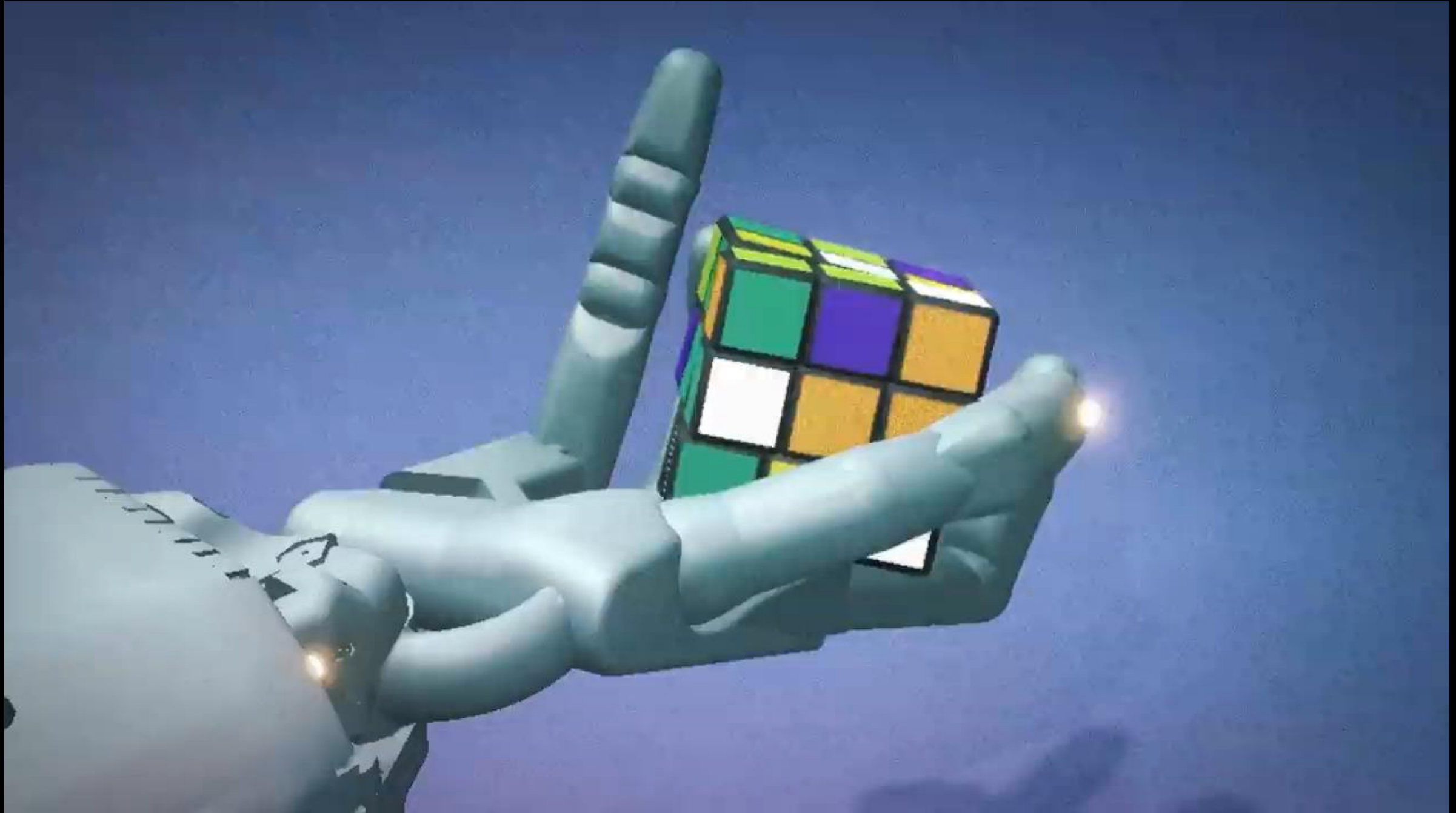
- Pre-train and finetune
- Direct Transfer
 - Make simulation very realistic
 - Make behavior robust



Solving a Rubik's Cube

- Deep RL in simulation





The international journal of science / 10 February 2022

nature



DRIVING FORCE

AI algorithm outcompetes human
champions in *Gran Turismo* racing game

Science Robotics

SEPTEMBER 2023



The international journal of science / 31 August 2023

nature

DRONE RACER

AI pilot beats human champions in aerial contest



Offset agreement
Overhaul pricing of carbon credits to help fund climate projects

Dining companions
Corals devour algal partners when food supplies run low

Sentence synthesis
Brain implants show promise in rendering speech from thoughts

THE COVER: ANDREW HARRIS

Will it transfer?

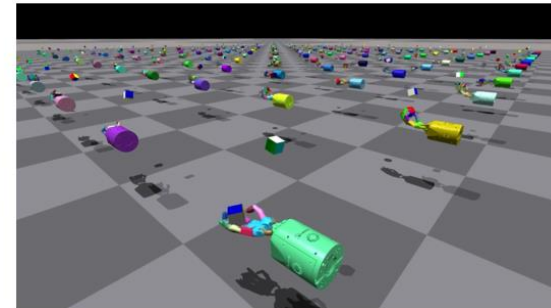
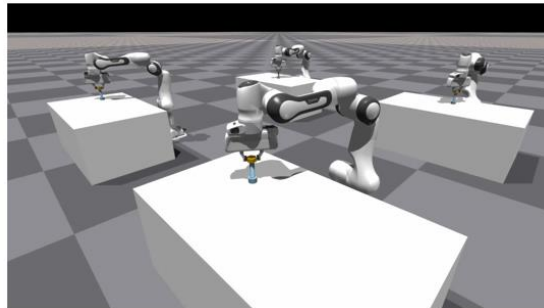
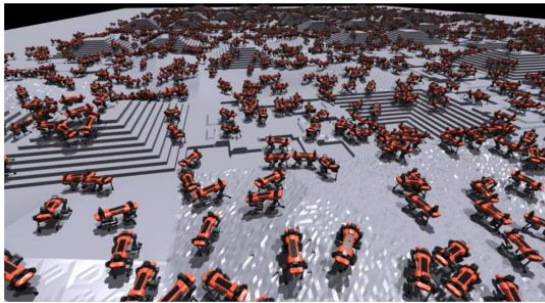


Guidelines for Sim2Real

- Minimize model mismatch
- Acknowledge all models are wrong
- Mitigate relevant mismatches

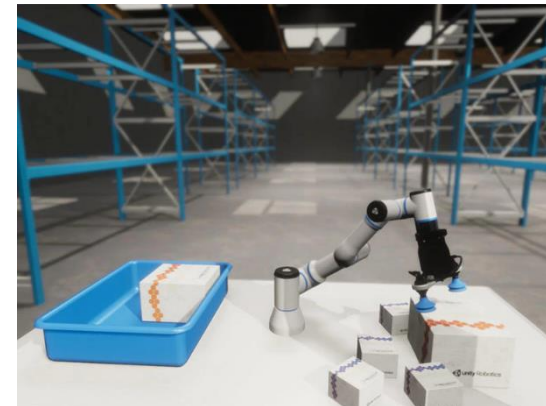
(None-Open) Tools

- Isaac Lab <https://developer.nvidia.com/isaac/lab>



- Unity <https://unity.com/solutions/automotive-transportation-manufacturing/robotics>

- Many more

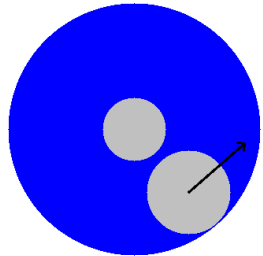


Common Pitfalls

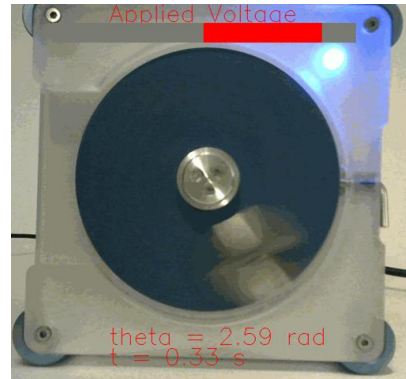
Inaccuracies in the Simulated Dynamics

Inaccurate Model

Applied Voltage



theta = 2.43 rad
t = 0.20 s

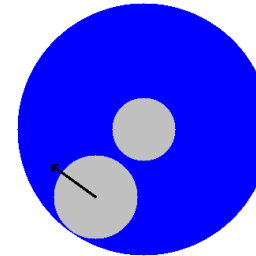


theta = 2.59 rad
t = 0.33 s

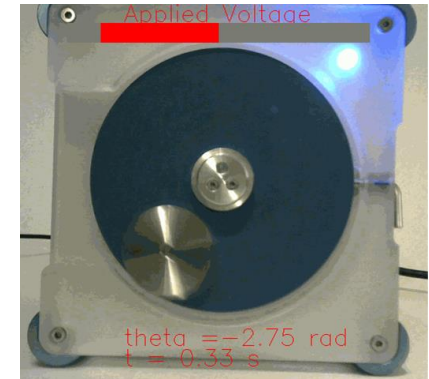


Domain Randomization

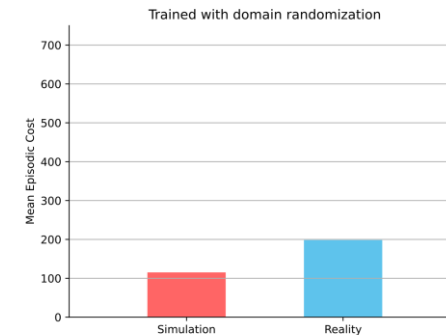
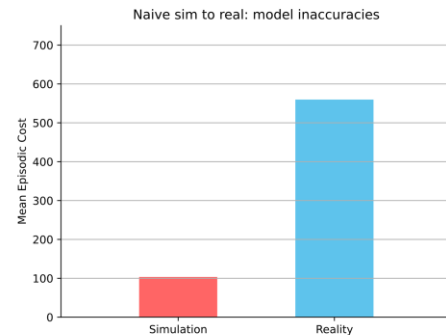
Applied Voltage



theta = -2.52 rad
t = 0.20 s



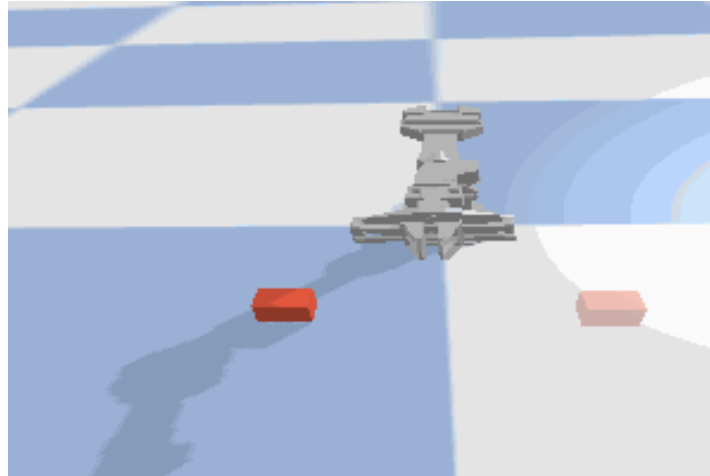
theta = -2.75 rad
t = 0.33 s



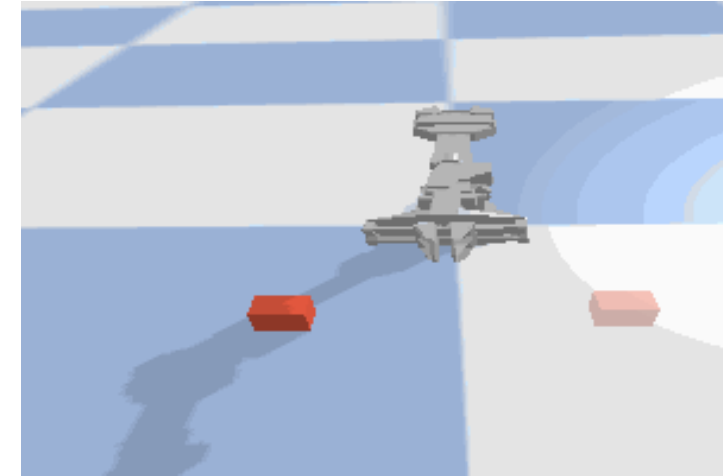
Unknowingly Changing the Markov-Decision Process

Training

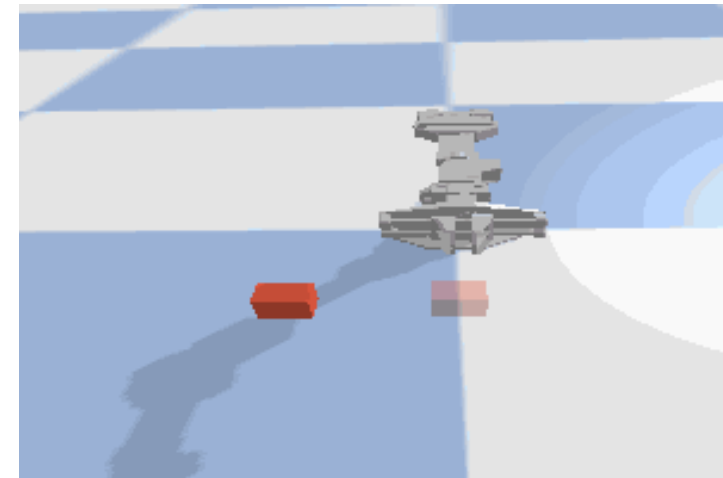
Trained **without**
evaluated **with**
collision avoidance.



Evaluation

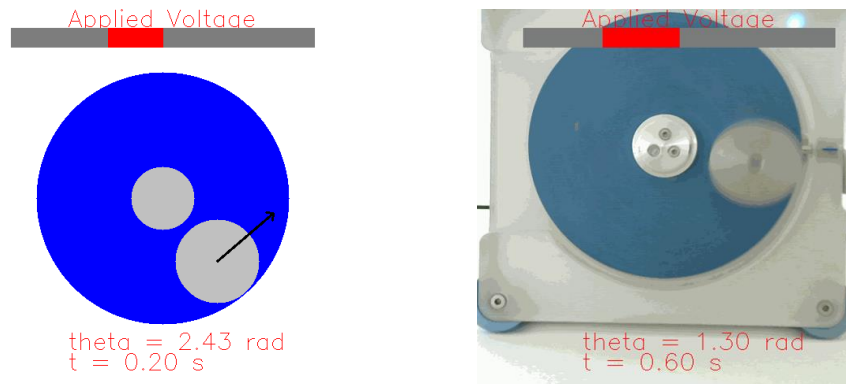


Trained **with**
evaluated **with**
collision avoidance.

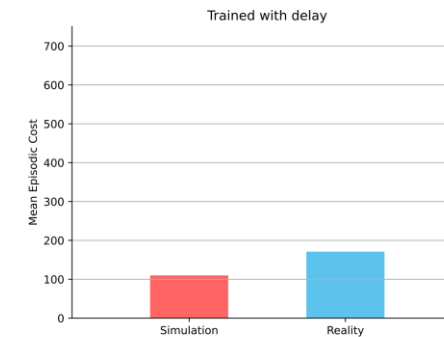
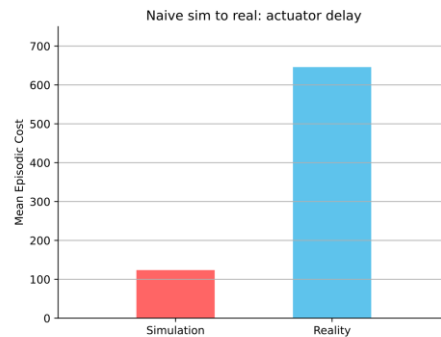
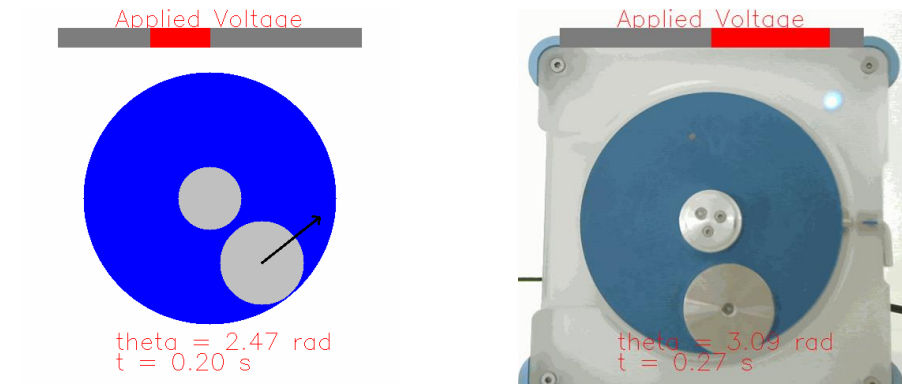


Delays

Ignored

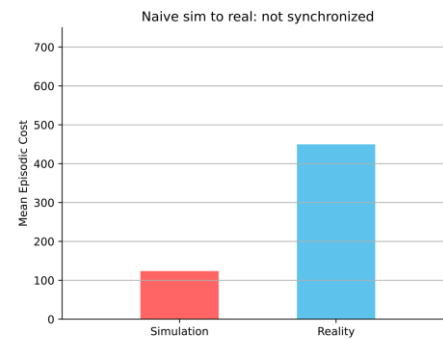
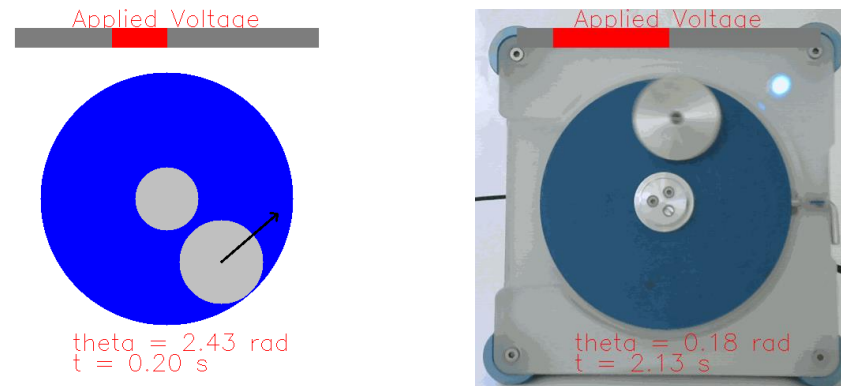


Included in Simulation



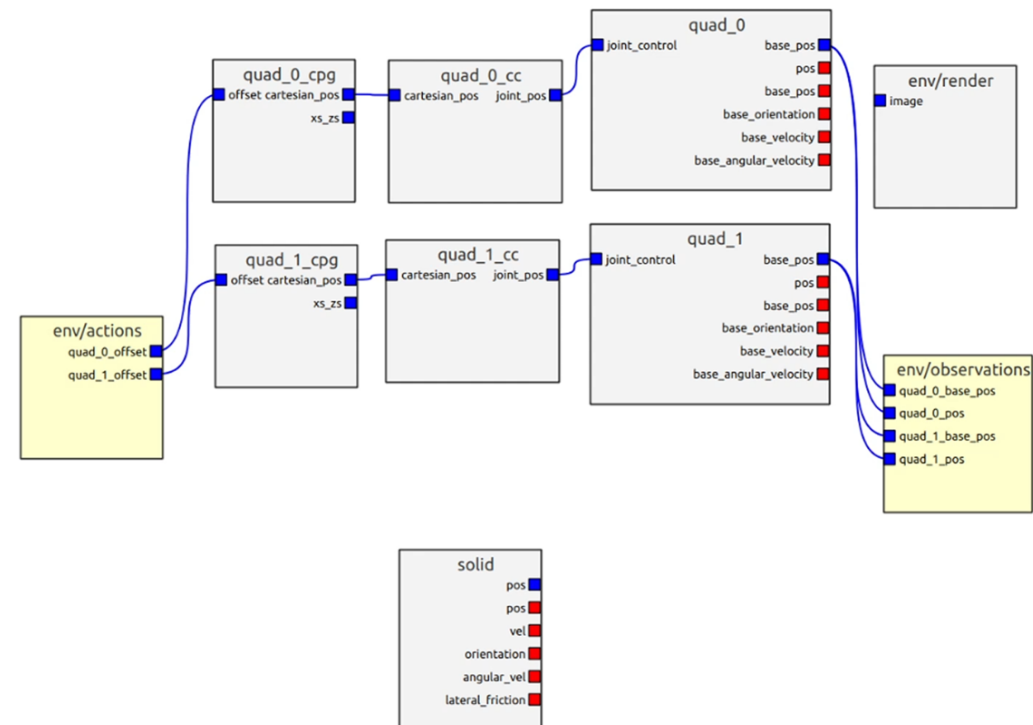
Asynchronous Frameworks

Naive



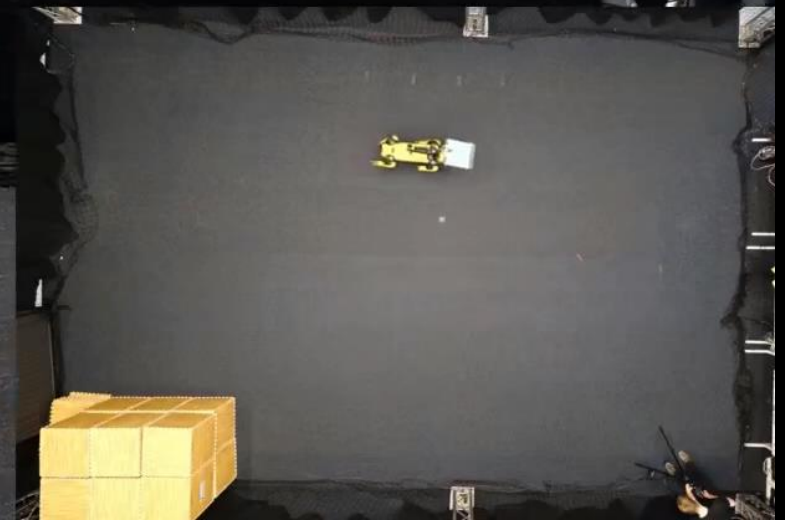
EAGERx (Engine Agnostic Graph Environments for Robotics)

“A Python3 framework that lets you easily define OpenAI gym compatible environments that work both in simulation and the real-world.”





Slow motion (0.33x)



Path following with $R=1.0$ meter
(see fig. 9a)



delay sim.



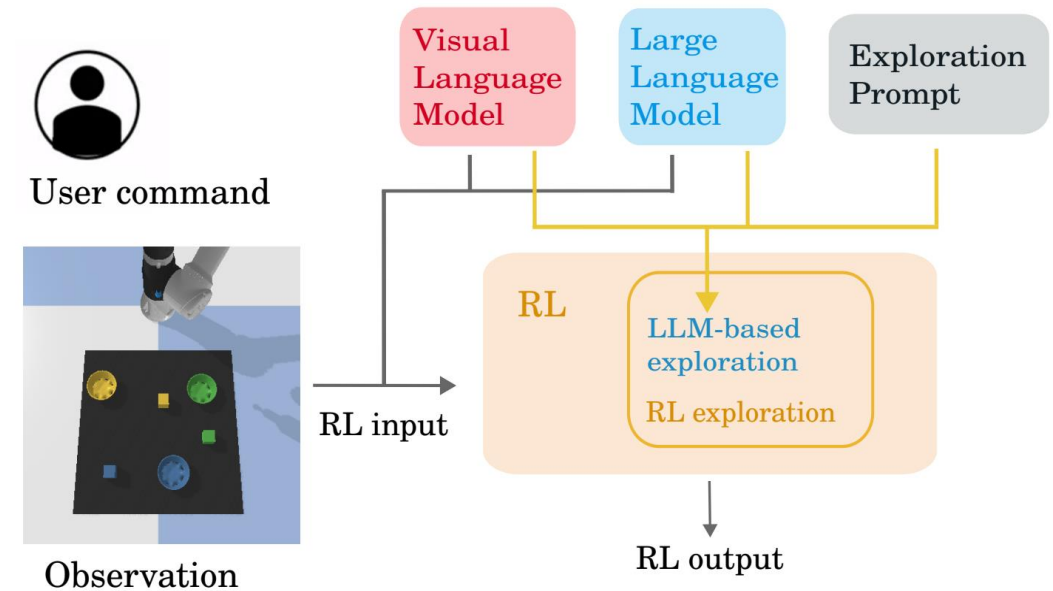
no delay sim.

Background lights turned
off for better visibility

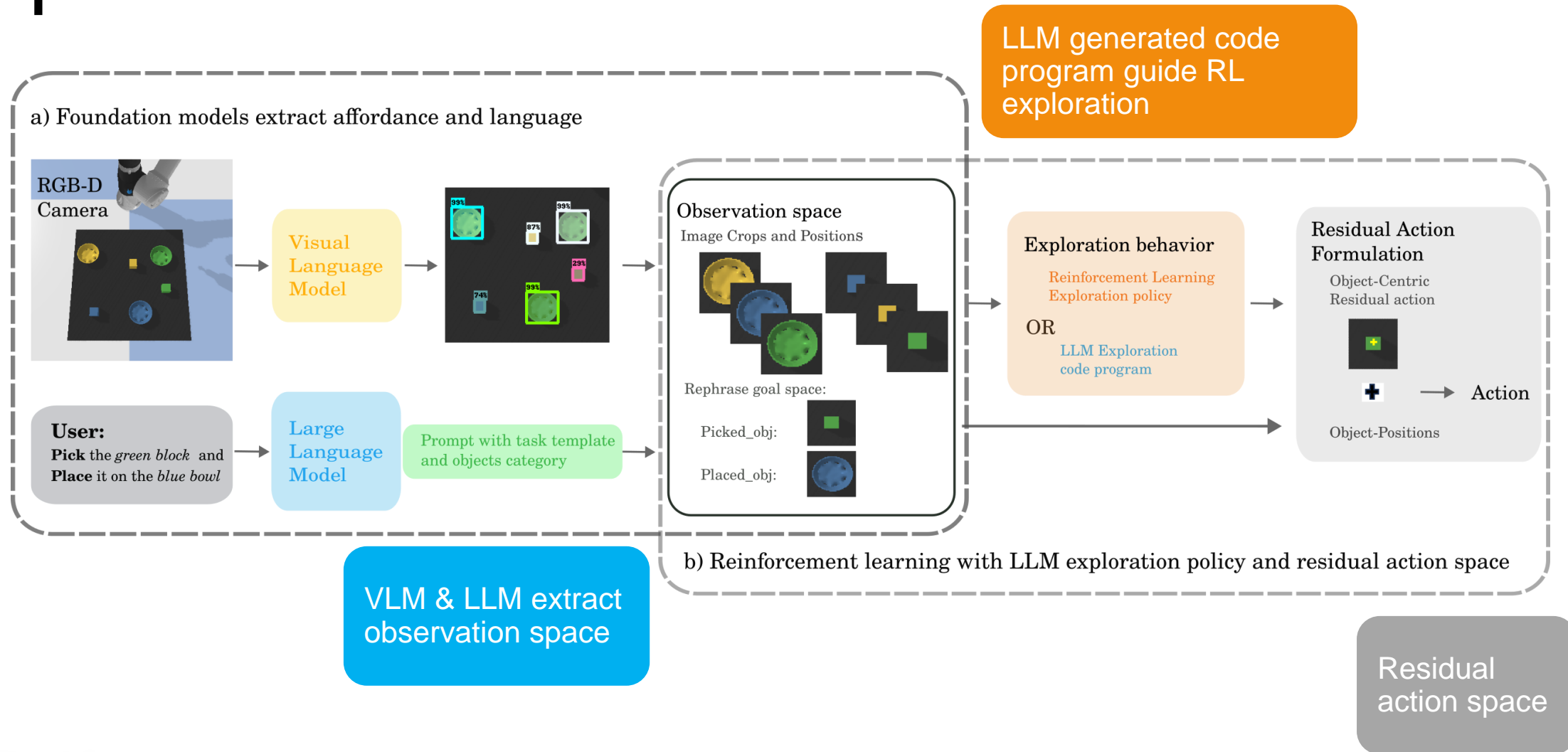
Other Ideas

Language conditioned robot policy learning

- Method: Foundation Models + RL
- Utilizing the commonsense knowledge in foundation models to help end-to-end reinforcement learning
- Foundation models guide exploration for reinforcement learning
- Foundation models extract observation space

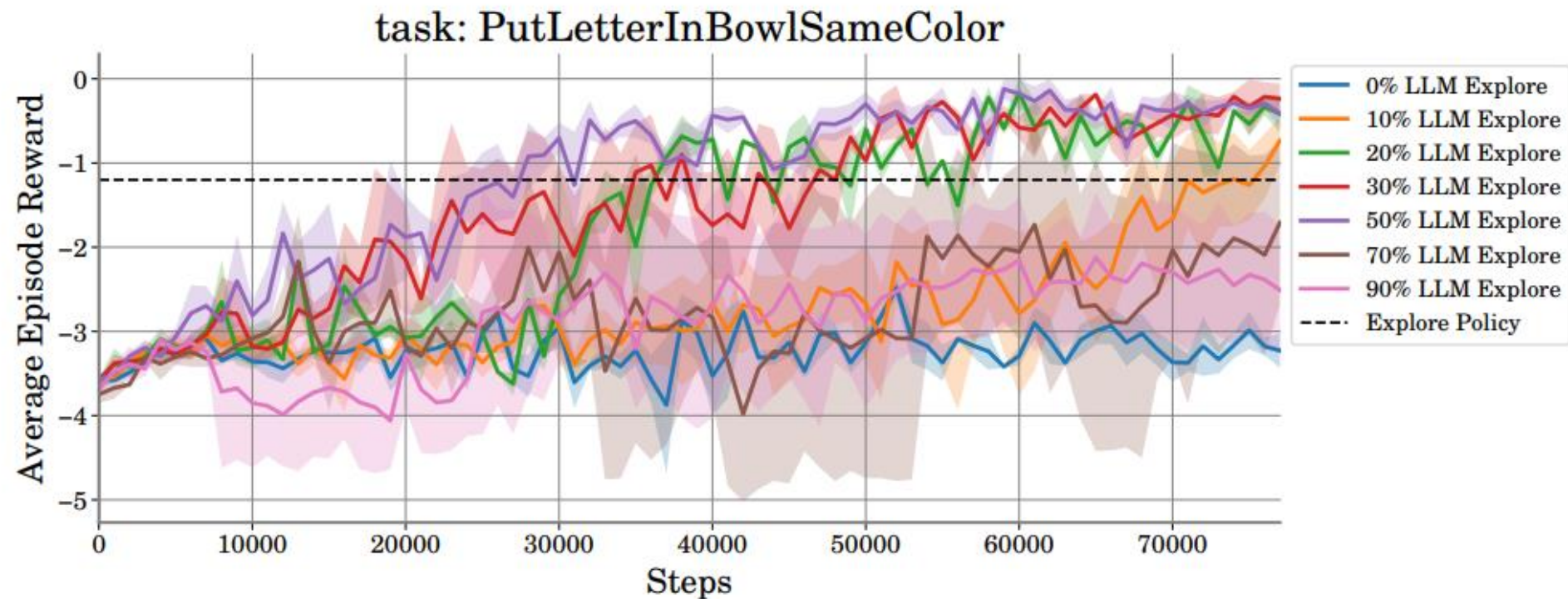


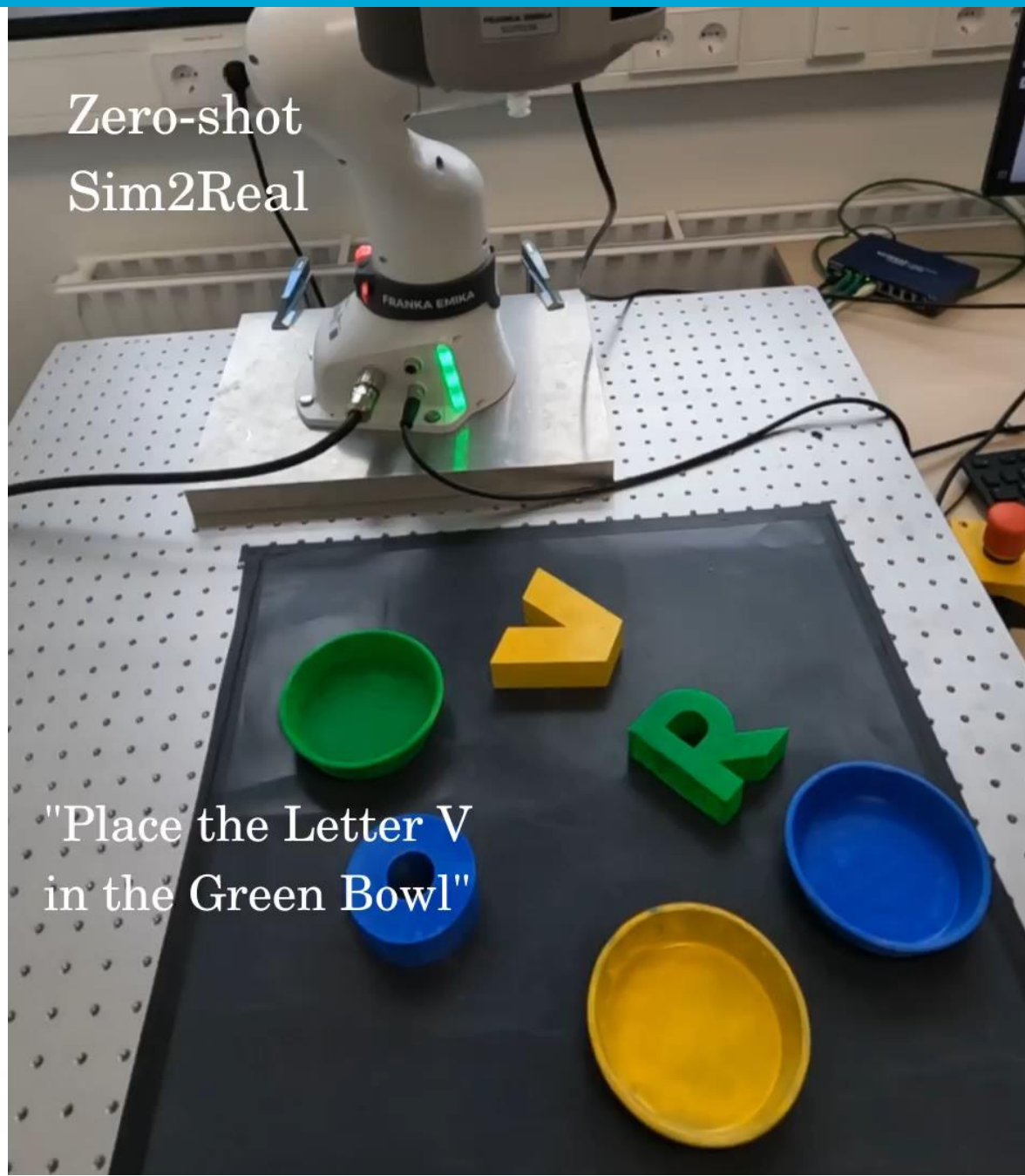
ExploRLLM



Experiment

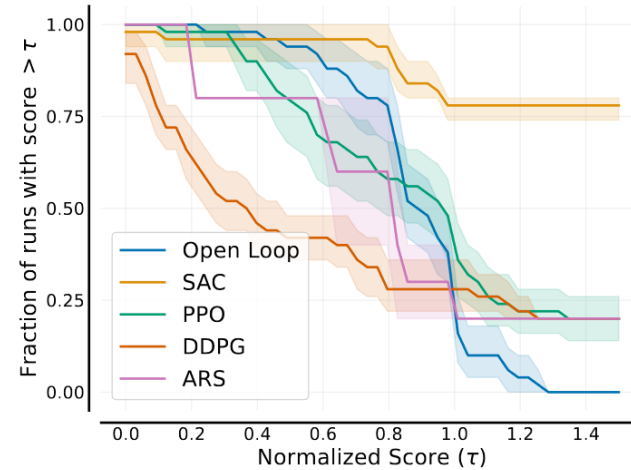
Analyze influence of LLM exploration frequency





Central Pattern Generators

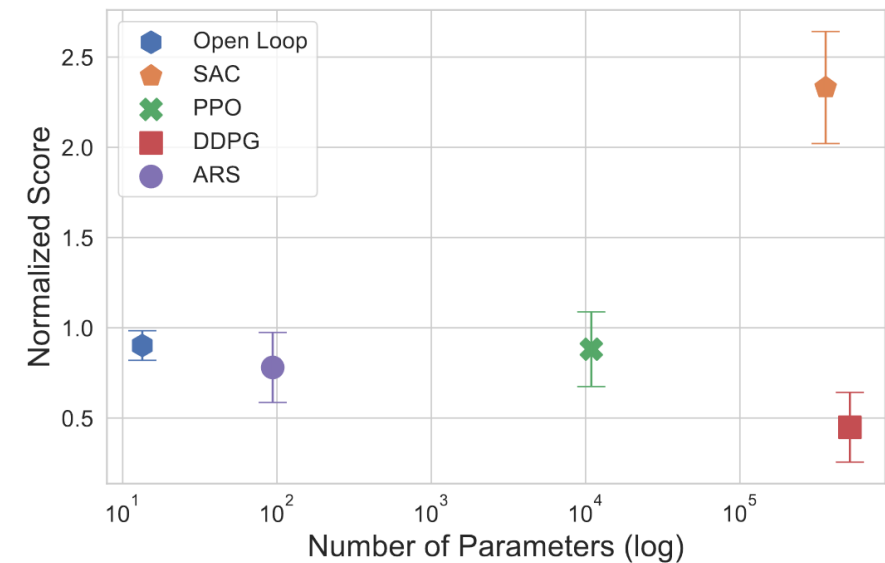
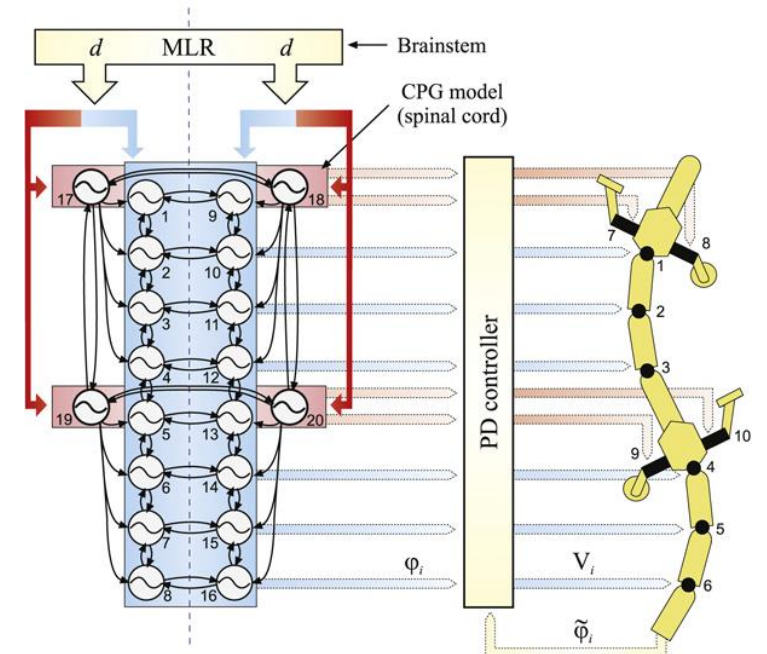
- Perform as well



- But are a lot more efficient

	SAC		PPO		DDPG		ARS		Open-Loop	
	CPU	GPU	CPU	GPU	CPU	GPU	CPU	GPU	CPU	GPU
Runtime (in min.)	80	30	10	14	60	25	5	N/A	2	N/A

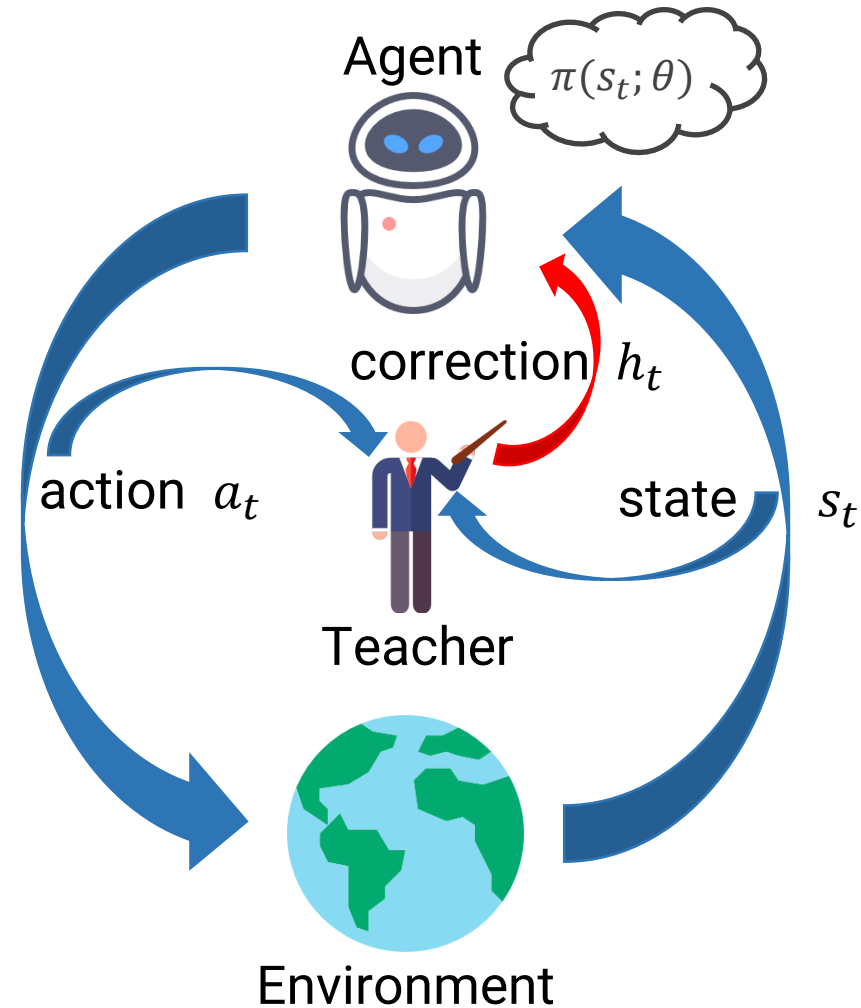
- And a lot more robust



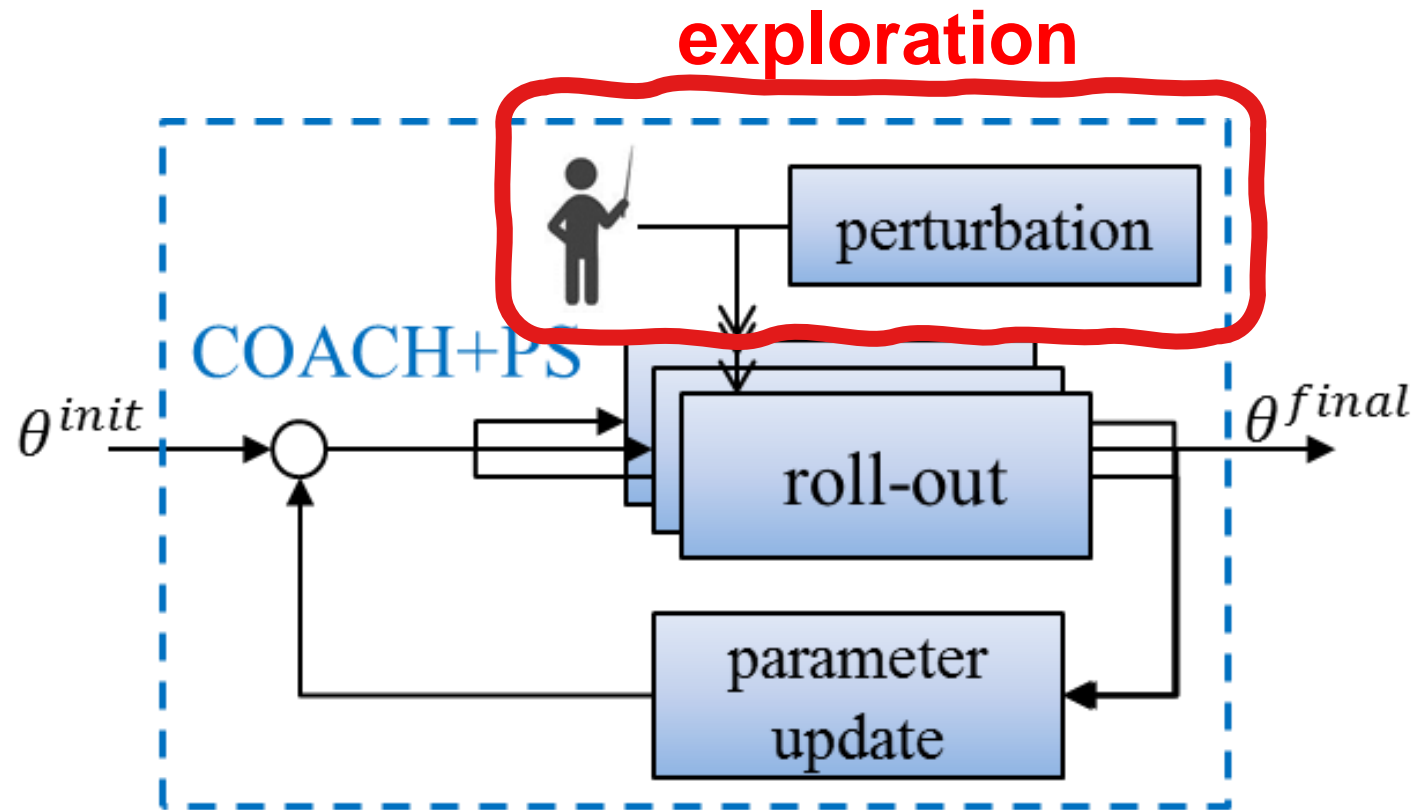
RL in Simulation

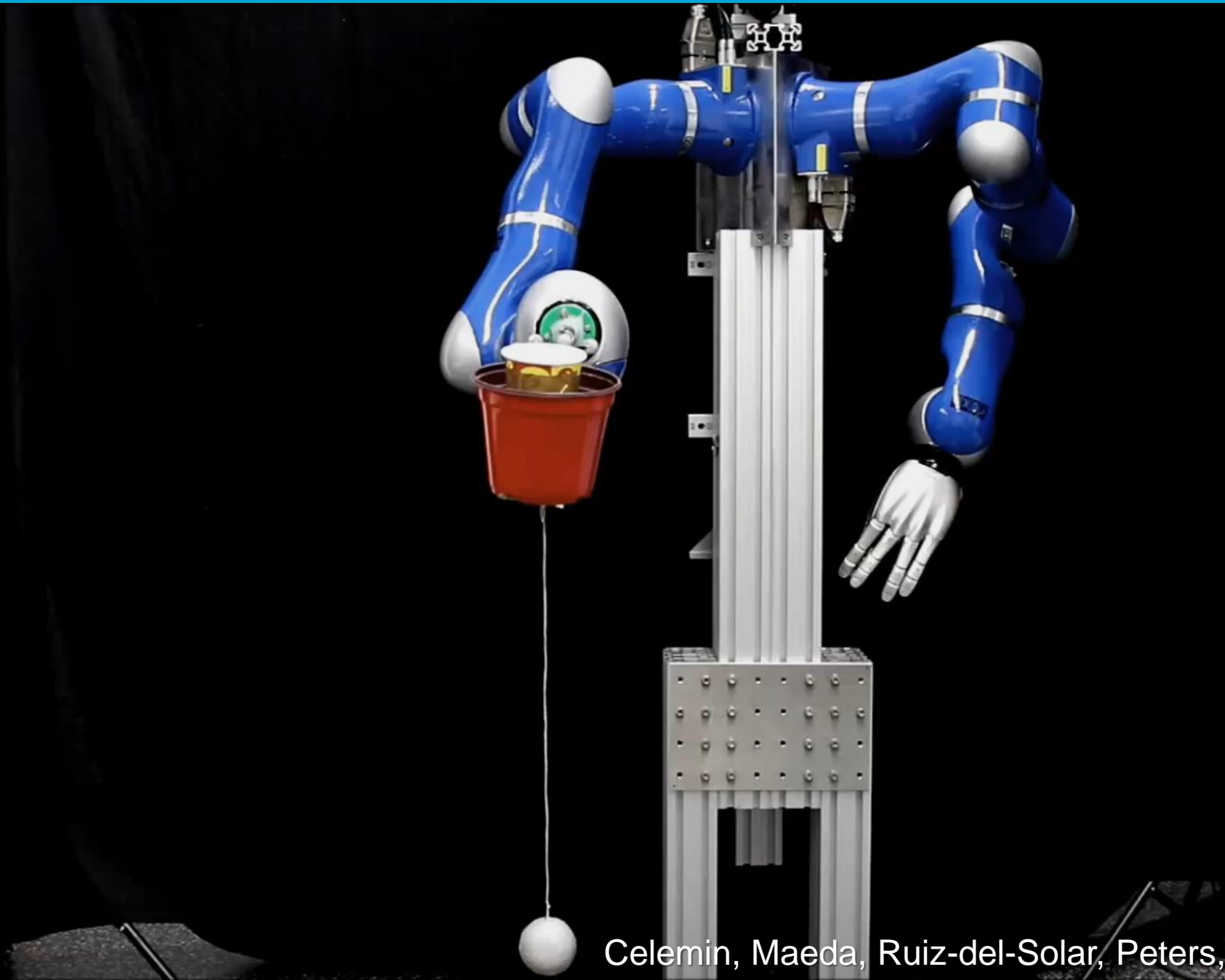
Interactive Reinforcement Learning

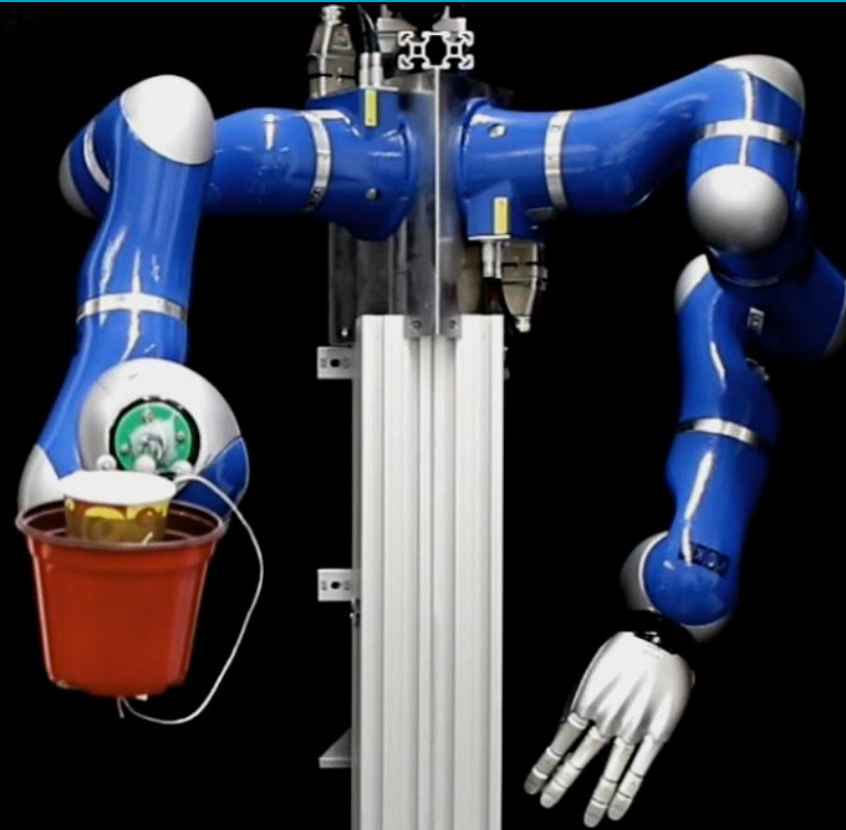
Interactive Learning



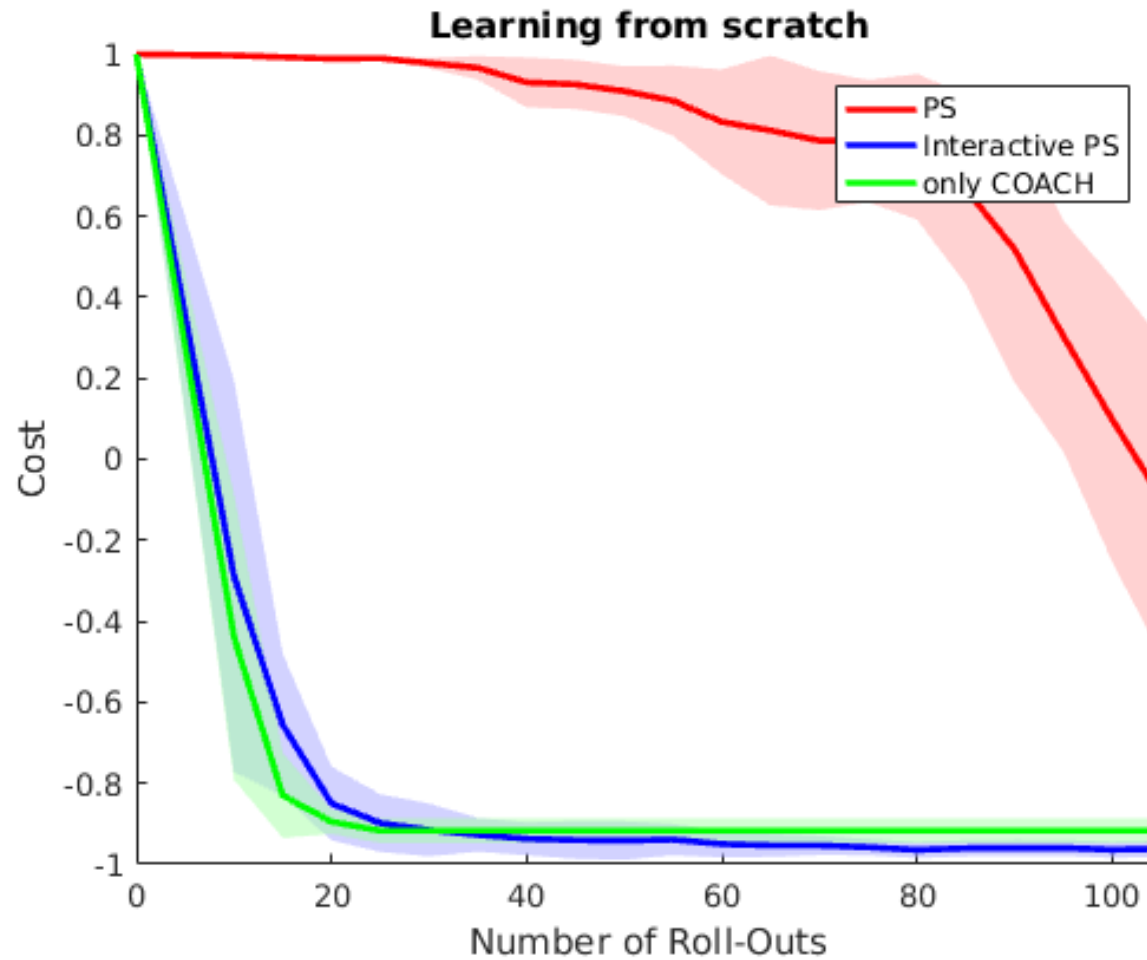
IL + RL: Simultaneous





0.5X

Results



Summary

Summary

- Specific challenges
 - Small (real) data
 - Safety
- Tractability through prior knowledge
 - Policy structure
 - Demonstrations
 - Simulations
 - Foundation models
 - Etc.