



Reinforcement and Imitation Learning to Control Contact-Rich Interactions

Roberto Martín-Martín, RobIn Lab, UT Austin



RobIn
ROBOT INTERACTIVE
INTELLIGENCE LAB

TEXAS
Robotics

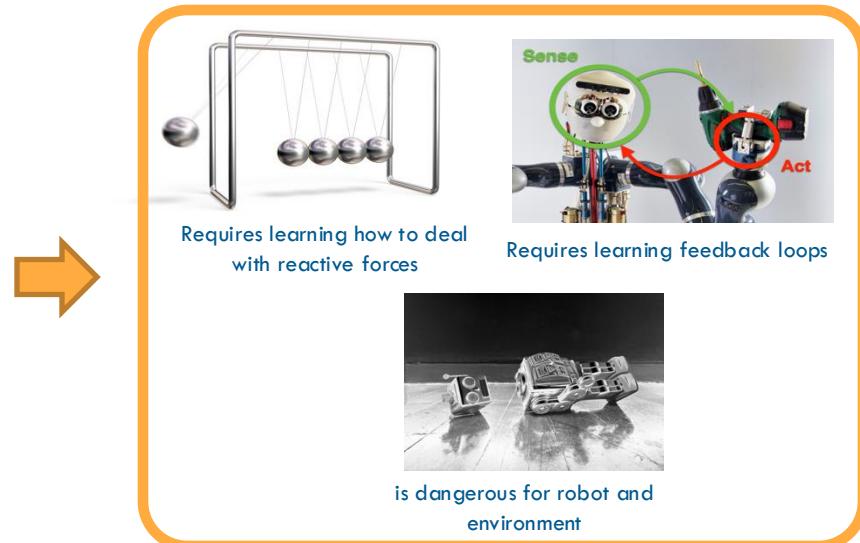
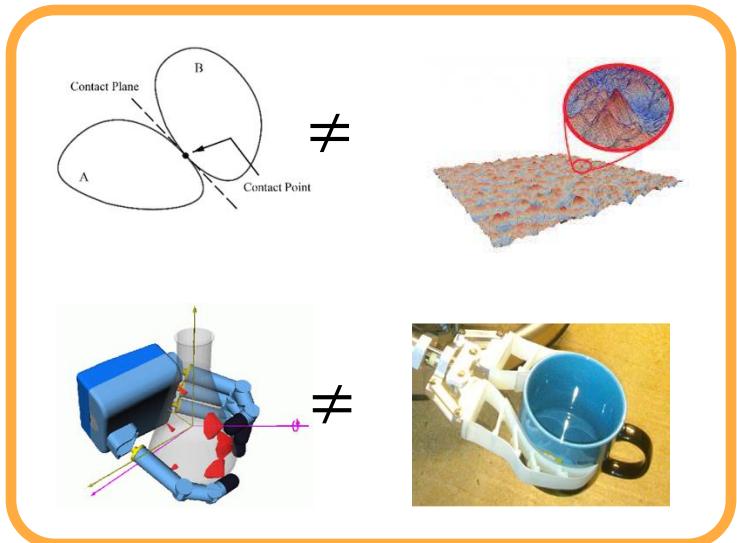


A young boy with short brown hair, wearing a blue t-shirt, is lying on his stomach on a green lawn. He is looking up and to the right with a joyful expression, his hands resting under his chin. To his left is a red plastic toy car with a smiling face, featuring large white eyes with black pupils and a small mouth. The car has a yellow front grille and two black wheels with white hubcaps. The background shows a wooden deck and some bushes.

Challenging Contact-Rich Interactions!

Why is hard to control Contact-Rich Manipulation?

Predicting the effect of contact interactions is difficult





Action

Perception

Learning Tasks in Different Action Spaces



Configuration/Joint Space



Task/Cartesian Space

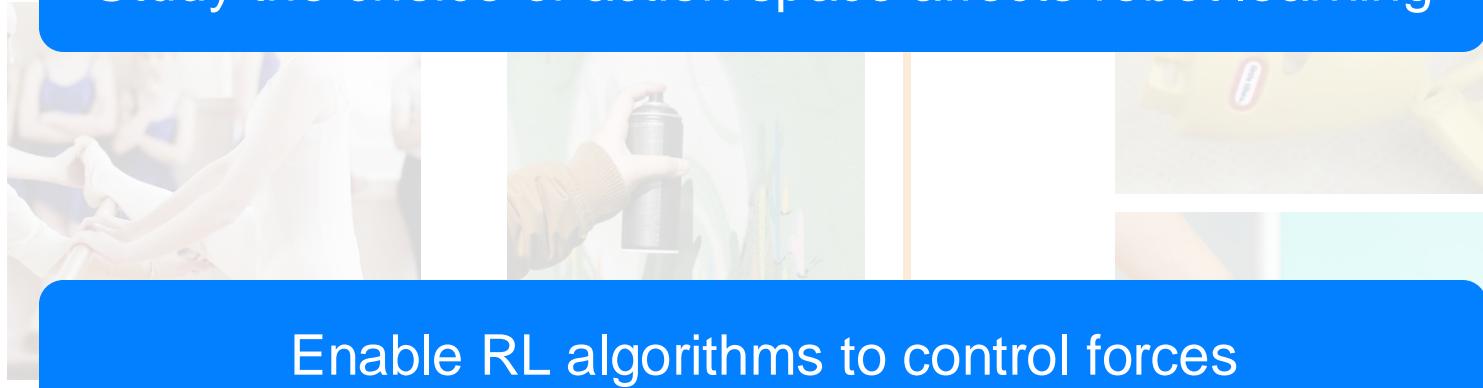


Motion

Forces

Learning Tasks in Different Action Spaces

Study the choice of action space affects robot learning



Enable RL algorithms to control forces

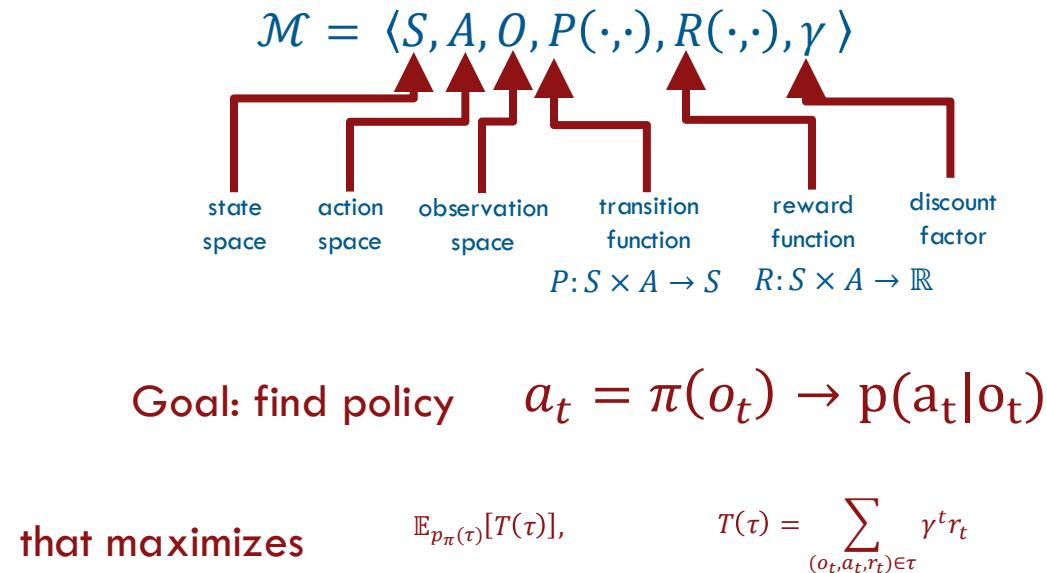
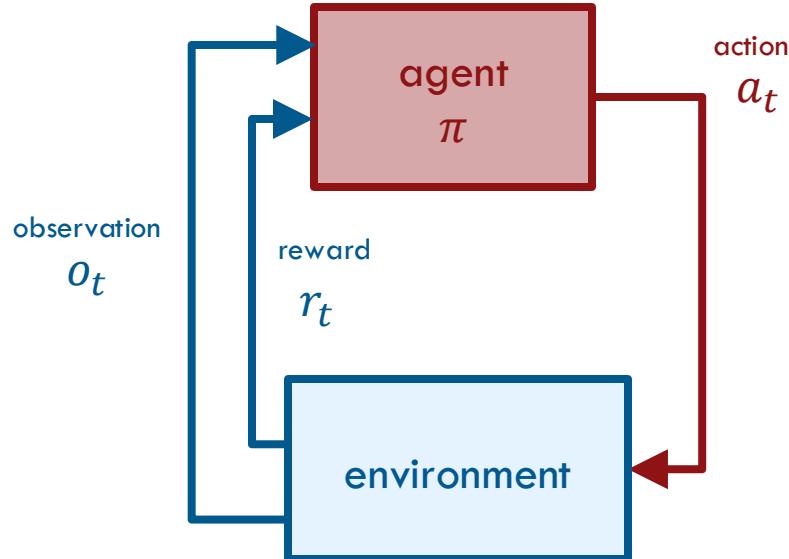
Configuration/Joint Space

Task/Cartesian Space

Motion

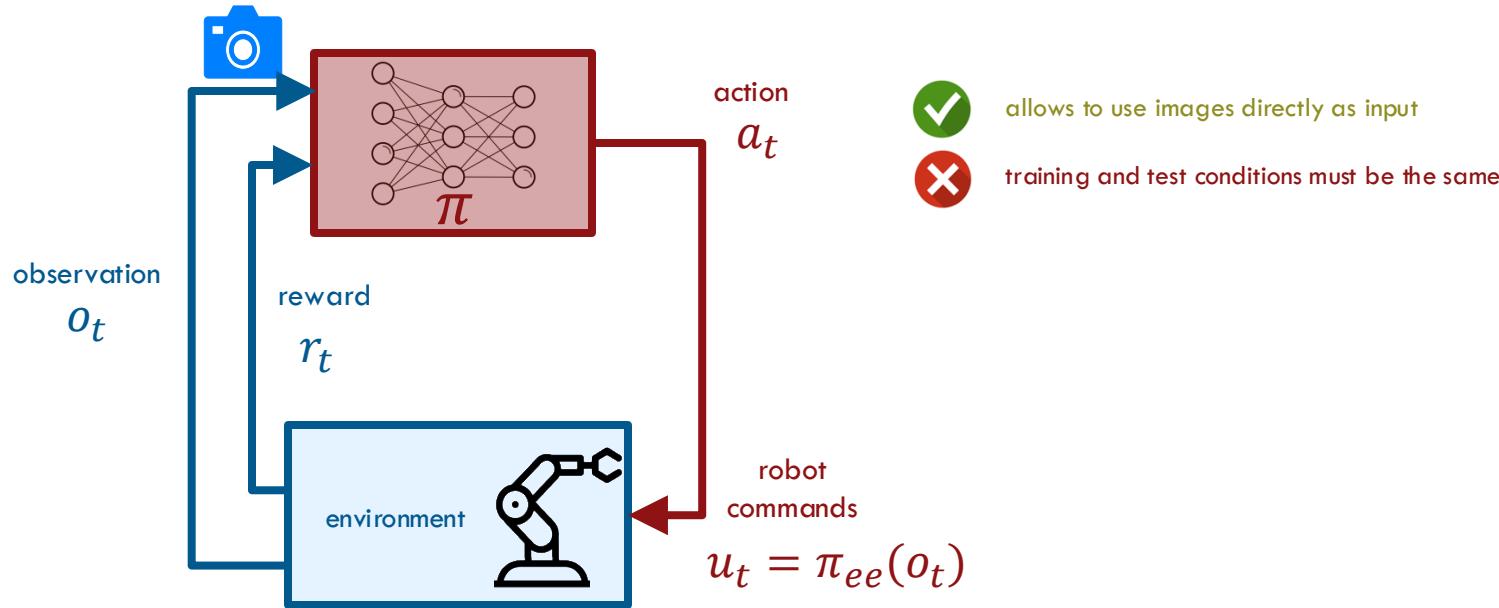
Forces

Reinforcement learning to control physical interactions



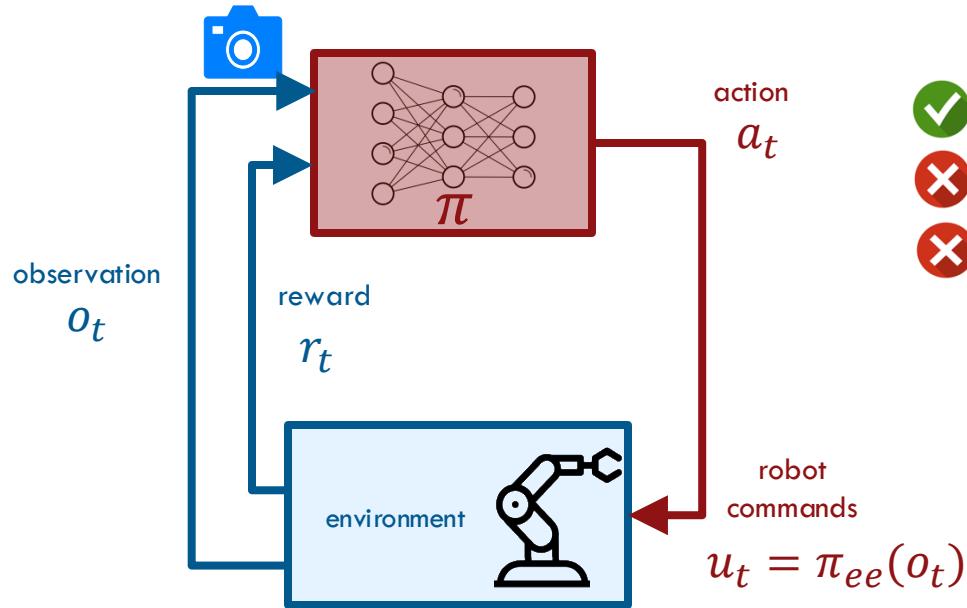
Deep reinforcement learning to control physical interactions

End-to-end learning from pixel to torques



Deep reinforcement learning to control physical interactions

End-to-end learning from pixel to torques



allows to use images directly as input



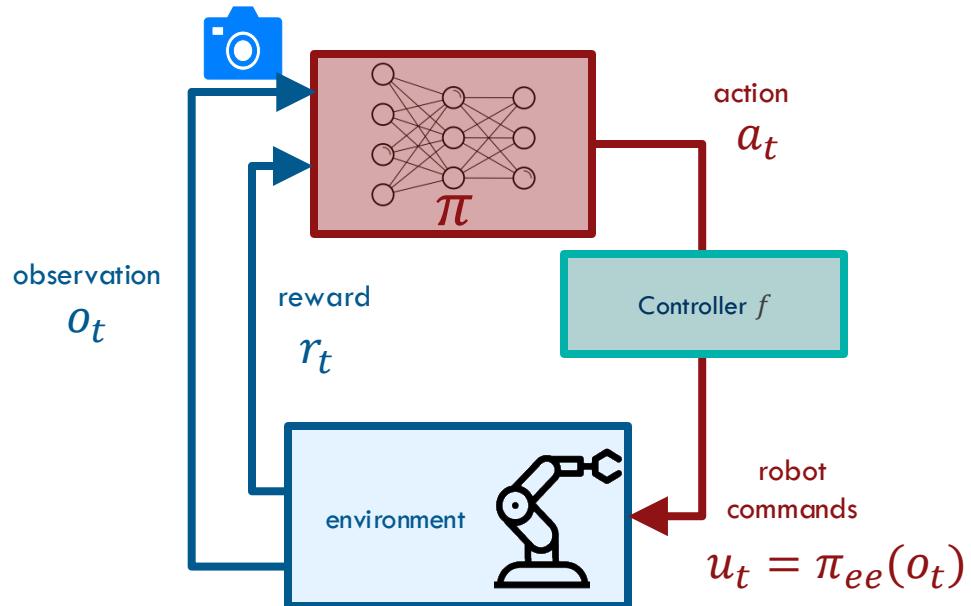
training and test conditions must be the same



data hungry: re-learn how to control robot's end-effector motion and forces for every task



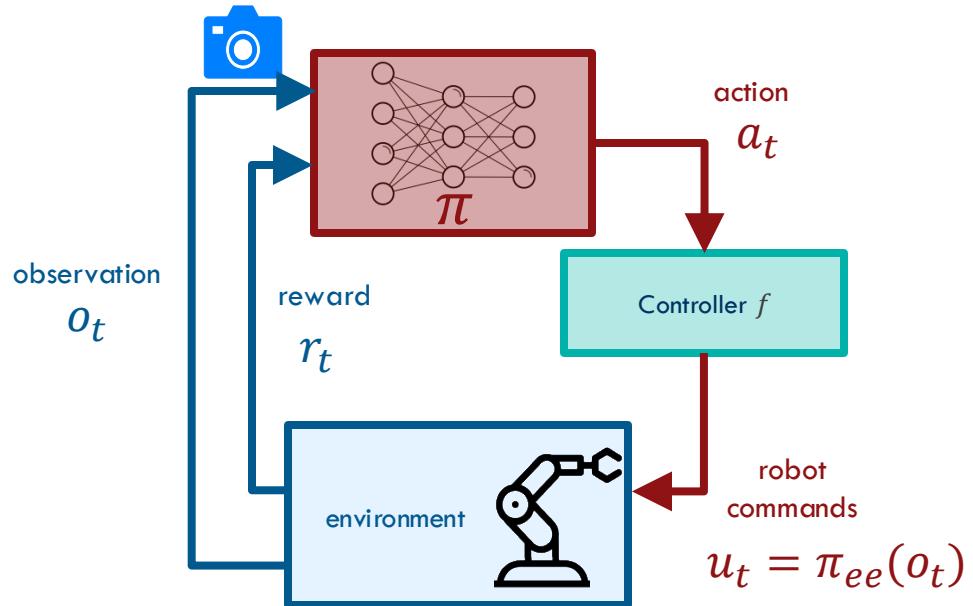
RL with Analytical Controller



$$a_t = \pi(o_t)$$

$$u_t = \pi_{ee}(o_t)$$

RL with Analytical Controller



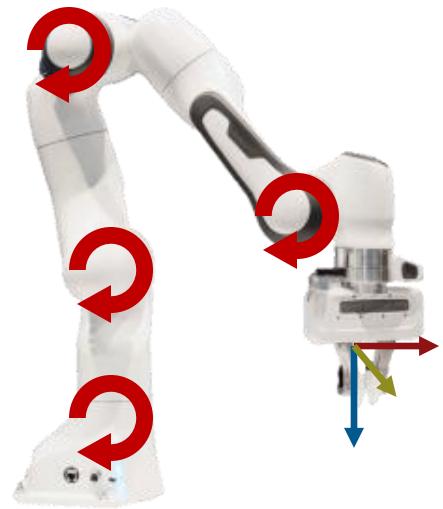
Reference Generator (learned)

$$u_t = f(\pi(o_t))$$

Robot Control (analytic)

$$u_t = \pi_{ee}(o_t)$$

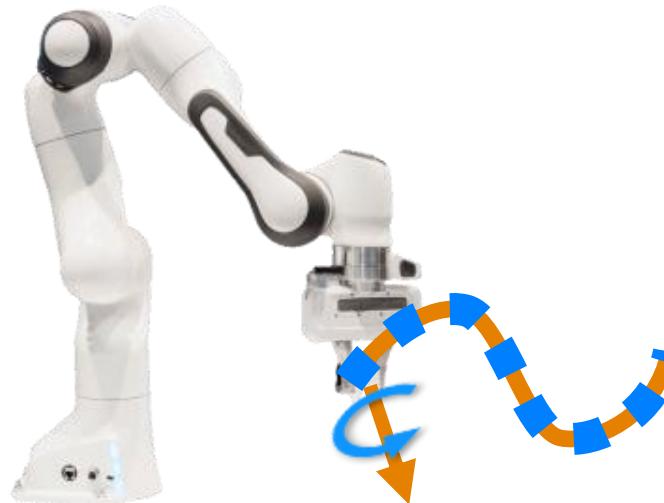
A Study of Action Spaces in RL



Joint Space

Cartesian/Task Space

A Study of Action Spaces in RL



Joint Space

Cartesian/Task Space

Motion: Positions/Velocities/Accelerations

Forces and Torques

VICES: Variable Impedance in End-Effector Space



VICES

Joint Space

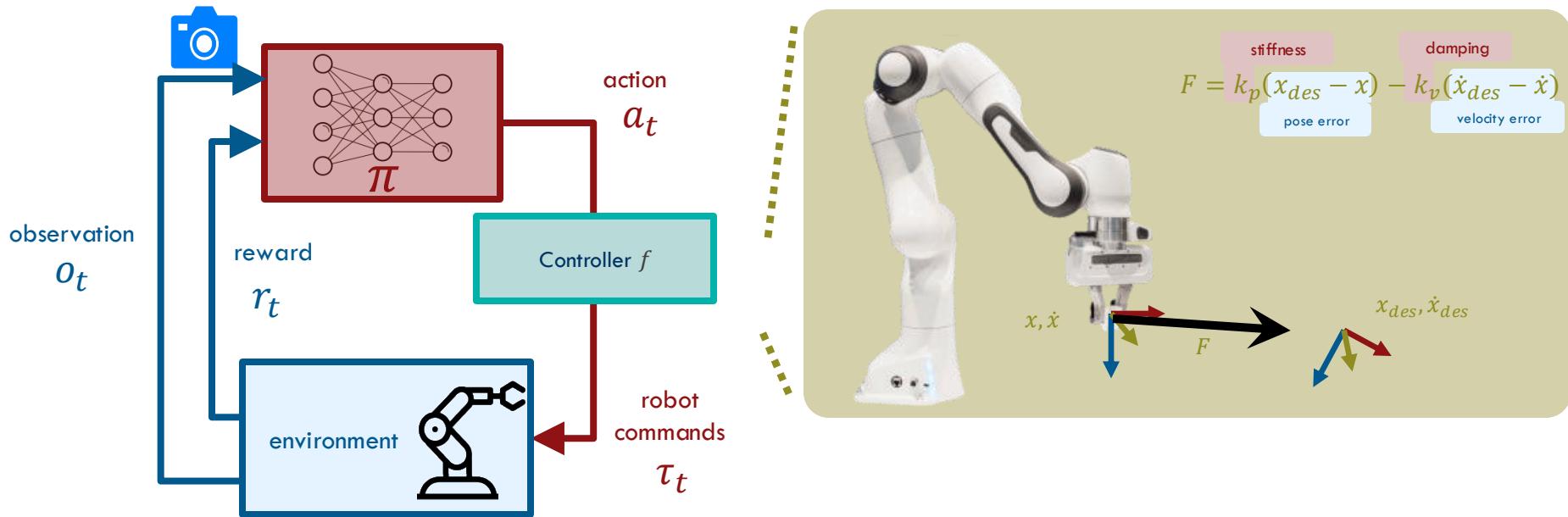
Motion: Positions/Velocities/Accelerations

Cartesian/Task Space

Forces and Torques

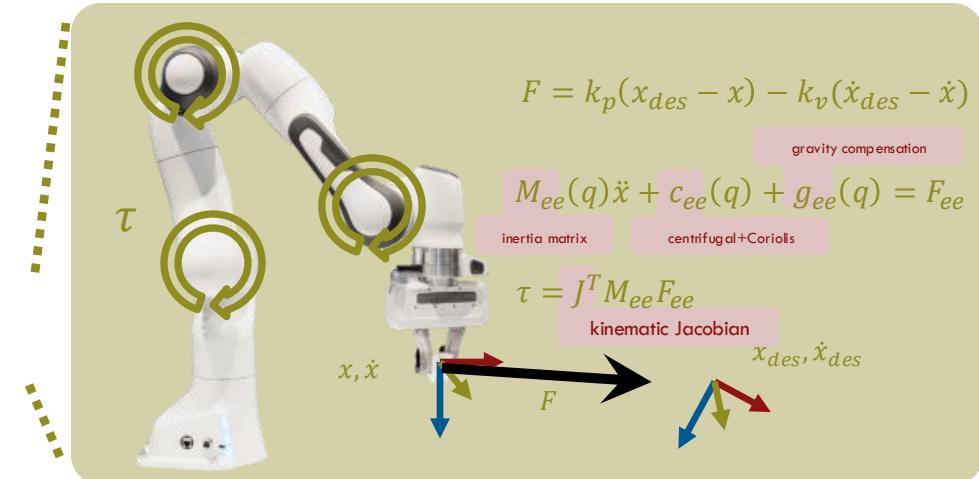
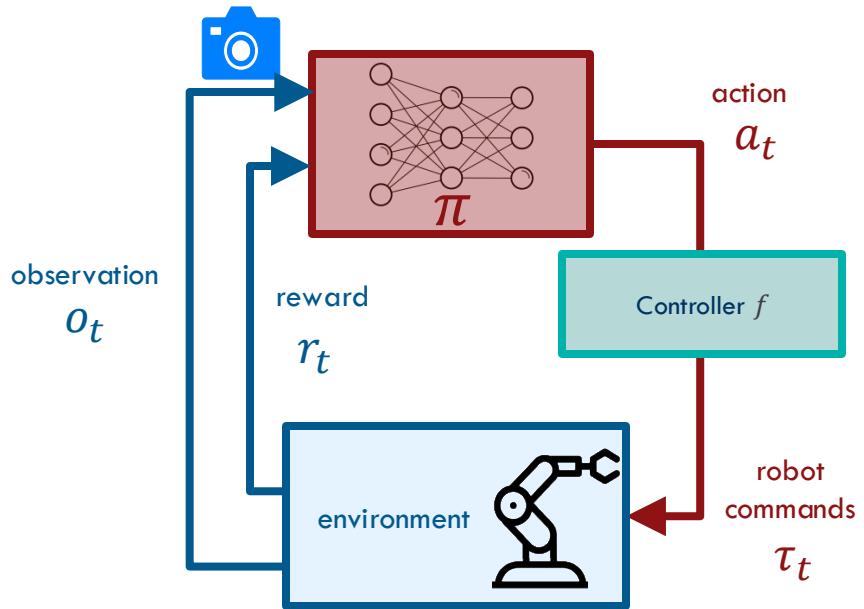
Interaction controller in end-effector space

Direct control of motion and forces of robot's hand



Interaction controller in end-effector space

Direct control of motion and forces of robot's hand



Studied Action Spaces

Action Space	$a = \pi(o) \in \mathcal{A}$
Joint Torques	$a = \tau_{des}$
Joint Velocities	$a = \dot{q}_{des}$
Joint Positions	$a = q_{des}$
Variable Joint Impedance	$a = (q_{des}, k_p, k_v)$
End-Effector Pose	$a = x_{des} \in SE(3)$
VICES	$a = (x_{des}, k_p, k_v)$

$$a = \pi(o) \in \mathcal{A}$$

$$a = \tau_{des}$$

$$a = \dot{q}_{des}$$

$$a = q_{des}$$

$$a = (q_{des}, k_p, k_v)$$

$$a = x_{des} \in SE(3)$$

$$a = (x_{des}, k_p, k_v)$$

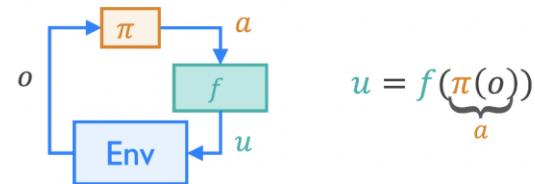
$$u = f(a) = \tau$$

$$u = a$$

$$u = k_v(a - \dot{q}_{cur})$$

$$u = M_j [k_p(a - q_{cur}) - k_v \dot{q}_{cur}]$$

$$u = J^T [M_{ee} [k_p(a - x_{cur}) - k_v \dot{x}_{cur}]]$$



$$u = f(\underbrace{\pi(o)}_a)$$

Experimental Evaluation

- effect of the **choice of action space** in the performance of RL

Experimental Evaluation

- effect of the **choice of action space** in the performance of RL
- **transferability** of policies in different action spaces across robots

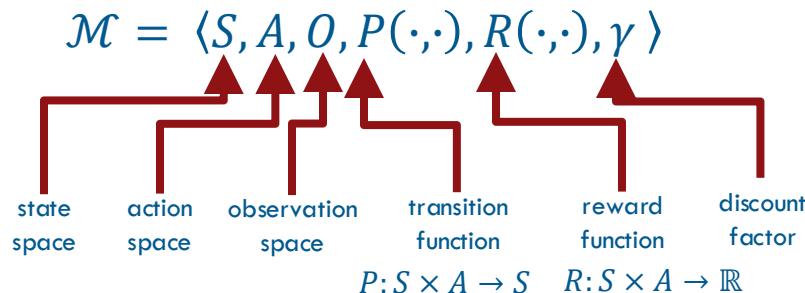
Experimental Evaluation

- effect of the **choice of action space** in the performance of RL
- **transferability** of policies in different action spaces across robots
- benefits of **learning to control forces**

Solving MDPs without Models

Model-free Reinforcement learning

When dynamics model is unknown and hard to estimate, we can **learn policy directly** from the agent's trajectories from interacting with an MDP



agent's experience

$$\tau = \{(s_i, a_i, r_i) \mid i = 0, \dots, H\}$$

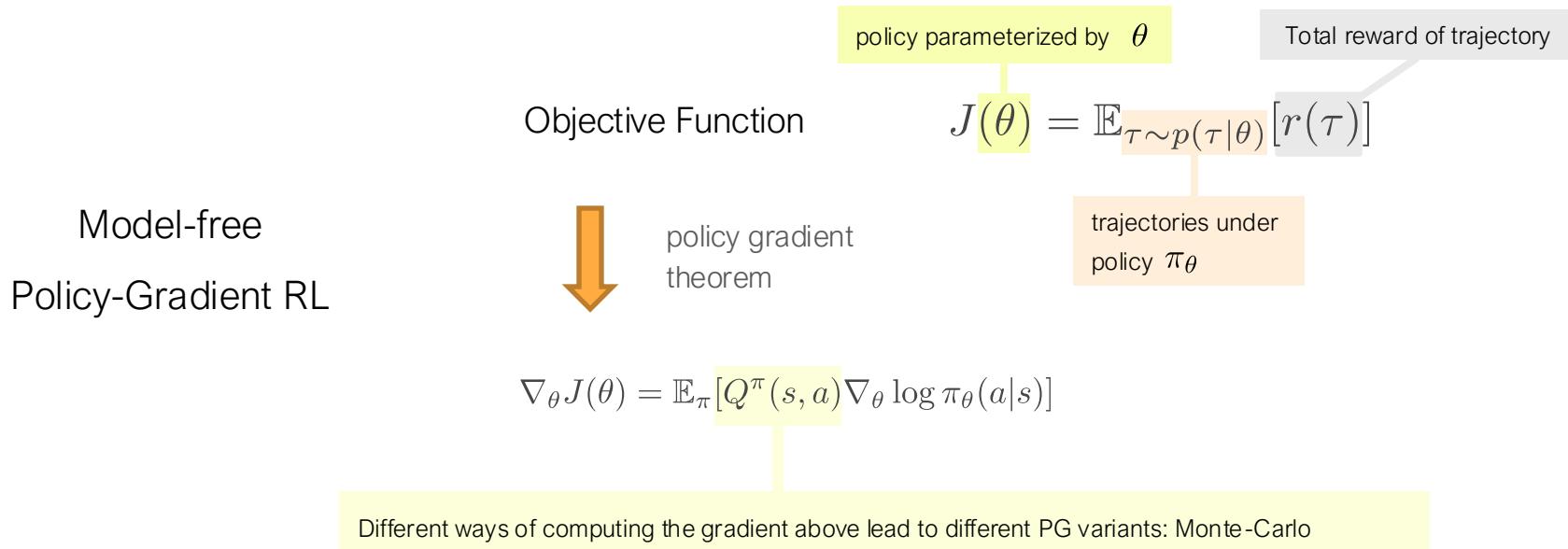


model-free
RL

optimal policy

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r(s_t, \pi(s_t)) \right]$$

Solving MDPs without Models – Policy Learning



Main problems: Data inefficiency and training instability!

Trust Region Methods

TRPO, Trust Region Policy Optimization

- Intuition: the new policy (after update) should not be “too different” from the previous one → Stability
 - TRPO measures how different two policies are by computing the KL divergence and imposing it as hard constraint to the optimization
- Additionally, we can reuse the experiences from the old policy, allowing us to reuse data! → Data efficiency

π_θ

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}(\theta_k, \theta)$$

s.t. $\bar{D}_{KL}(\theta || \theta_k) \leq \delta$

$$\mathcal{L}(\theta_k, \theta) = \underset{s, a \sim \pi_{\theta_k}}{\text{E}} \left[\frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a) \right]$$

$$D_{KL}(\theta || \theta_k) = \underset{s \sim \pi_{\theta_k}}{\text{E}} [D_{KL}(\pi_\theta(\cdot|s) || \pi_{\theta_k}(\cdot|s))]$$

New problem: Computationally expensive

Proximal Policy Optimization

- Convert the constrained optimization problem into a simpler one
 - Use a “clipped” surrogate objective:

$$L(s, a, \theta_k, \theta) = \min \left(\frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a), \text{clip} \left(\frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A^{\pi_{\theta_k}}(s, a) \right)$$

- Implicit constraint:

$$(1 - \epsilon)\pi_{\theta_{\text{old}}} \leq \pi_\theta \leq (1 + \epsilon)\pi_{\theta_{\text{old}}}$$

Improved Stability!

Experimental Setup

Embodiments

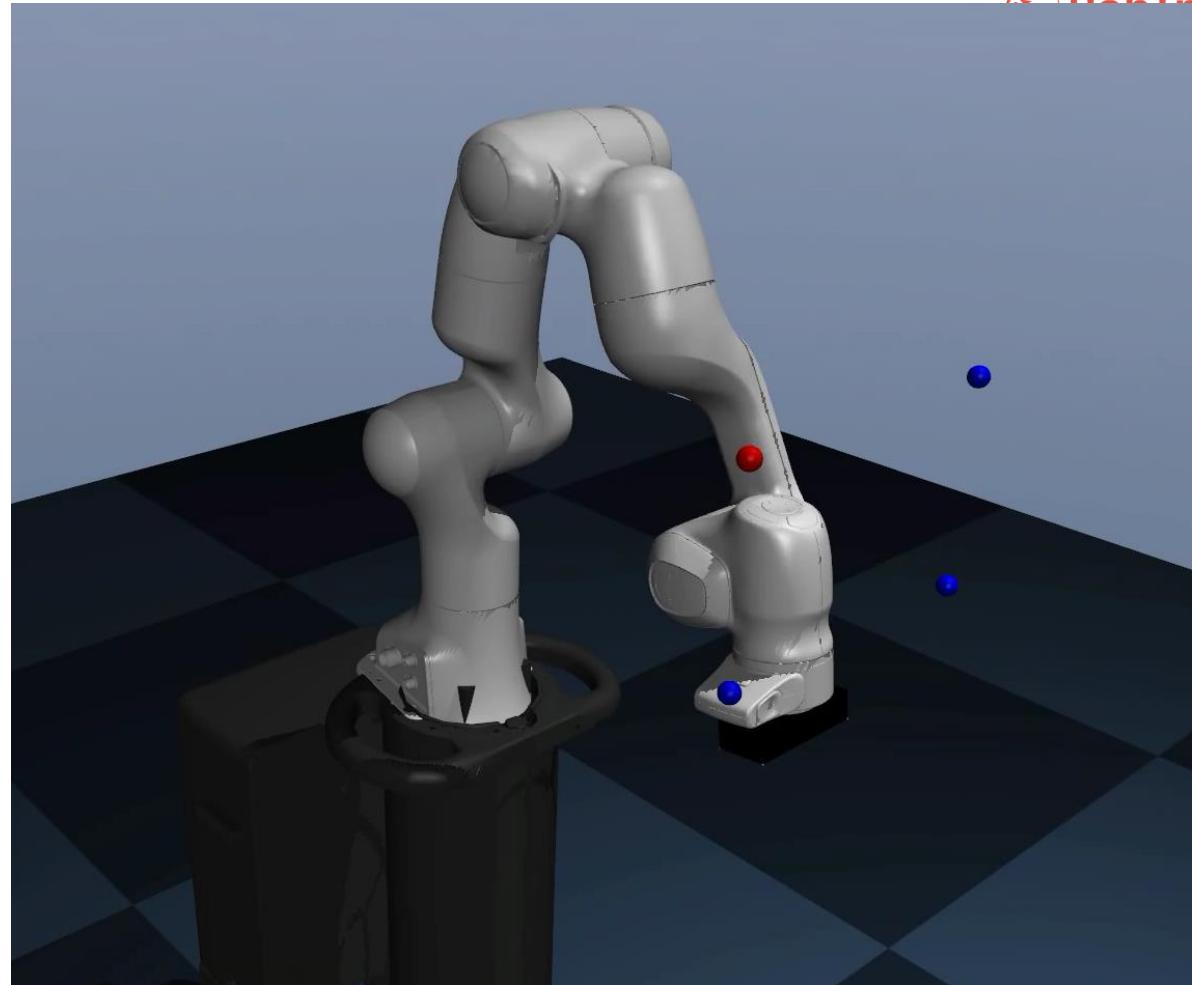


Path Following

No Contact

Input: Proprioception

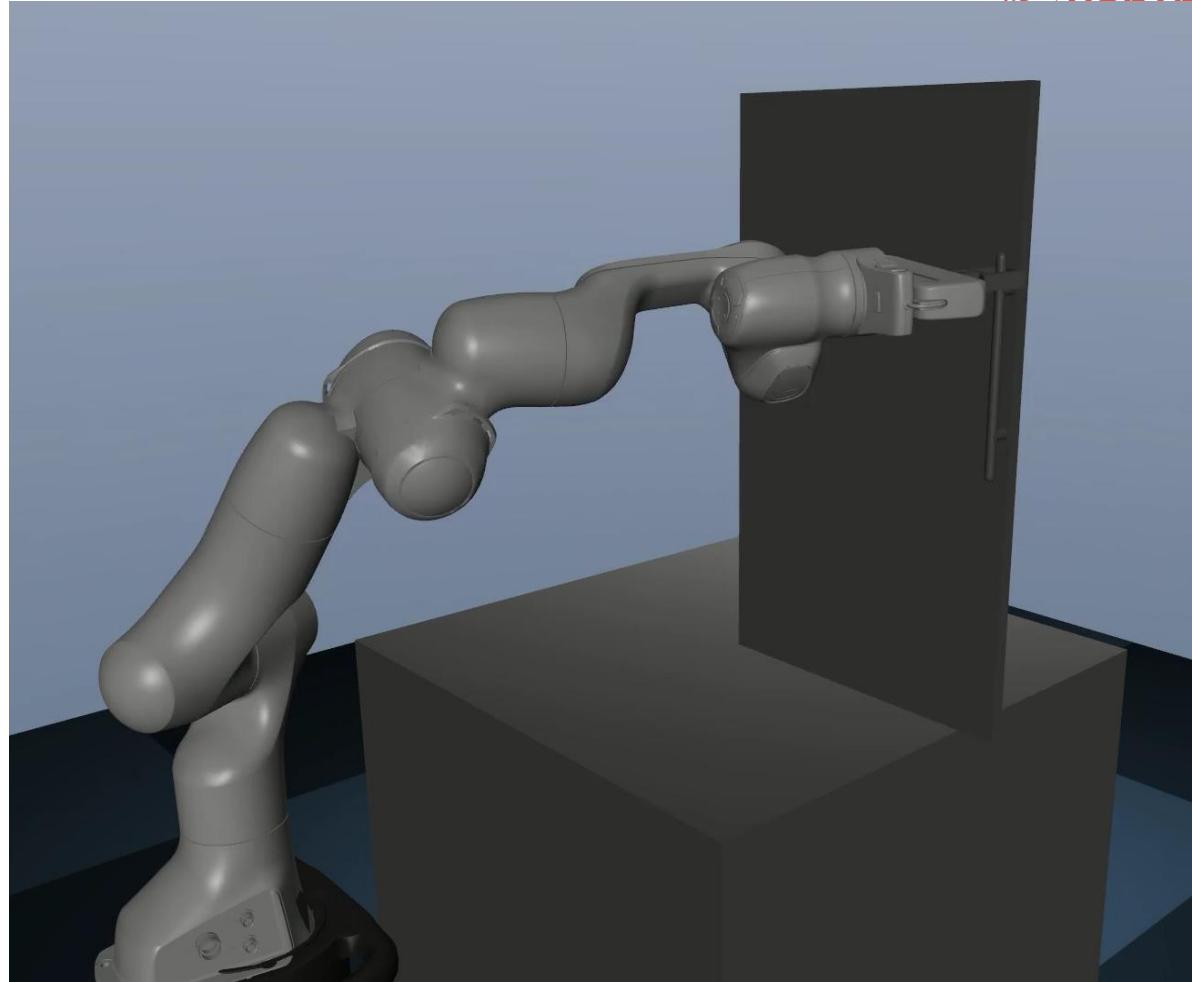
Trajectory State



Door Opening

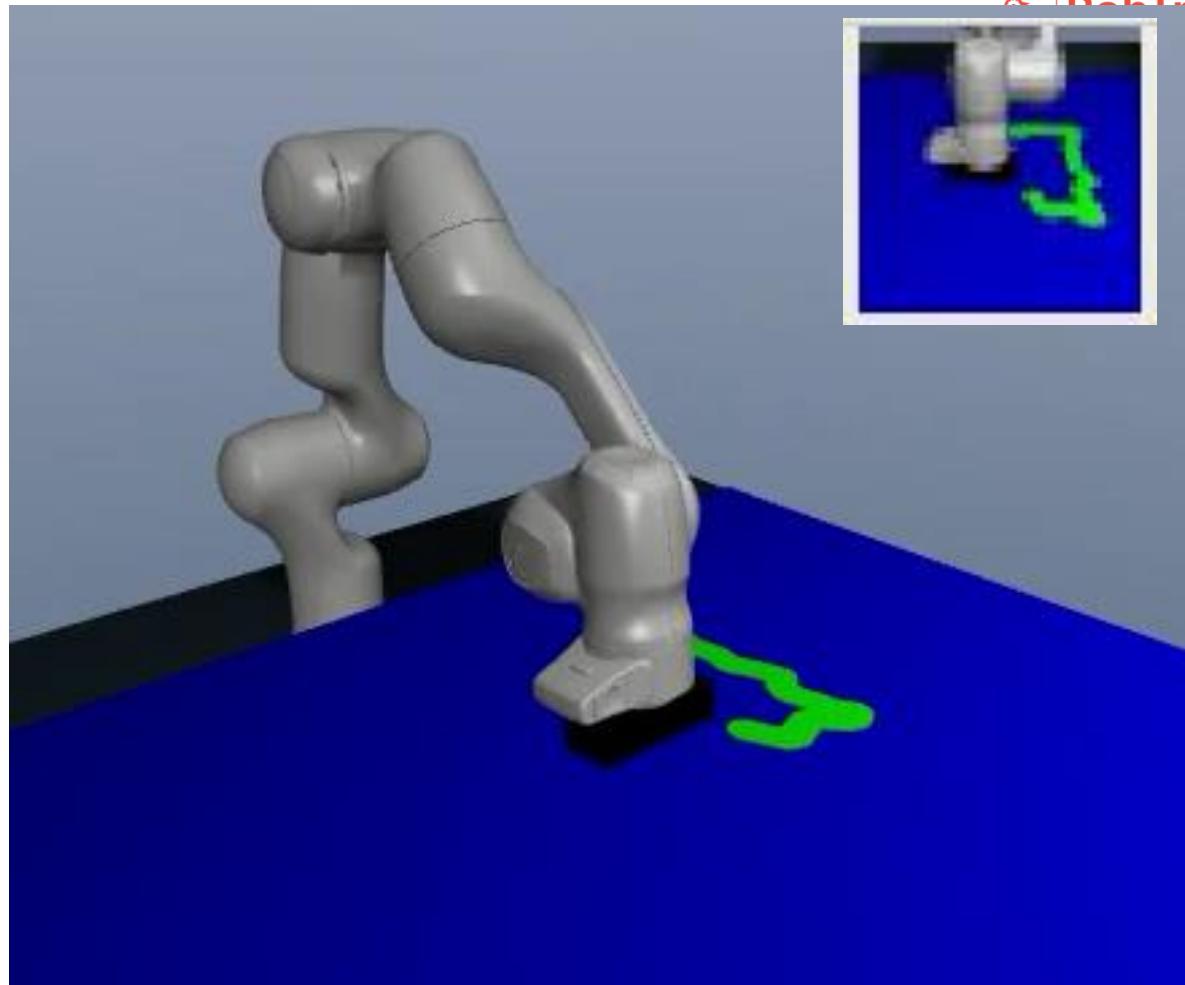
Intermittent Contact

Input: Proprioception and
Door state

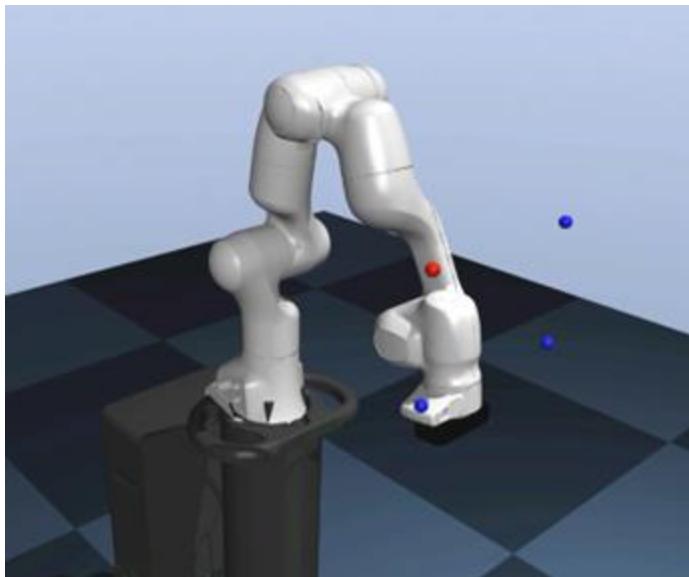


Surface Wiping

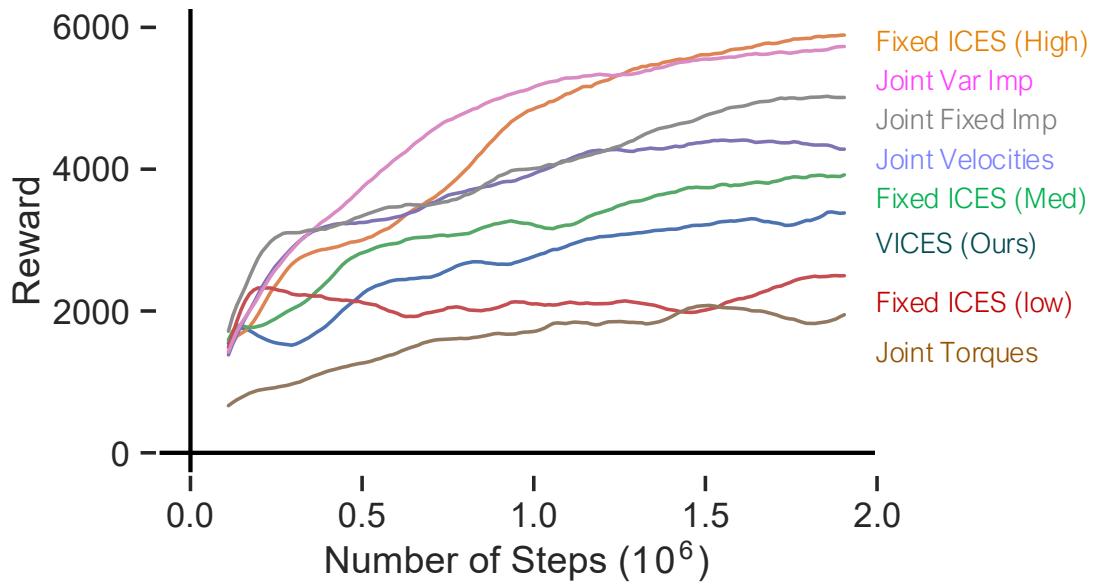
Contact Rich
Input: Proprioception and
Images



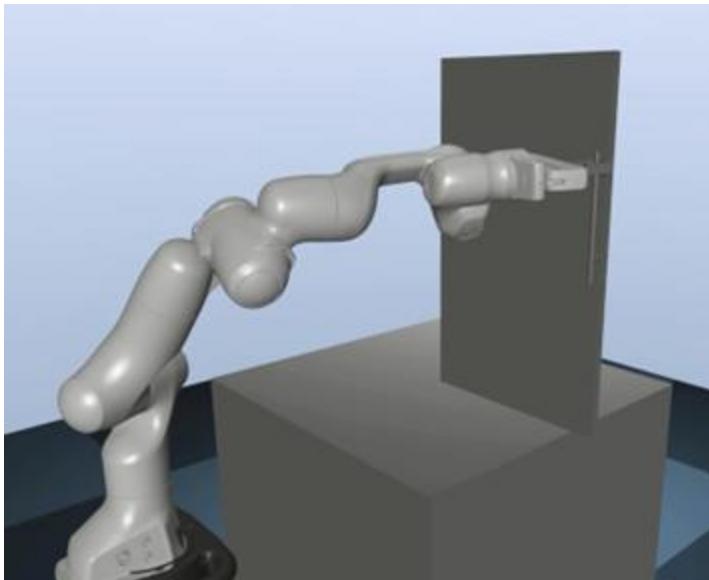
Sample Efficiency: Path Following



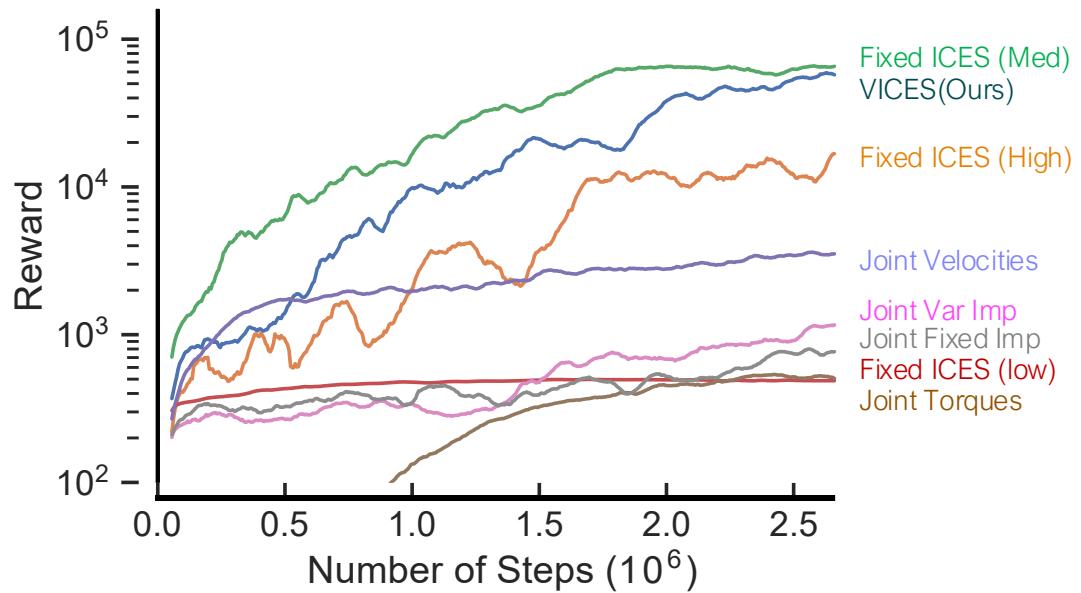
Trained Policy Rollout (VICES, Ours)



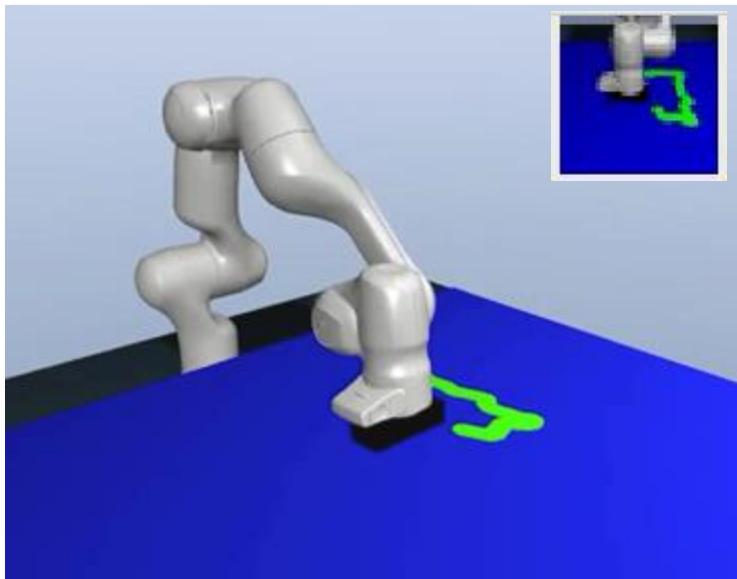
Sample Efficiency: Door Opening



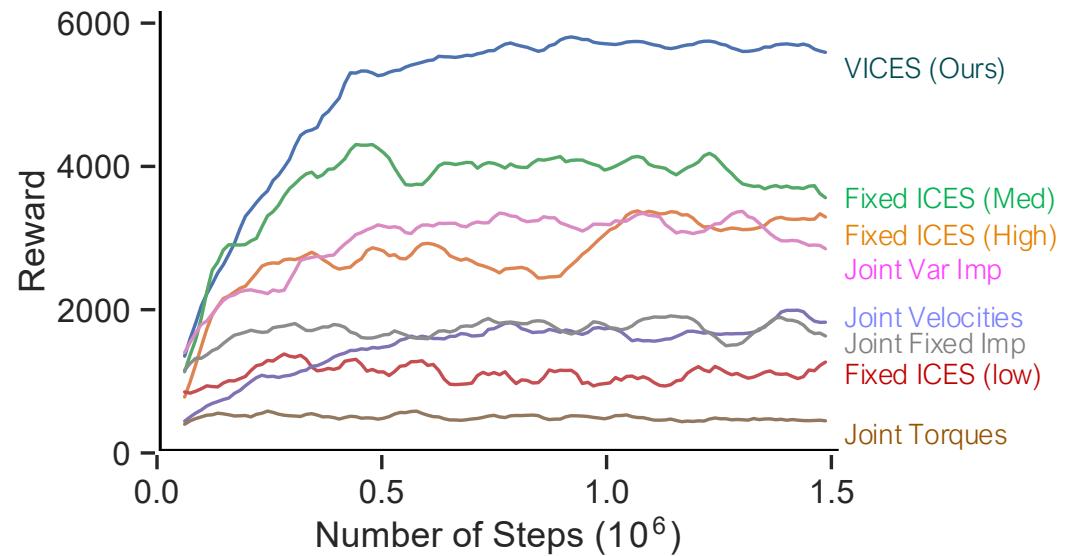
Trained Policy Rollout (VICES, Ours)



Sample Efficiency: Surface Wiping



Trained Policy Rollout (VICES,Ours)



Transfer in Sim: Path Following

Trained on Panda



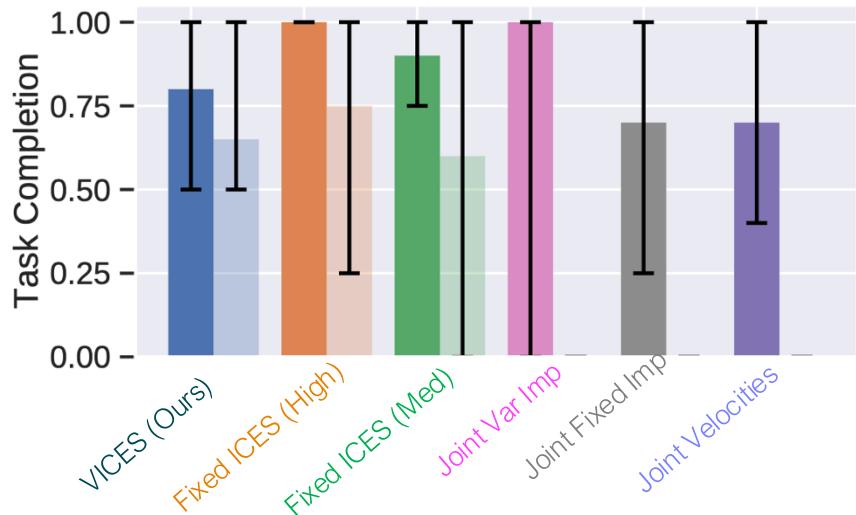
Tested on Sawyer



$$u = f_{\text{panda}}(\pi_{\text{panda}}(o_t))$$

$$u = f_{\text{sawyer}}(\pi_{\text{panda}}(o_t))$$

Trained Policy Rollout (VICES,Ours)

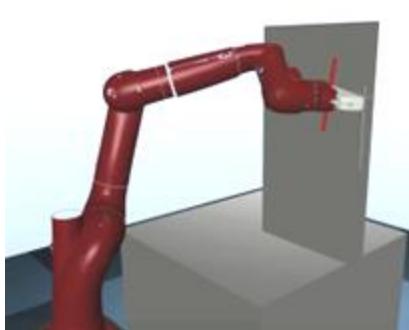


Transfer in Sim: Door Opening

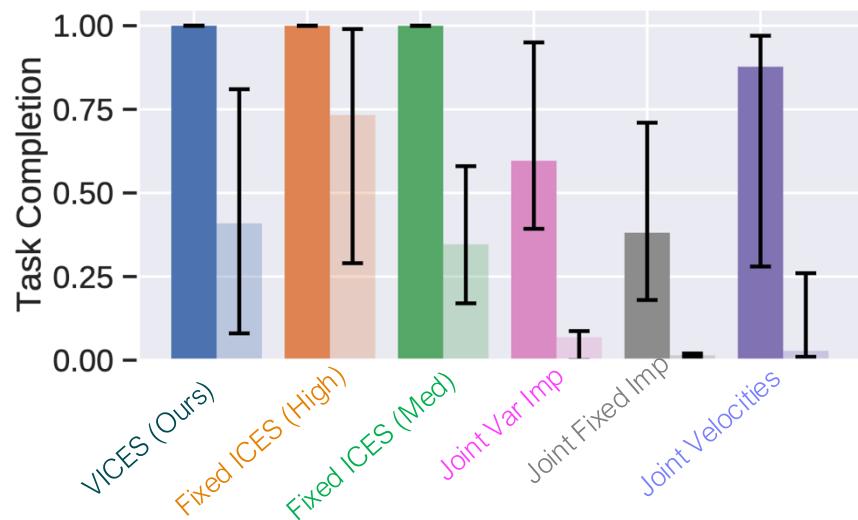
Trained on Panda

Tested on Sawyer

$$u = f_{\text{panda}}(\pi_{\text{panda}}(o_t)) \quad u = f_{\text{sawyer}}(\pi_{\text{panda}}(o_t))$$



Trained Policy Rollout (VICES,Ours)

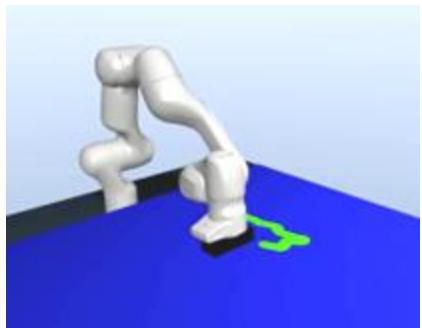


Transfer in Sim: Surface Wiping

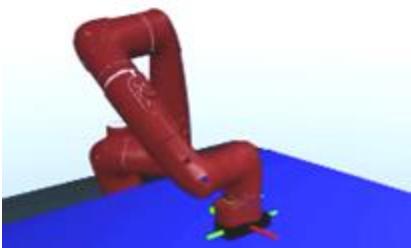
Trained on Panda

Tested on Sawyer

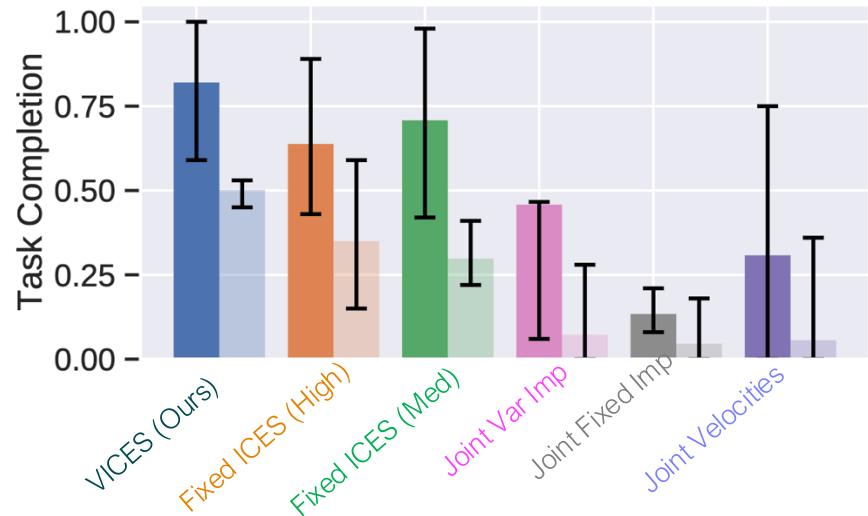
$$u = f_{\text{panda}}(\pi_{\text{panda}}(o_t))$$



$$u = f_{\text{sawyer}}(\pi_{\text{panda}}(o_t))$$



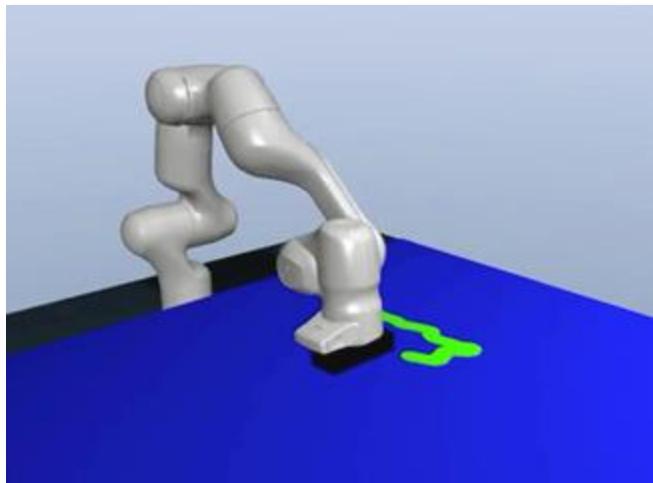
Trained Policy Rollout (VICES,Ours)



Transfer to Real: Surface Wiping

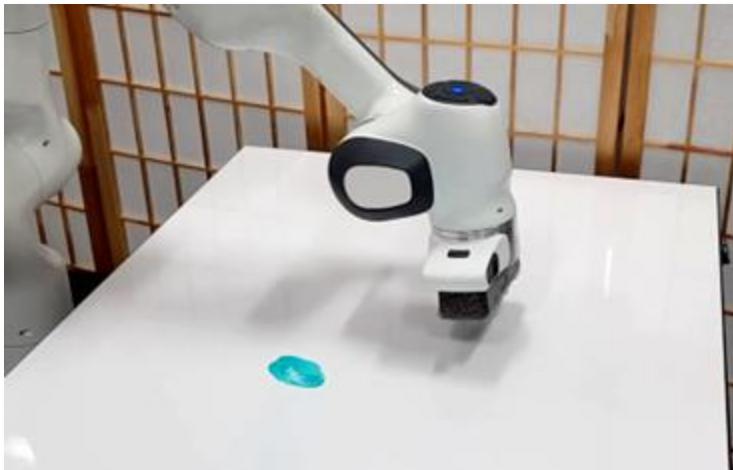
Trained on Simulation

$$u = f_{sim}(\pi_{sim}(o_t))$$



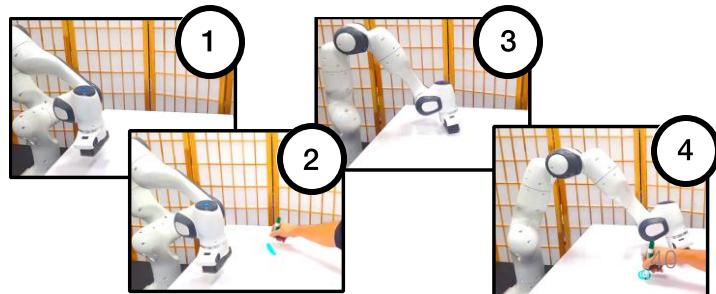
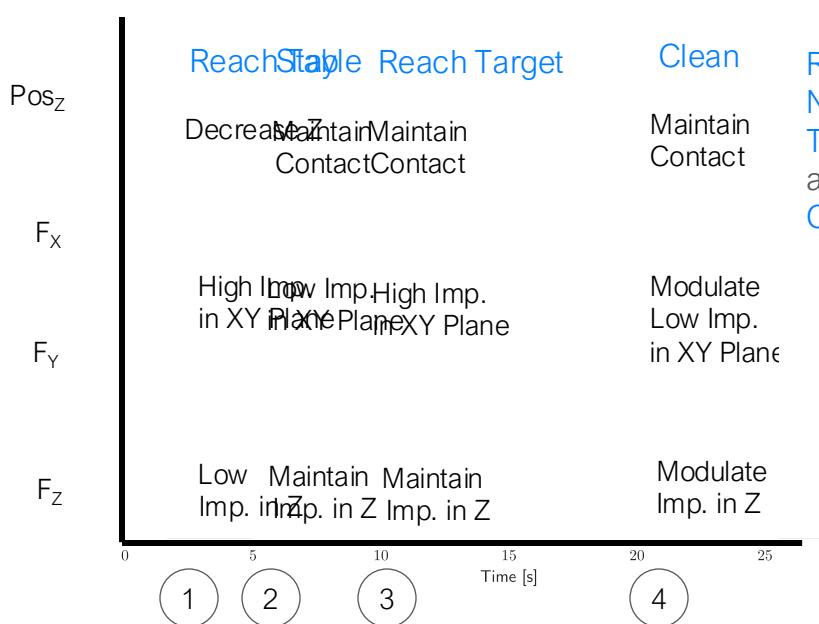
Tested on Real World

$$u = f_{real}(\pi_{sim}(o_t))$$



Success 80% (10 Trials)

Varying Impedance for Wiping



Key Takeaways

- the action space of the RL policy is crucial (more than the algorithm)

Key Takeaways

- the action space of the RL policy is crucial (more than the algorithm)
- transfer between robots is easier in task-space

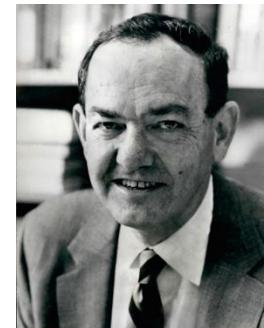
Key Takeaways

- the action space of the RL policy is crucial (more than the algorithm)
- transfer between robots is easier in task-space
- in contact-rich tasks the policy needs to control forces → VICES

Robotics: Representations

“Solving a problem simply means representing it so as to make the solution transparent.”

Herbert A. Simon, Sciences of the Artificial



Usually referring to the state representation but the action representation is key!

Learning action spaces?

LASER: Learning a Latent Action Space for Efficient Reinforcement Learning

Arthur Allshire*,†, Roberto Martín-Martín*,‡, Charles Lin‡, Shawn Manuel‡, Silvio Savarese‡, Animesh Garg†○

Abstract— The process of learning a manipulation task depends strongly on the action space used for exploration: posed in the incorrect action space, solving a task with reinforcement learning can be drastically inefficient. Additionally, similar tasks or instances of the same task family impose latent manifold constraints on the most effective action space: the task family can be best solved with actions in a manifold of the entire action space of the robot. Combining these insights we present LASER, a method to learn latent action spaces for efficient reinforcement learning. LASER factorizes the learning problem into two sub-problems, namely action space learning and policy learning in the new action space. It leverages data from similar manipulation task instances, either from an offline expert or online during policy learning, and learns from these trajectories a mapping from the original to a latent action space. LASER is trained as a variational encoder-decoder model to map raw actions into a disentangled latent action space while maintaining action reconstruction and latent space dynamic consistency. We evaluate LASER on two contact-rich robotic tasks in simulation, and analyze the benefit of policy learning in the generated latent action space. We show improved sample efficiency compared to the original action space from better alignment of the action space to the task space, as we observe with visualizations of the learned action space manifold.

Additional details: pair.toronto.edu/laser

I. INTRODUCTION

Deep Reinforcement Learning (RL) has fueled rapid progress in robot manipulation by enabling learning of closed loop visuomotor control policies that integrate perception and control in a single system [1]. However, the focus of end-to-end policy learning has been on the complexity of

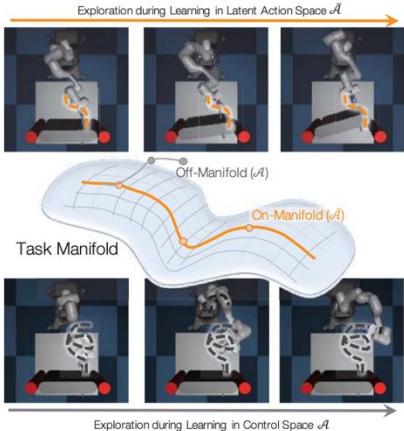
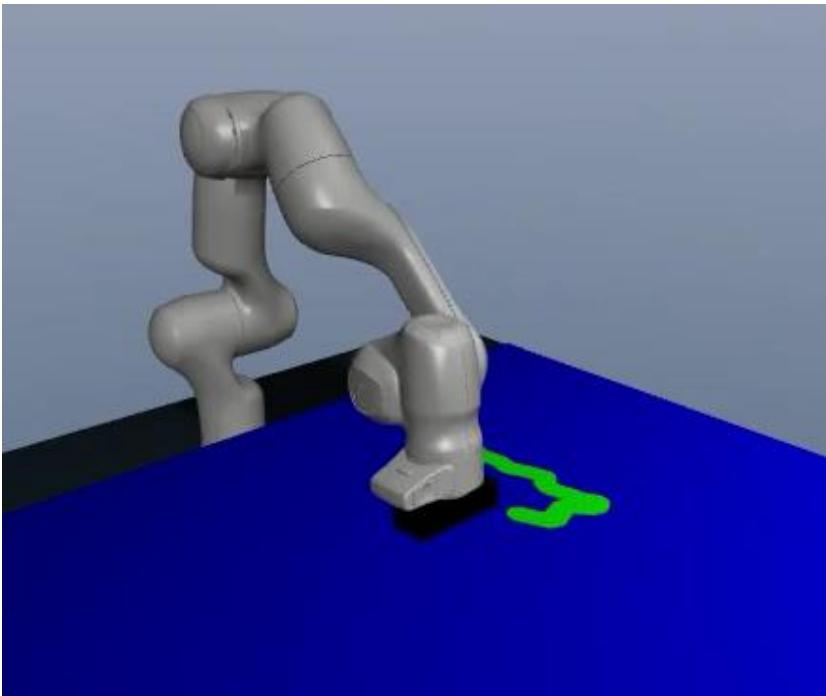


Fig. 1: Learning Latent action spaces for efficient reinforcement learning. Manipulation tasks, such as opening a door, are often structured and do not require exploration in the entire action space, only on certain manifold. LASER learns this action space manifold from data, either offline (expert) or online (training with LASER), enabling faster learning in subsequent novel instances of the task by transferring the knowledge via an efficient latent action space.

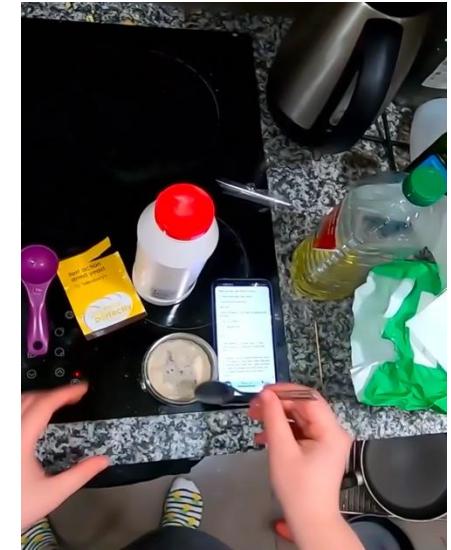
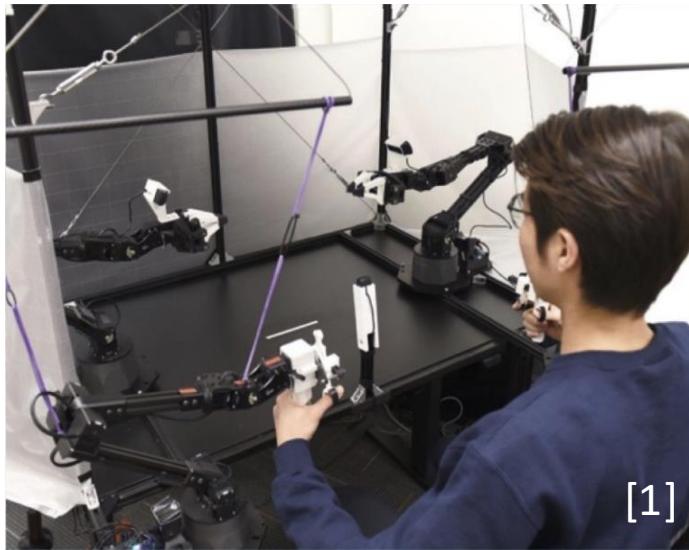
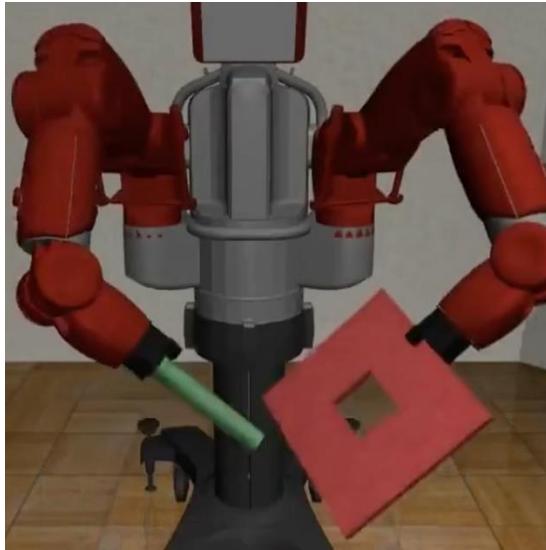
- Can we directly learn good action spaces?

Next Step: Generalizing Contact-Rich Manipulation



• • •





Exploring to Learn

✗ Large action space

Using Robot Demos to Learn

✗ Difficult to collect

Learning from Human Demonstration

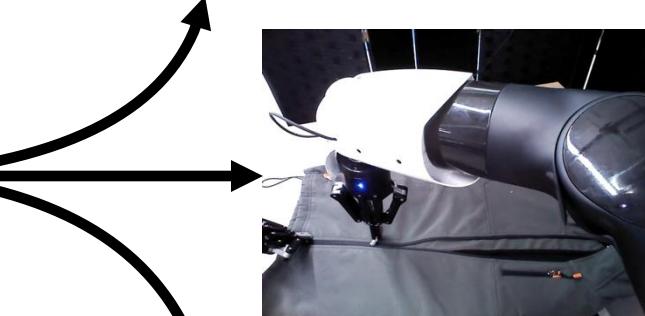
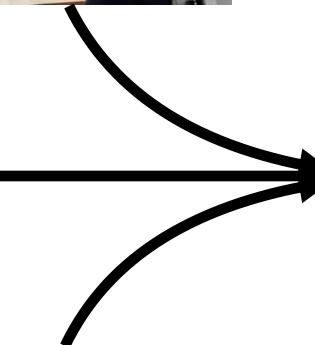


Bimanual Manipulation as Screw Motions



- ✓ Helps to cleanly parse human demonstrations for bimanual tasks
- ✓ Provides a more efficient action space for exploration

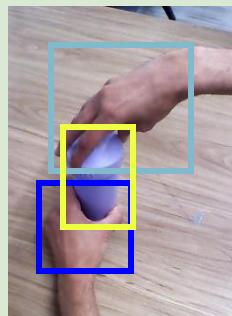
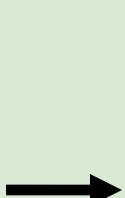




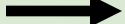
Extracting Screw Action from a Human Demonstration



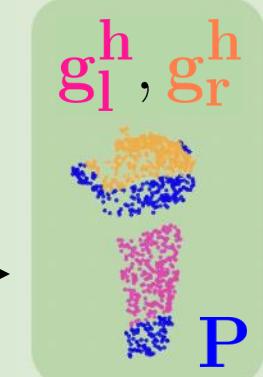
Human Video
Demo (RGBD)



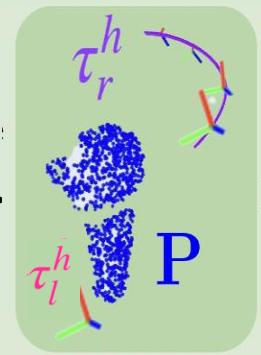
Hand-Object Detection



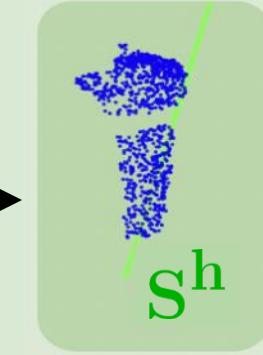
Hand Pose Detection



Contact Point Extraction

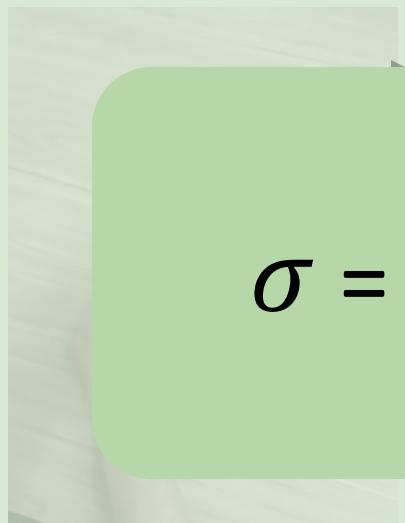


Wrist Pose Extraction



Axis Extraction

Extracting Screw Action from a Human Demonstration



Human Video
Demo (RGBD)



Hand Pose Detection



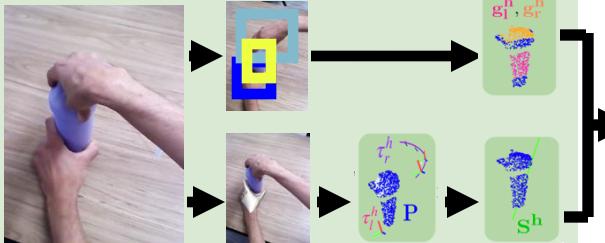
Wrist Pose Extraction



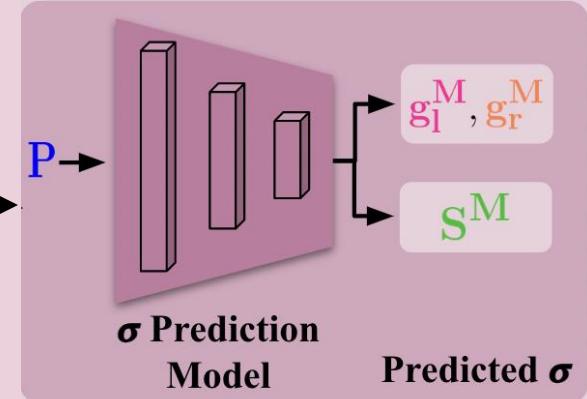
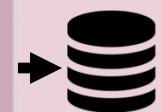
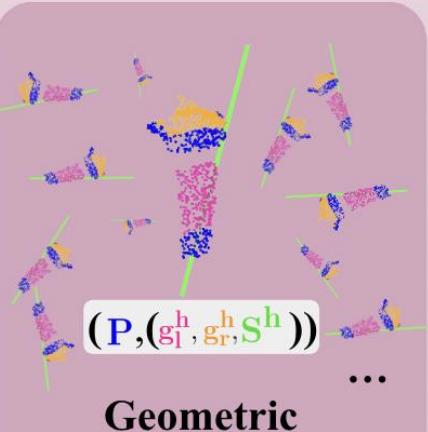
Axis Extraction

$$\sigma = (\quad g_l^h, g_r^h \quad S^h \quad \tau_l^h \quad)$$

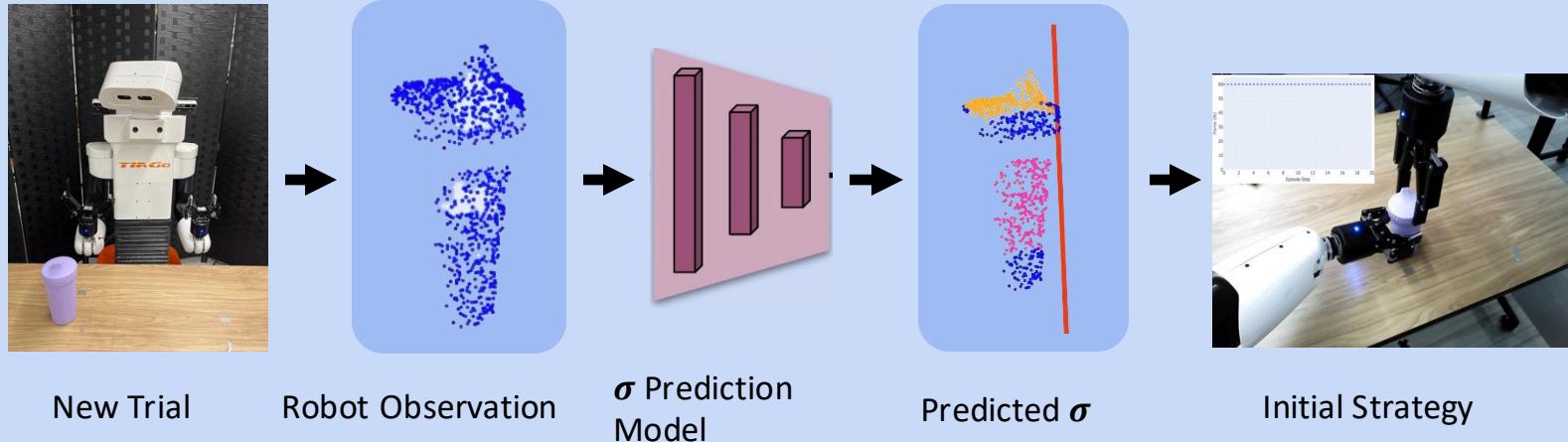
Extracting Screw Action from a Human Demonstration

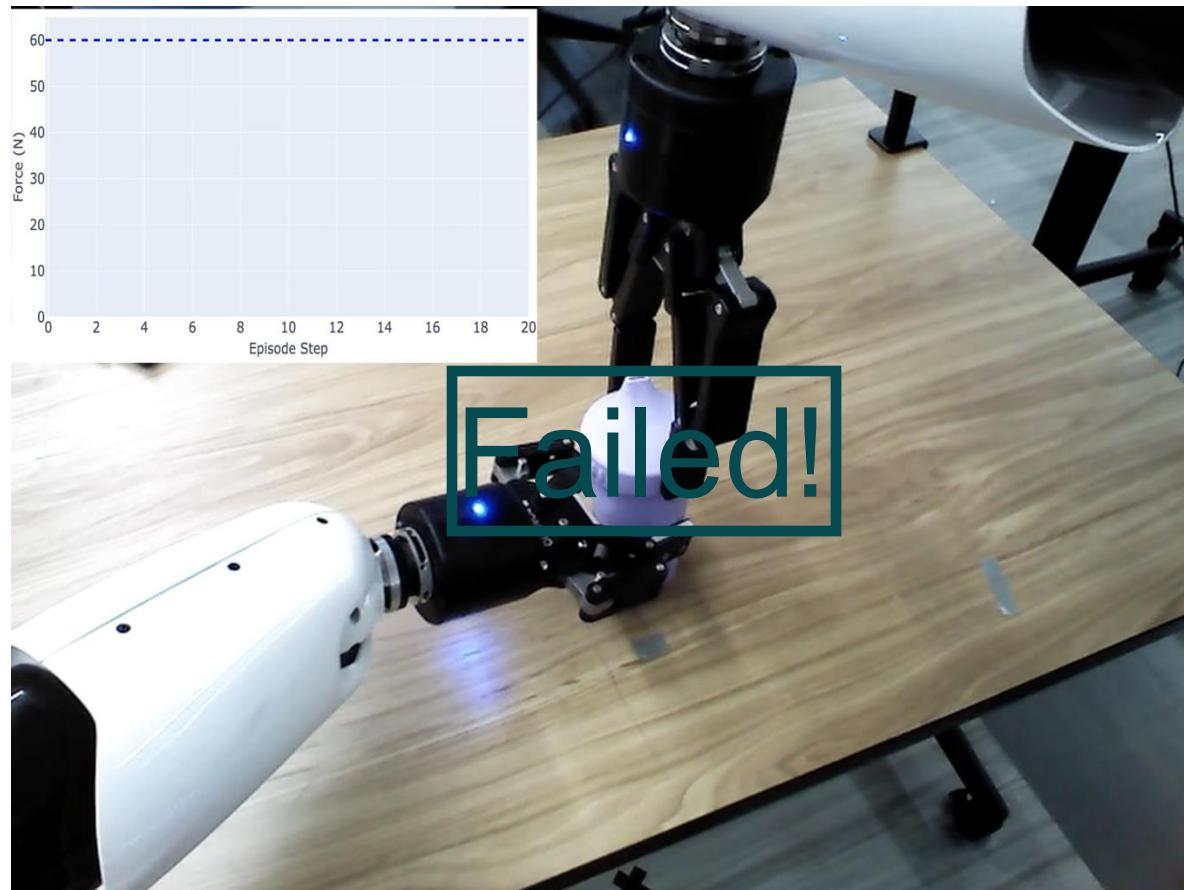


Learning to Predict a Screw Action from a Point Cloud

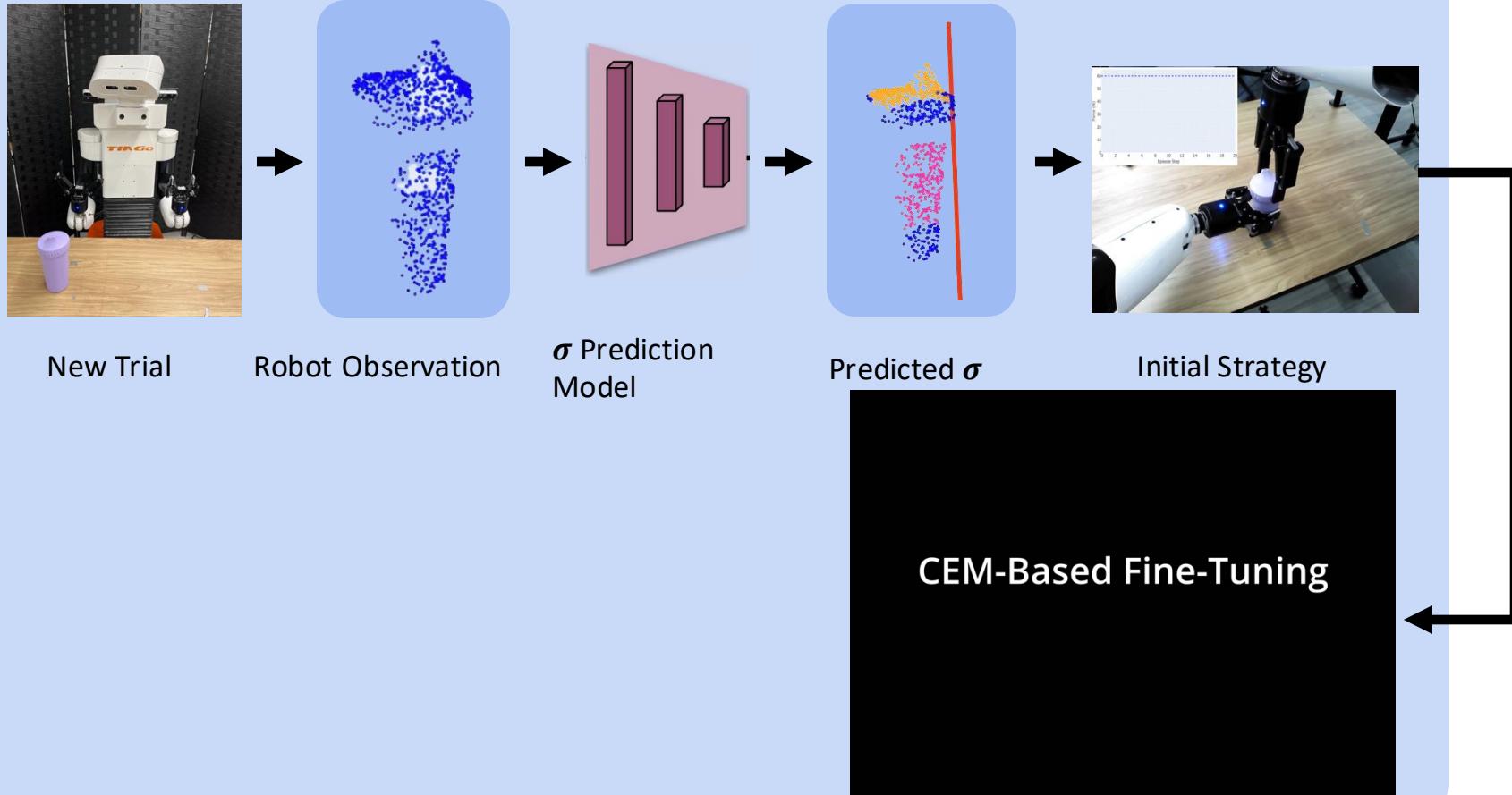


Self-Supervised Screw-Action Policy Fine-Tuning



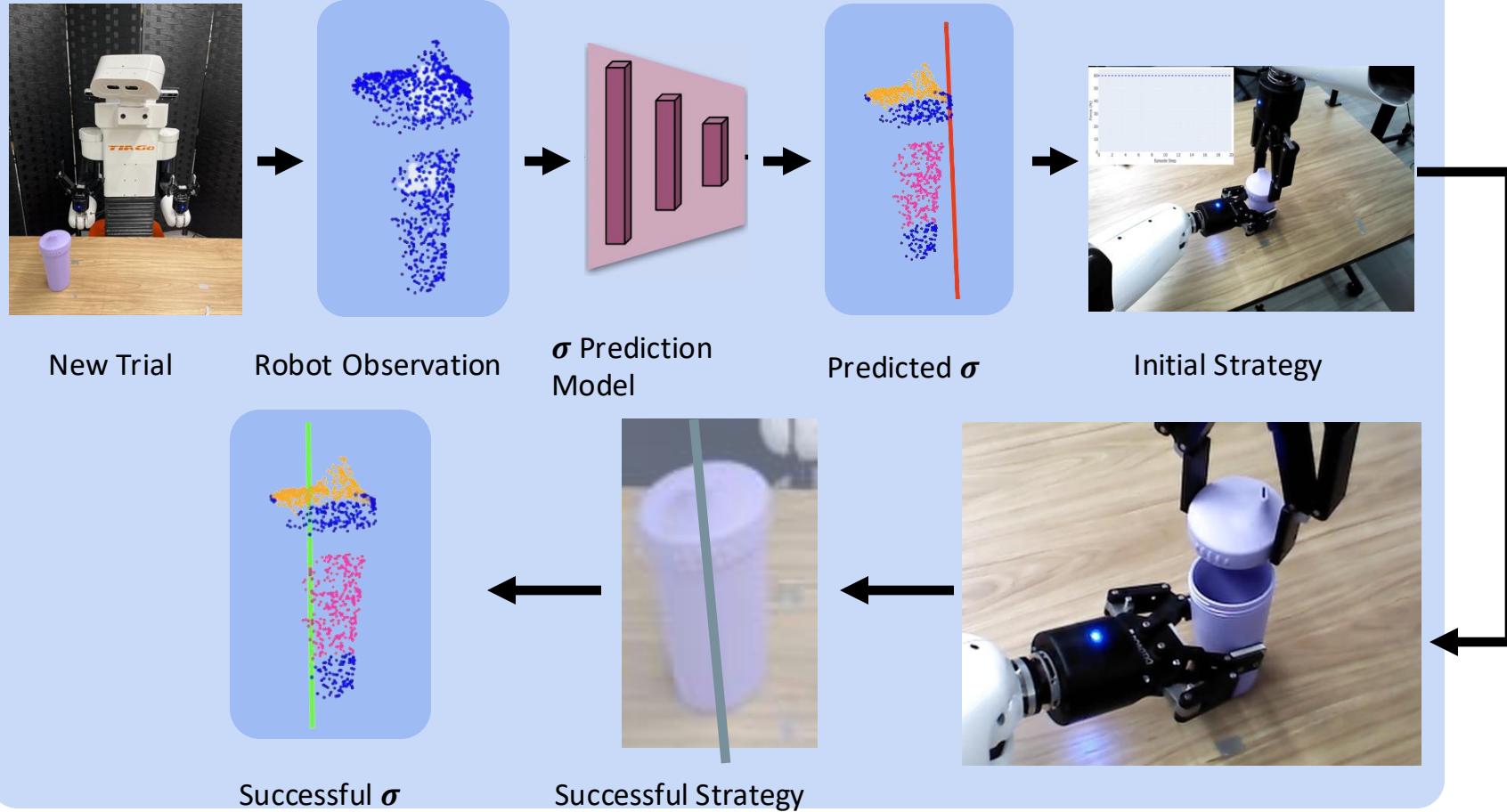


Self-Supervised Screw-Action Policy Fine-Tuning

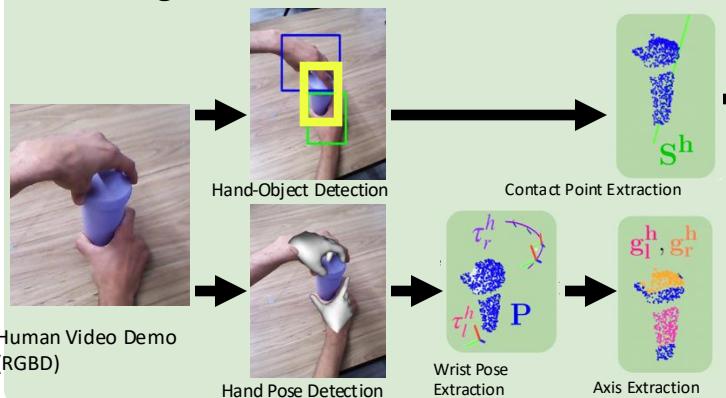


CEM-Based Fine Tuning
Success!

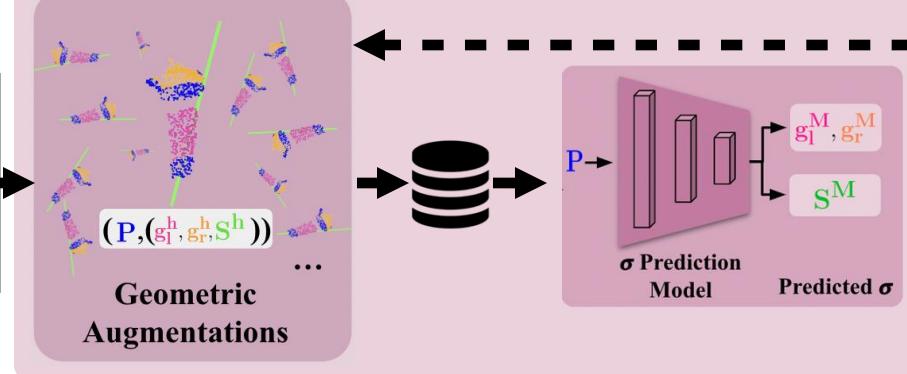
Self-Supervised Screw-Action Policy Fine-Tuning



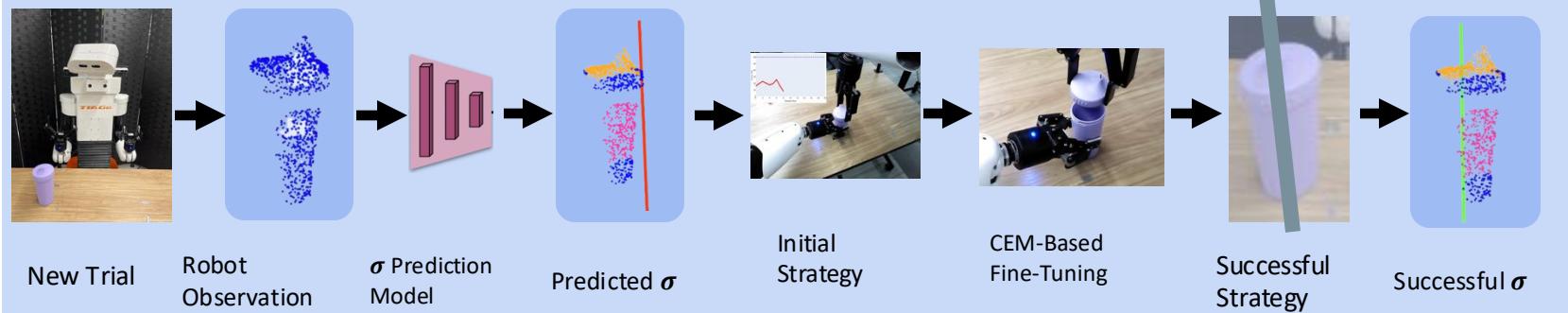
Extracting Screw Action from a Human Demo



Predicting a Screw Action from a Point Cloud



Self-Supervised Screw-Action Policy Fine-Tuning



Open Bottle



Close Zipper



Insert Roll



Close Laptop



Stir



Cut

Success

Initial
Strategy

Open Bottle

Success

Epoch 0

Success

Epoch 0

Close Laptop

Close Zipper

Success

Epoch 0

Success

Epoch 0

Stir

Insert Roll

Success

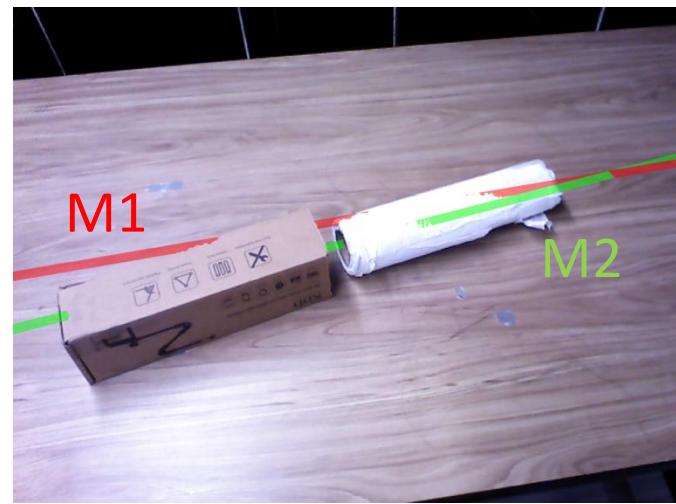
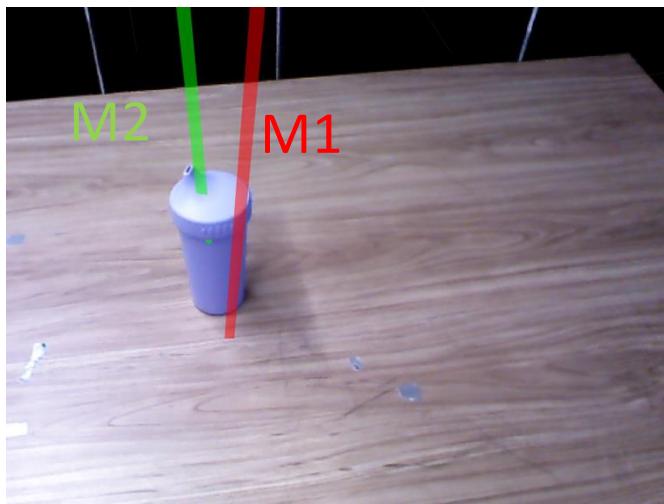
Epoch 0

Success

Epoch 0

Cut

Does the screw action prediction model improve from the interactions?



Can ScrewMimic generalize to new objects?

Success

Epoch 0

Success

Epoch 0

Does the Screw Action generalize when there is left-hand motion?

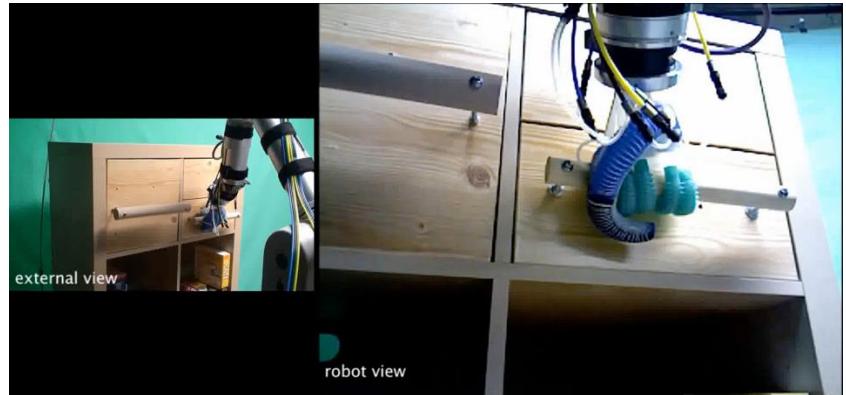
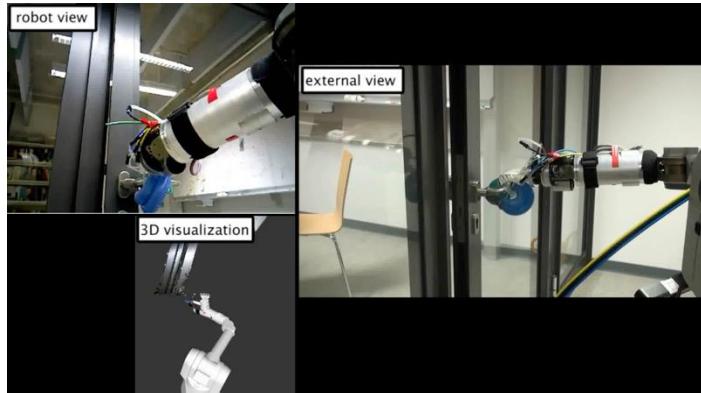


Key Takeaways of ScrewMimic



- Many bimanual contact-rich tasks can be represented by a simple 1-DoF screw joint
- Our screw action representation helps in more accurate human demo interpretation and efficient exploration
- ScrewMimic can continually expand its manipulation capabilities to new objects

The role of embodiment in contact-rich manipulation



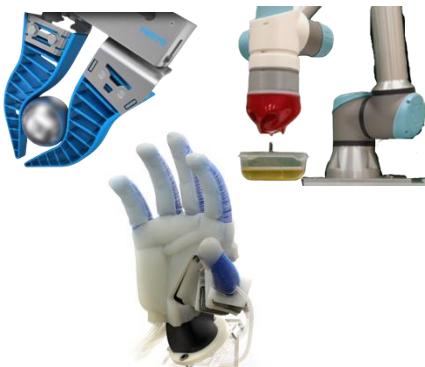
Existing robotic hand designs

Covering both extremes of the spectrum rigid-soft

✓ Compliance

✓ Collision Tolerance

✗ Precision



BaRiFlex



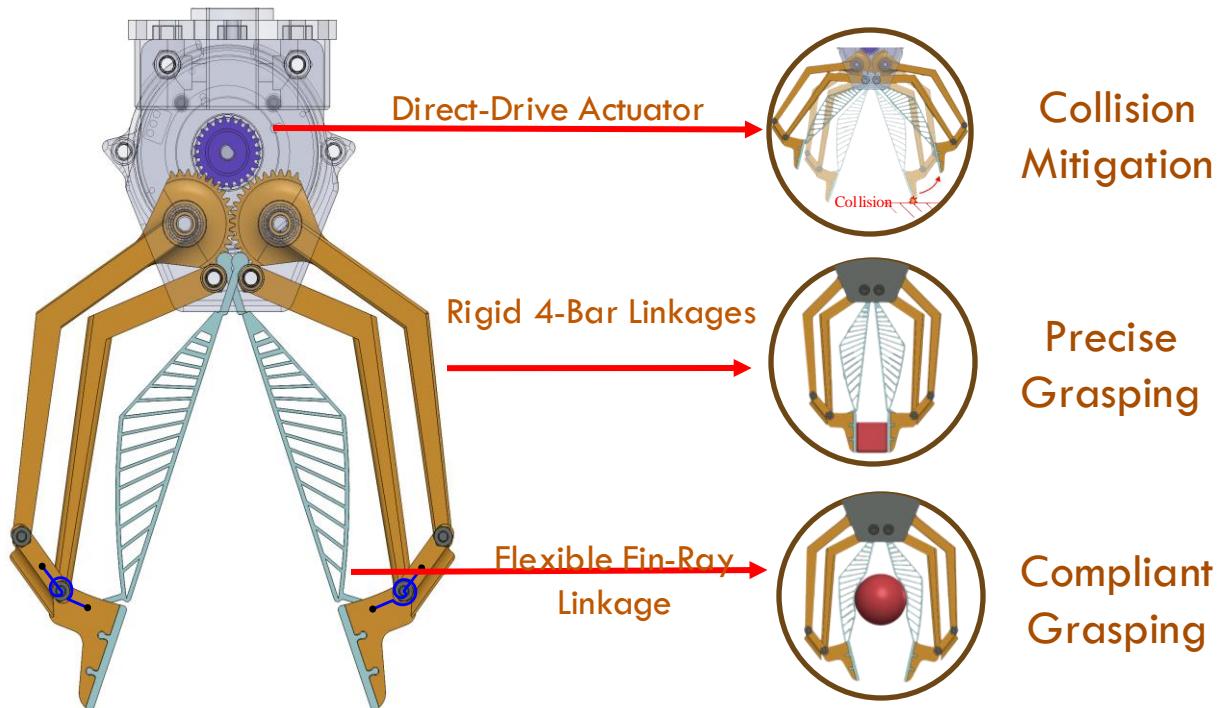
✗ Compliance

✗ Collision Tolerance

✓ Precision

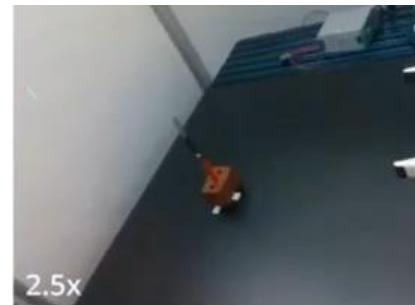
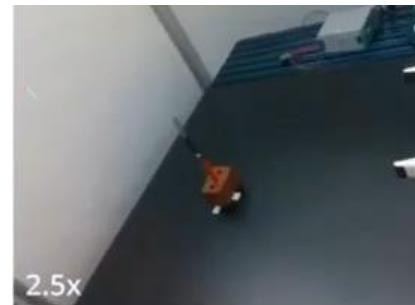
✓ Strength

BaRiFlex Design



less than \$500!!

Grasping Versatility Test (BaRiFlex)



Task Versatility Test



Manipulating Heavy Objects

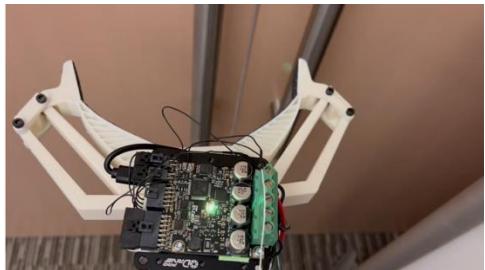


High-Speed Catching

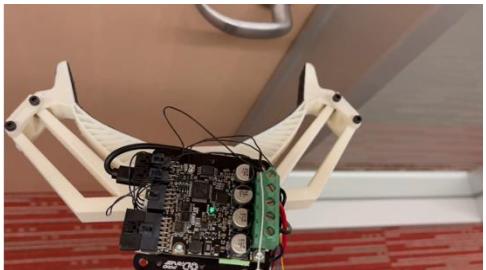


Long-Horizon Manipulation

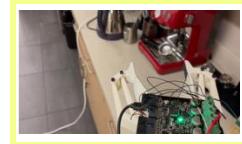
Task Versatility Test



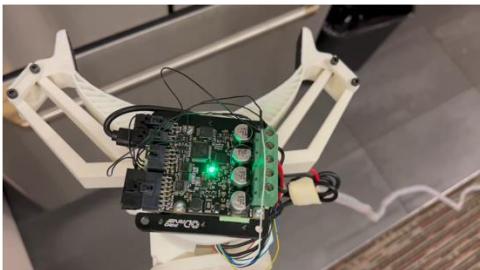
Open Sliding Door



Open Door with Lever Handle
Manipulating Heavy Objects



Long-Horizon Manipulation



Open Fridge



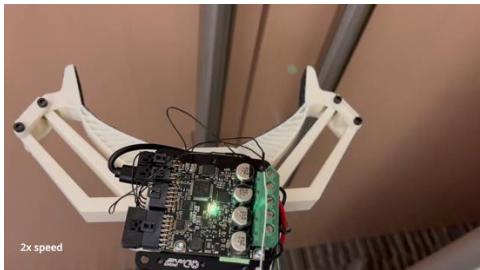
Catch Tennis Ball



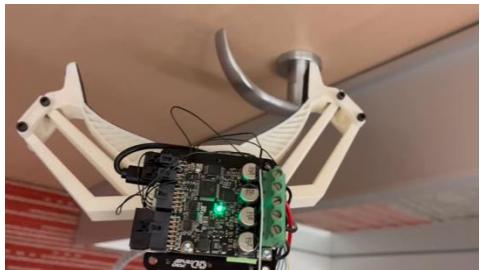
Catch Softball

High-Speed Catching

Task Versatility Test



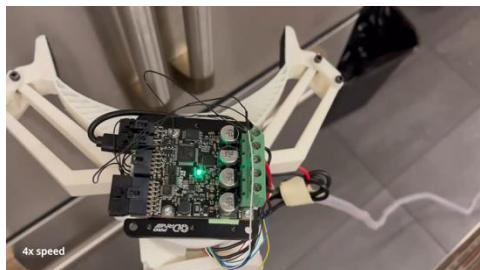
Open Sliding Door



Open Door with Lever Handle
Manipulating Heavy Objects



Long-Horizon Manipulation



Open Fridge



Catch Tennis Ball



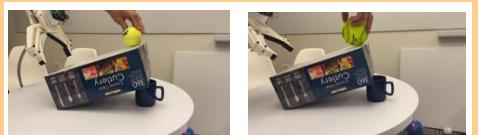
Catch Softball

High-Speed Catching

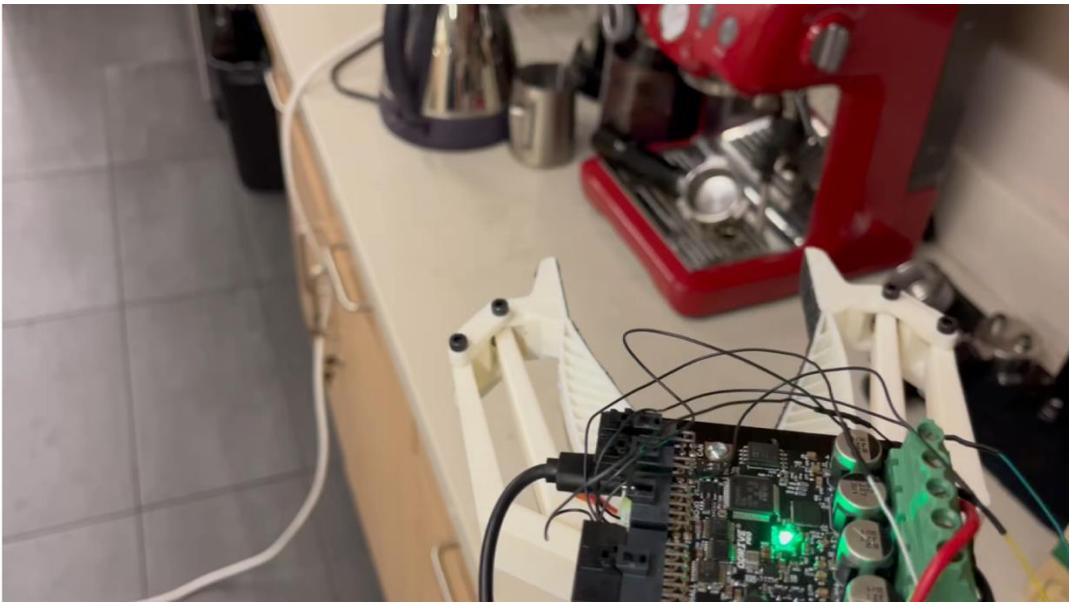
Task Versatility Test



Manipulating Heavy Objects



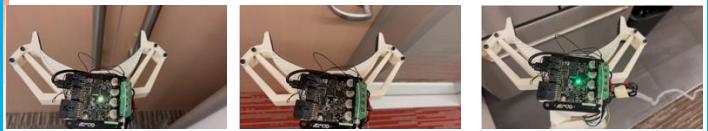
High-Speed Catching



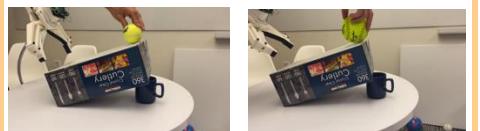
Make Instant Coffee

Long-Horizon Manipulation

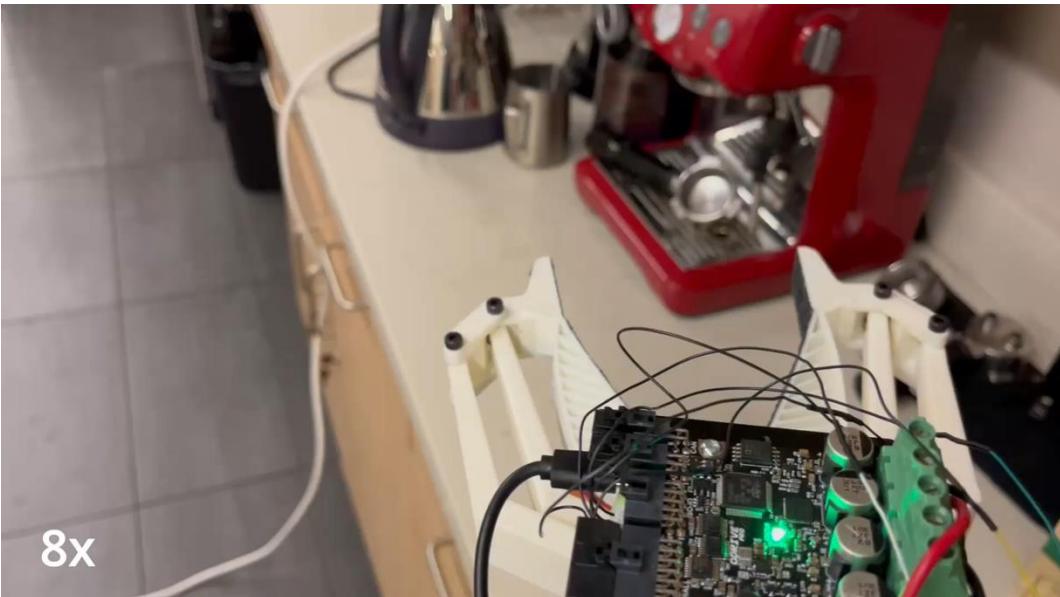
Task Versatility Test



Manipulating Heavy Objects



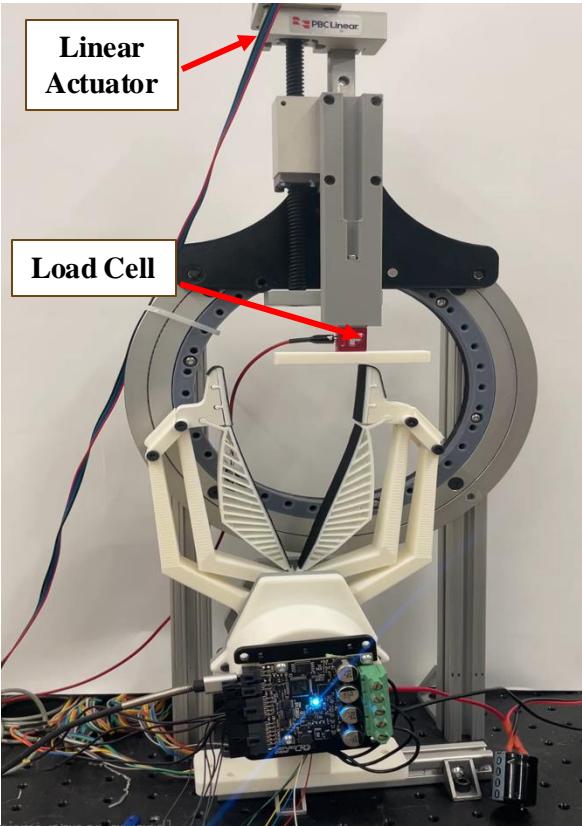
High-Speed Catching



Make Instant Coffee

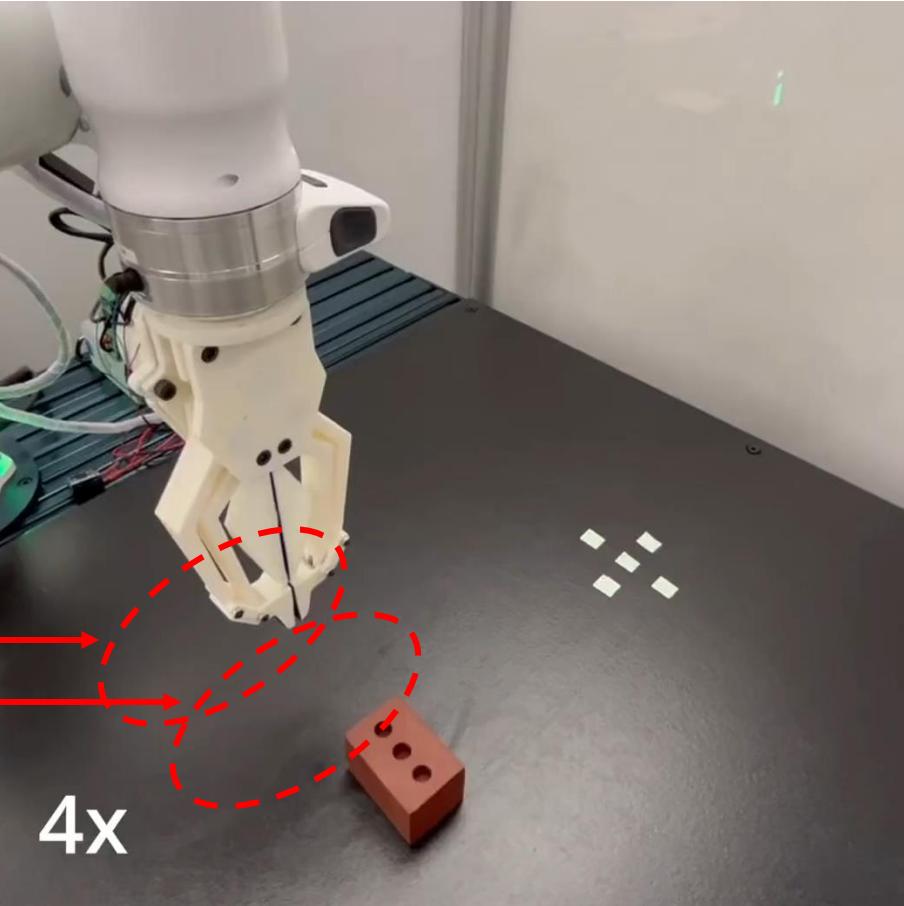
Long-Horizon Manipulation

Robustness and Durability Tests



Real World Reinforcement Learning Test

Collision!
Collision!



Gu-Cheol
Jeong



Arpit Bahety



Ashish
Despande

Final Thoughts

- Last challenge in robotics is to control contact-rich interactions
- Variable impedance facilitates learning to control contact-rich manipulation
- Bimanual manipulation of contact-rich tasks can be approximated by 1DoF screw axes





Thank you!

robertomm@cs.utexas.edu



RobIn
ROBOT INTERACTIVE
INTELLIGENCE LAB

