



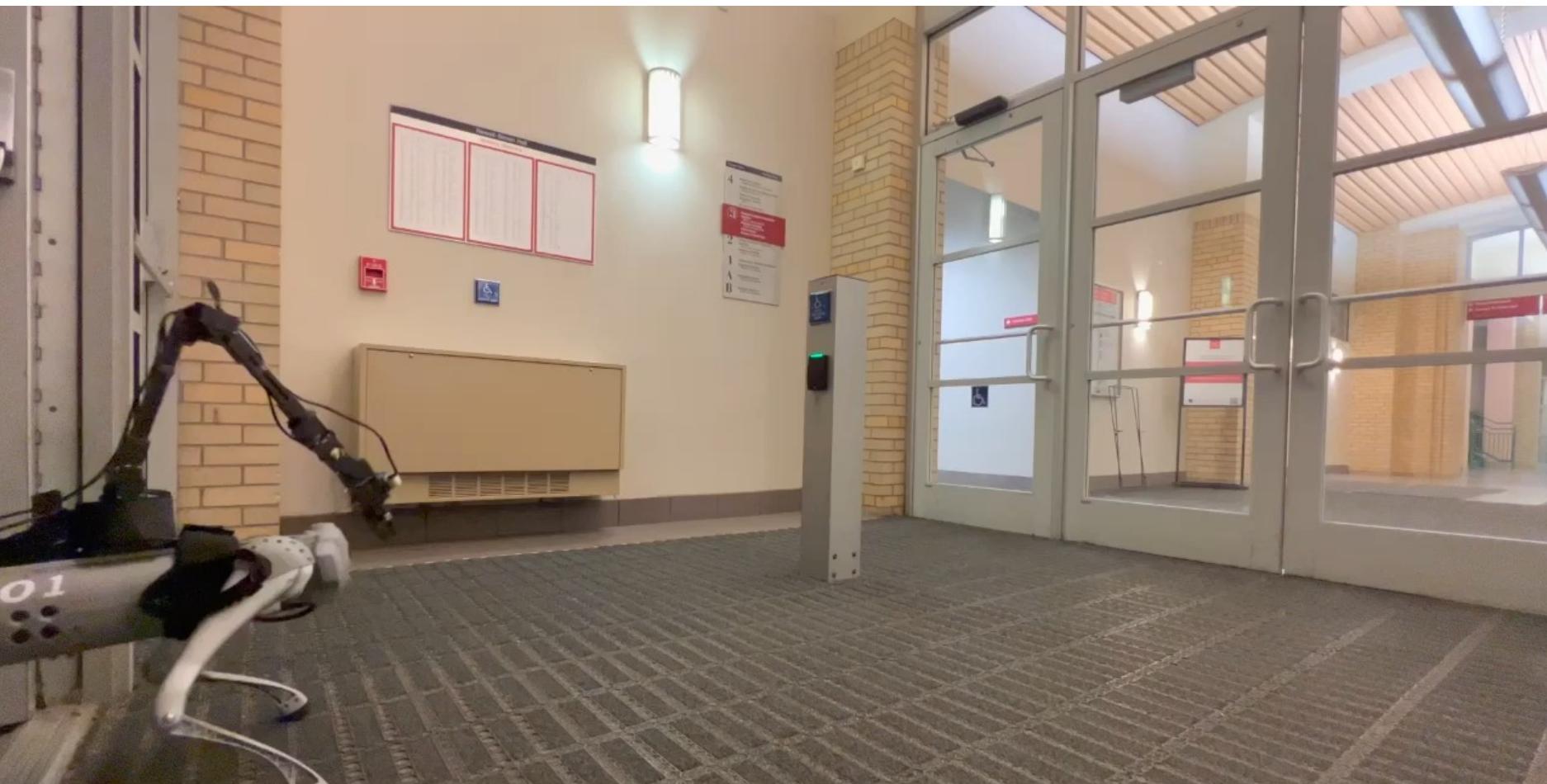
UC San Diego

CSE-291

# Boosting Reinforcement Learning and Planning with Demonstrations: A Survey

Present by Tongzhou Mu

# RL and Planning in Deep Learning Era



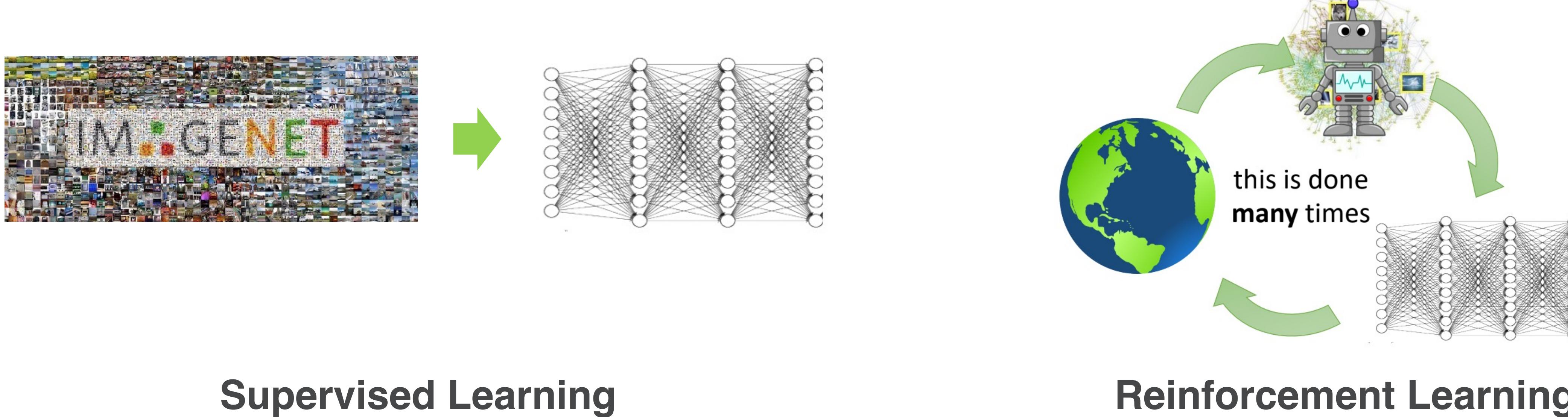
Pressing door-open button to go inside the building



ChatGPT: Optimizing  
Language Models  
for Dialogue

# RL Can Be Costly

- Deep RL algorithms usually require tremendous training samples
  - Impractical or inefficient in complex environments

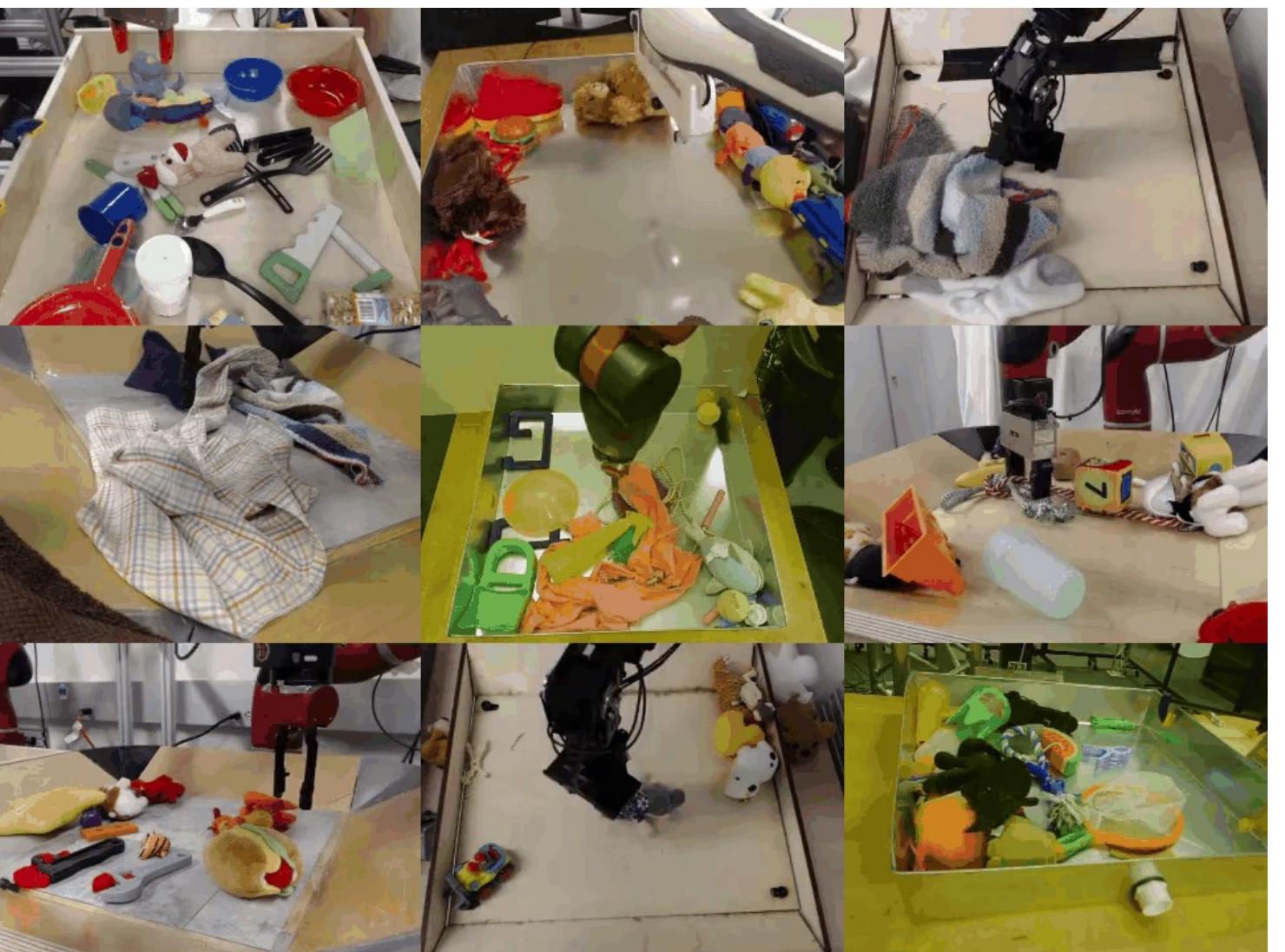


- Supervised Learning**
- Dataset is collected beforehand
  - Fit the labeled data

- Reinforcement Learning**
- Dataset is collected during interaction
  - Find a good policy by trial-and-error

# Datasets for Decision Making

- Pre-collected demonstrations in many domains
  - Robotics
  - Autonomous Driving
  - ...



RoboNet



Waymo Open Dataset

# Topics of This Talk

- Key questions to discuss
  - How to utilize demonstrations in RL and planning?
  - How to collect demonstrations that are useful for RL and planning?

# Outline

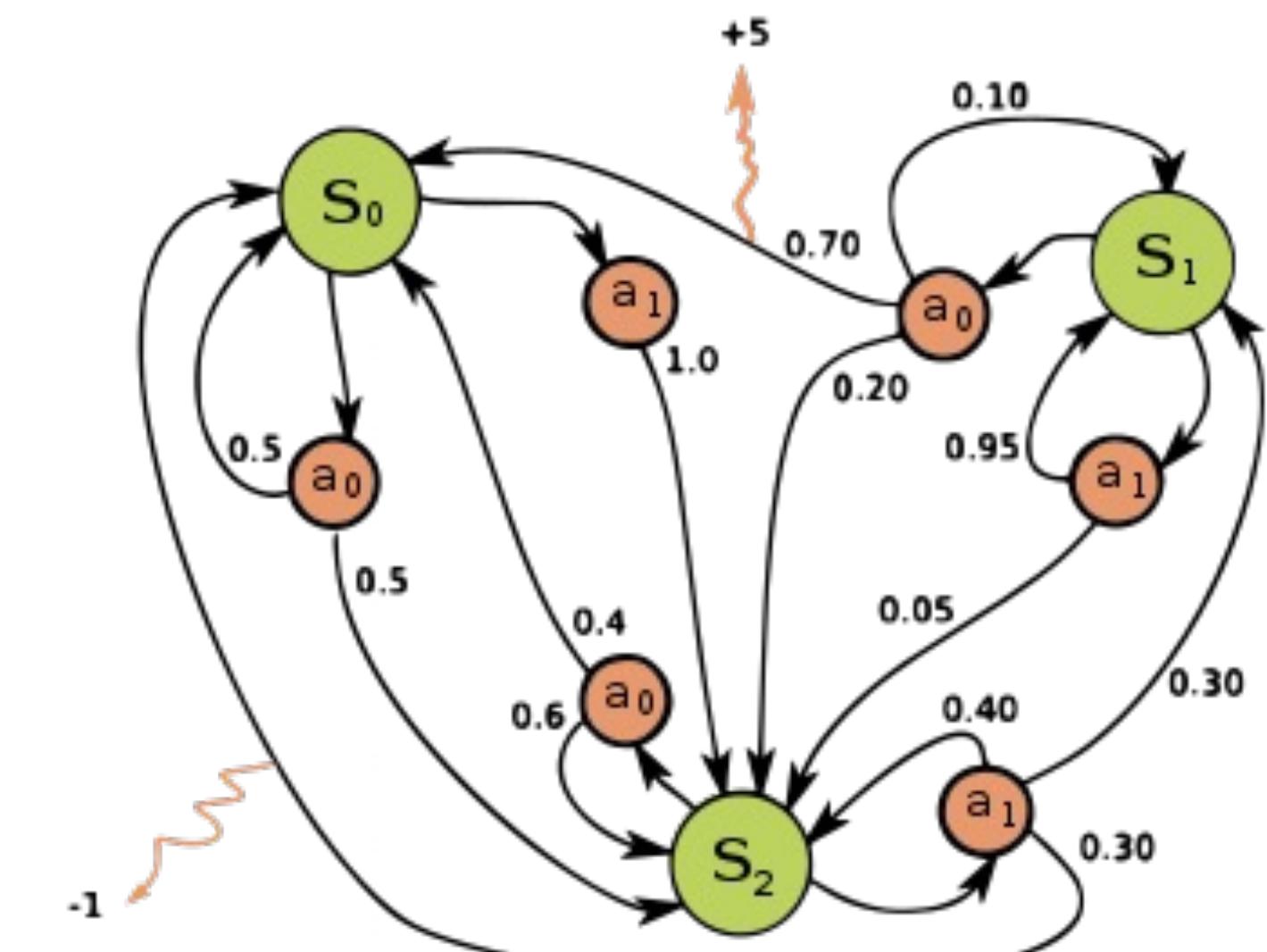
- **Background**
- **Use Demonstrations**
  - Offline – without interactions with environments
  - Online – with interactions with environments
- **Collect Demonstrations**
- **Case Study**
- **Future Directions**

# Outline

- **Background**
- **Use Demonstrations**
  - Offline – without interactions with environments
  - Online – with interactions with environments
- **Collect Demonstrations**
- **Case Study**
- **Future Directions**

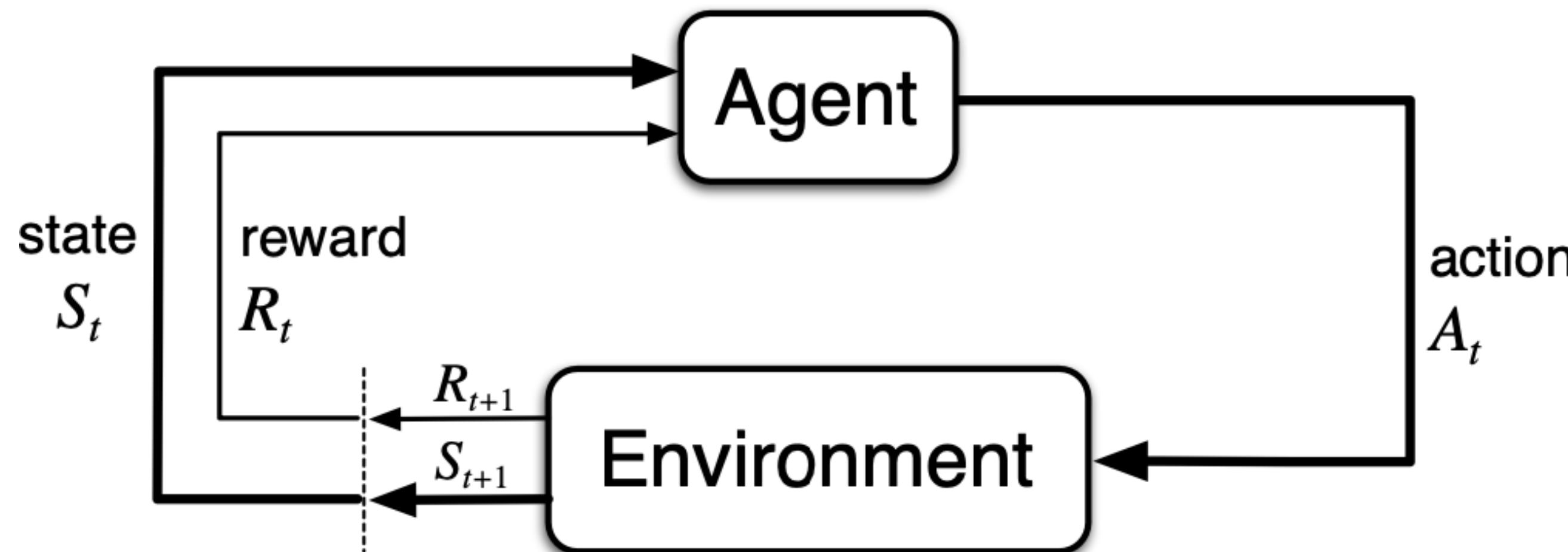
# Markov Decision Process (MDP)

- Markov decision process (MDP)
  - Probabilistic description of environment
  - **Markov process with rewards and decisions**
- A **Markov decision process** is a tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R})$ 
  - $\mathcal{S}$  is a set of states (discrete or continuous)
  - $\mathcal{A}$  is a set of actions (discrete or continuous)
  - $\mathcal{P}$  is a state transition probability function
    - $\mathcal{P}_{s,s'}^a = P(s'|s, a) = \Pr(S_{t+1} = s' | S_t = s, A_t = a)$
  - $\mathcal{R}$  is a reward function
    - $\mathcal{R}_s^a = R(s, a) = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$
  - Sometimes, an MDP also includes an initial state distribution  $\mu$



# Agent-Environment Interface

- **Agent:** learner and decision maker
- **Environment:** the thing agent interacts with, comprising everything outside the agent.
- **Action:** how agent interacts with the environment.



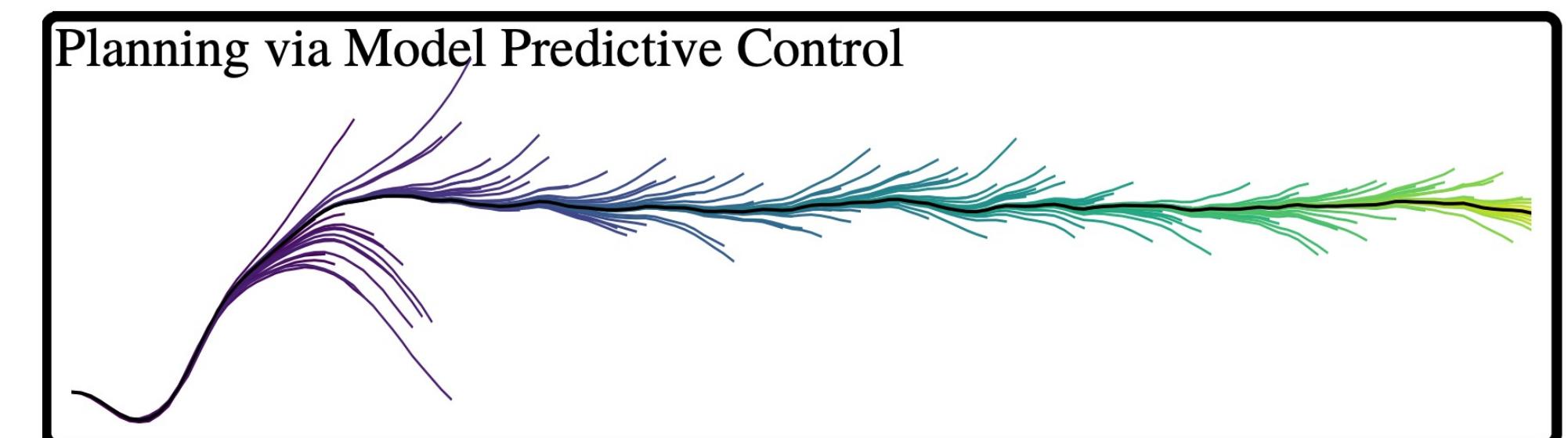
# Reinforcement Learning (RL)

- Goal of RL
  - To find an **optimal policy** in an MDP by **interacting** with it.
  - i.e., it offers an **online way** to solve an MDP
- Policy
  - A function mapping the current state to the next action choice(s)
  - Deterministic:  $\pi : S \rightarrow A$
  - Stochastic:  $\pi : S \rightarrow \text{Prob}(A)$
- Learning Objective
  - Maximize the expected return (accumulative rewards)  $E_{s,\pi(s)} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) | s_0 \right]$

# Planning (with Models)

- Goal of Planning
  - Solve the MDP with the environment model (including transition and reward)
  - The environment model can be either given or learned from data
- Model Predictive Control (MPC)
  - Continuous action space
  - At each time step, MPC samples multiple sequences and executes the first action of the best sequence to the environment.

$$\max_{a_{t:t+\tau}} \mathbb{E}_{s_{t'+1} \sim p(s_{t'+1}|s_{t'}, a_{t'})} \left[ \sum_{t'=t}^{t+\tau} r(s_{t'}, a_{t'}) \right]$$



# Planning (with Models)

- Monte Carlo Tree Search (MCTS)
  - Discrete action space
  - At each time step, MCTS incrementally extends a search tree from the current environment state.
  - Each node in the tree corresponds to a state, which will be evaluated by the return obtained after rollouts in the model with a random policy.

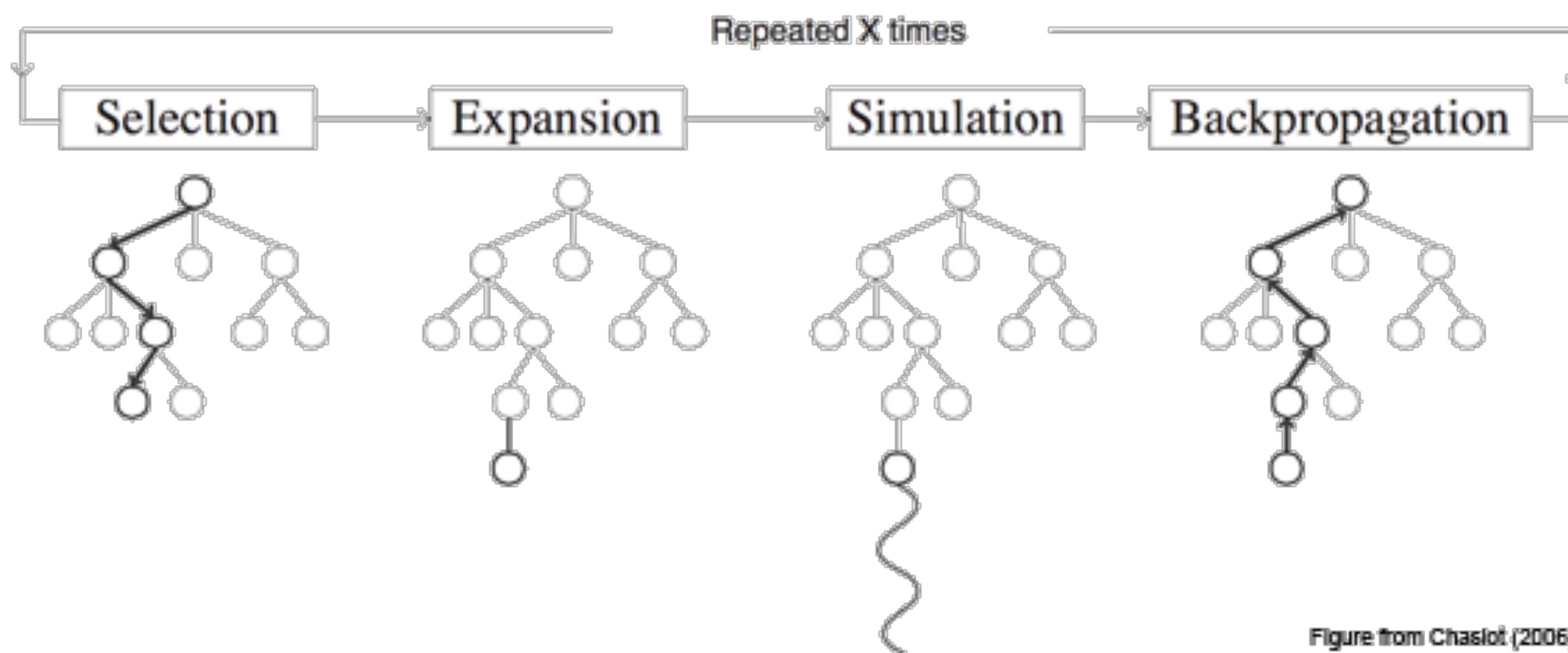


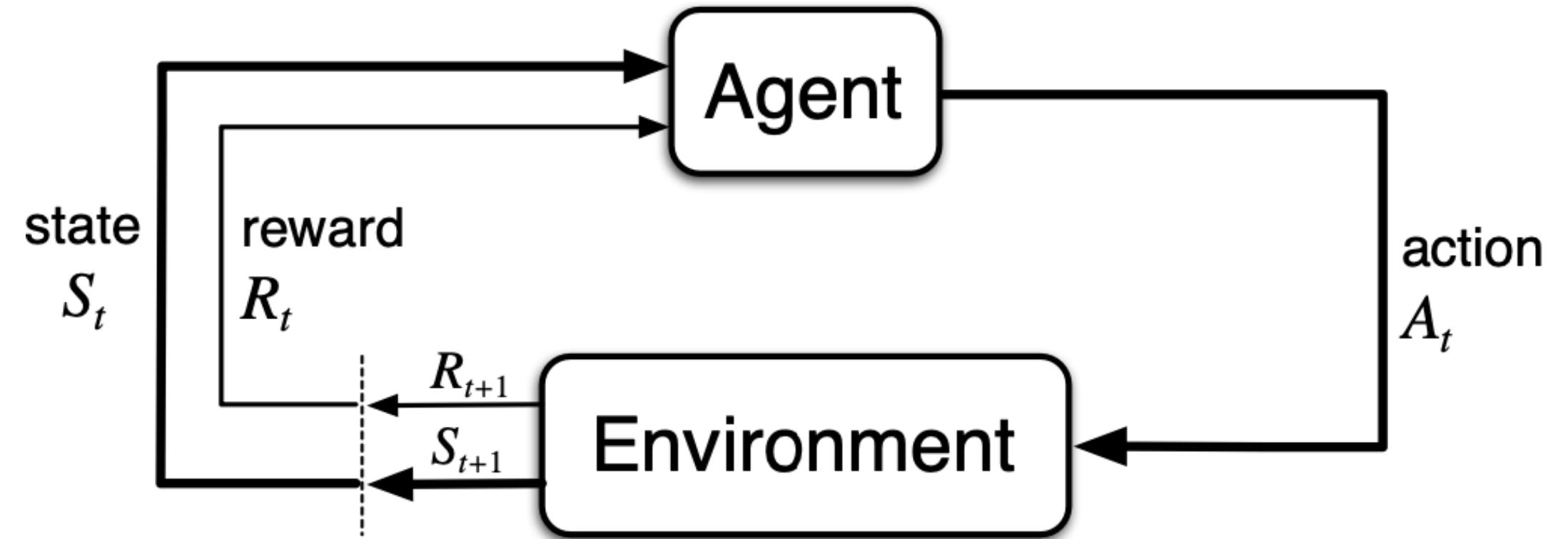
Figure from Chaslot (2006)

# Outline

- **Background**
- **Use Demonstrations**
  - Offline – without interactions with environments
  - Online – with interactions with environments
- **Collect Demonstrations**
- **Case Study**

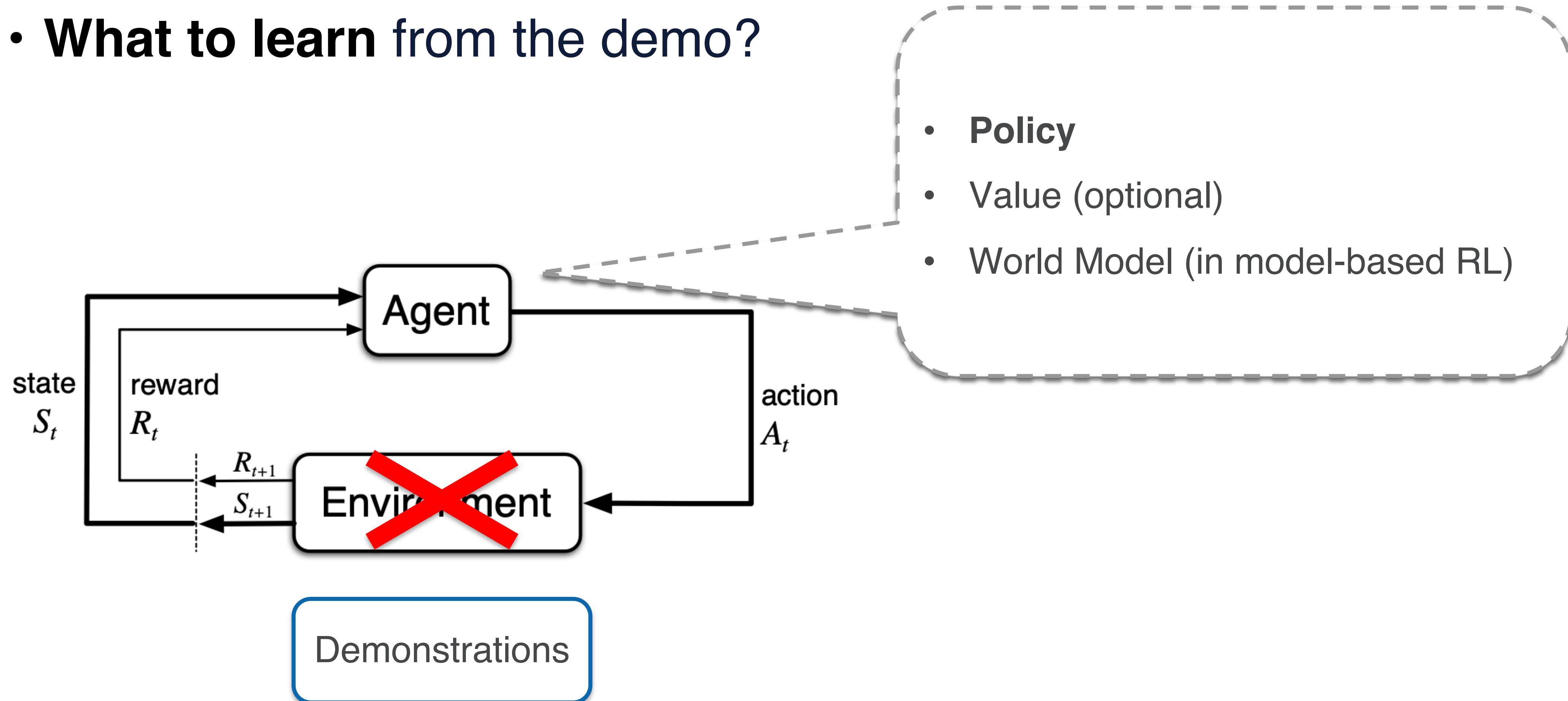
# Demo Use – Offline & Online

- Offline Stage
  - Agent **cannot** access the environment
  - Only learn from the demonstrations
- Online Stage
  - Agent **can** interact with the environment
  - Learn from the environment while still leveraging demonstrations
- Note
  - Both stages are optional
  - They can also be combined together



# Use Demo Offline

- What to learn from the demo?

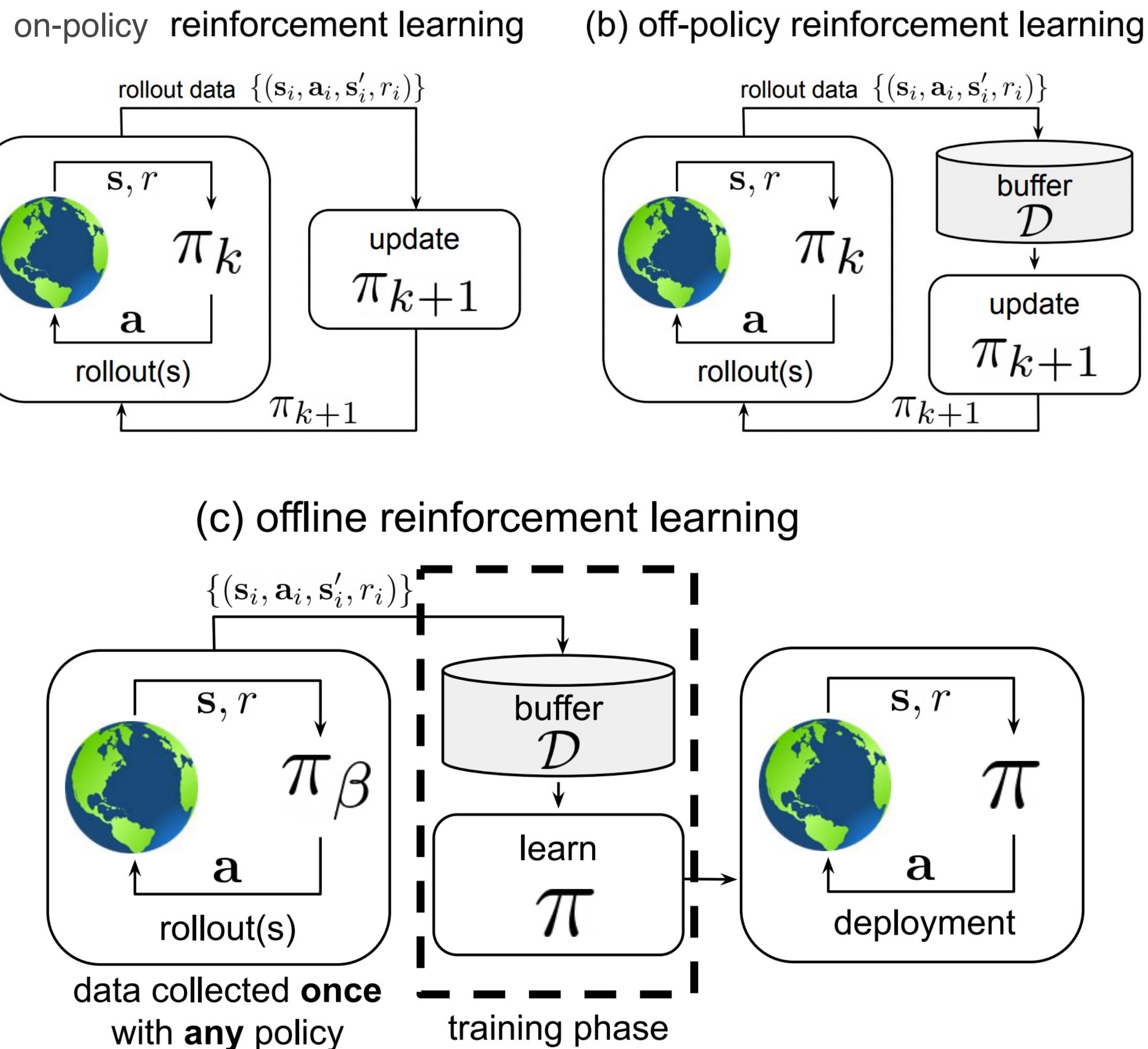


# Learn Policy

- Imitation Learning
  - Imitate the behaviors in demonstrations
  - Two types of methods: **behavior cloning** and **inverse RL**
- Behavior Cloning (BC)
  - Supervised learning, i.e., clone the expert's actions at each state in demo
  - $$\underset{\theta}{\text{maximize}} \sum_{(s,a) \in \rho_D} \ln \pi_{\theta}(a|s)$$
- Inverse RL
  - Recover a reward function that can induce the behaviors in demonstrations
  - We will talk more about this later
- Limitations of Imitation Learning
  - Requires **expert** demonstrations

# Learn Policy

- Offline RL



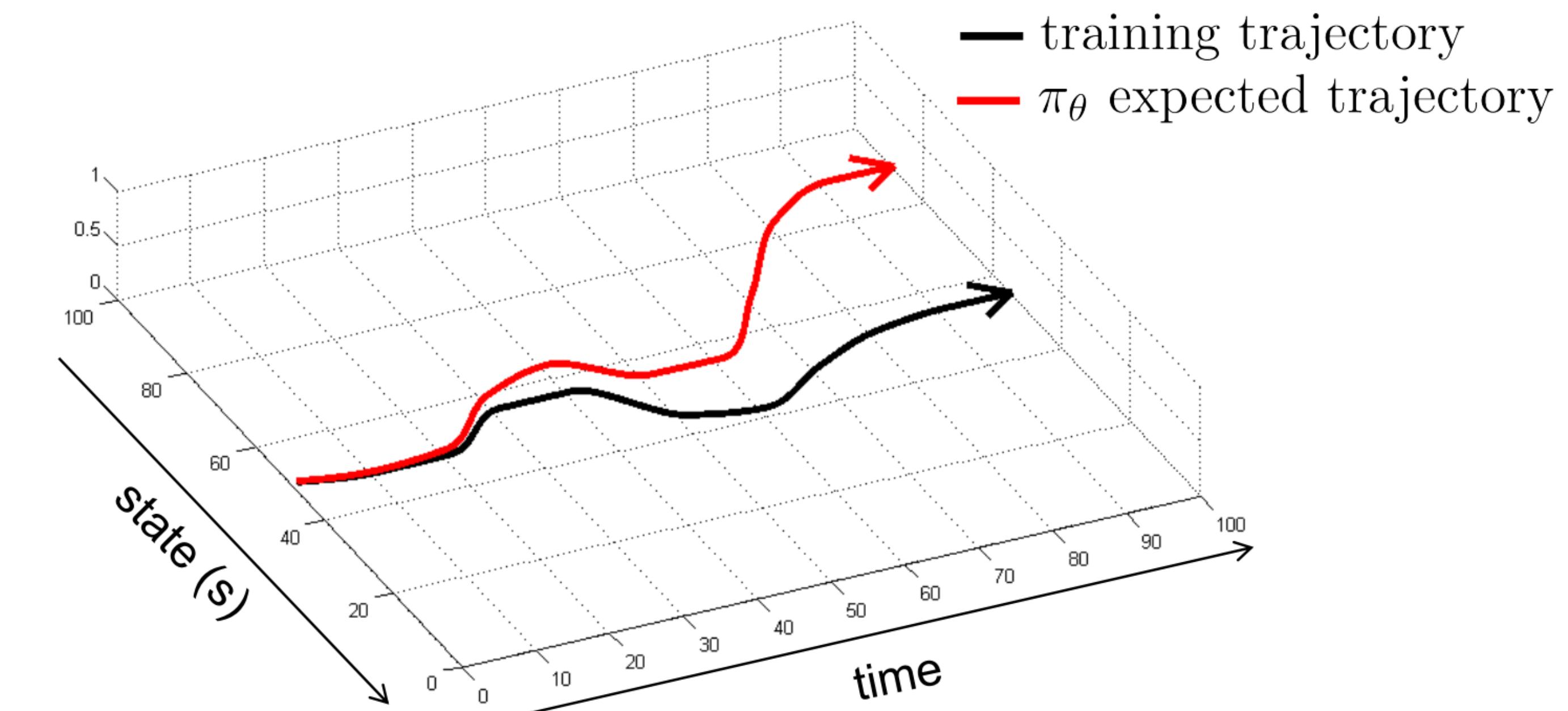
# Learn Policy

- Is offline learned policy good enough?
  - No, due to the **distribution shift**
  - When we use a learned policy to act, it changes what we see

**Simple example:**

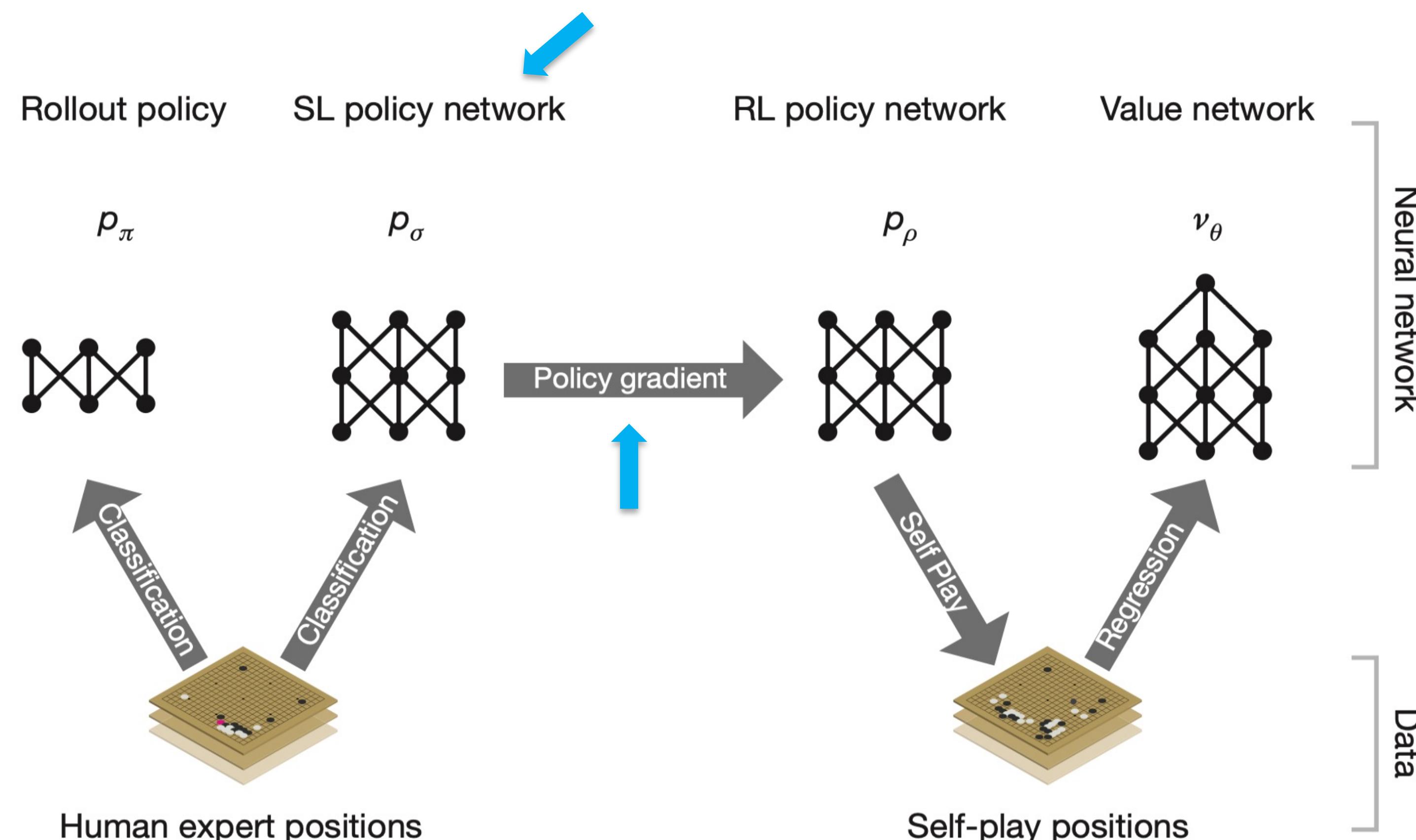
behavioral cloning  
train  $\pi_\theta$  to *copy*  
assume data is *optimal*

error scales as  $O(T^2)$



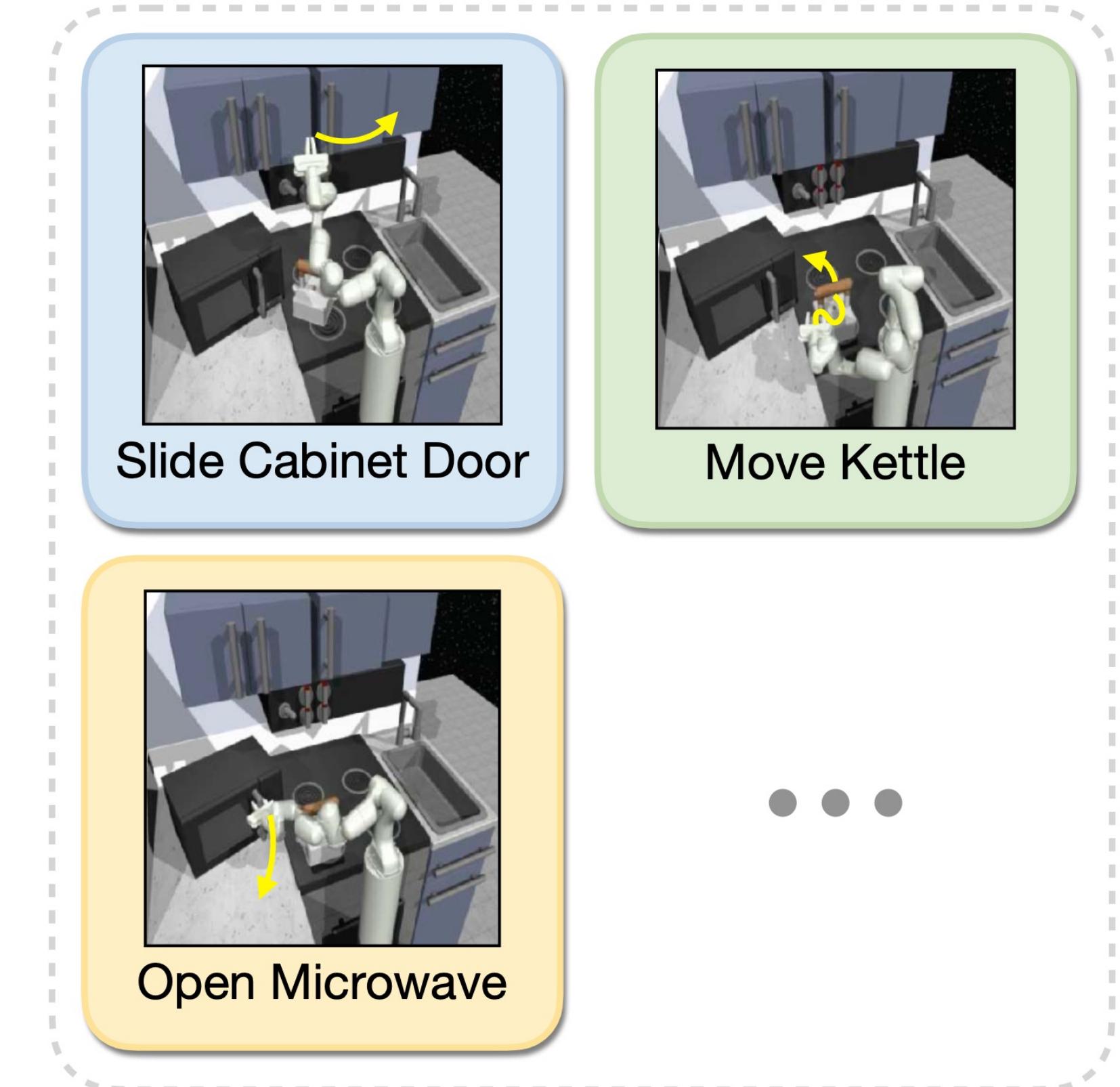
# Learn Policy

- Common Practice: Fine-tune it by online RL
  - Example: AlphaGo



# Skill Discovery

- Policy can decomposed into **skills**
  - A cooking policy can be decomposed into “slide door”, “move kettle”, etc.
  - **Skills** to refer to some short-horizon policies shared between many tasks.
- Demo may come from many different tasks
  - Learn a set of shared skills from demo
  - Utilize the learned skills in downstream tasks



# Skill Discovery

- Minimize the description lengths

The description length of skill encodings

$$\mathcal{L}_{\text{Love}} = n_s \cdot c_{\text{avg}}$$

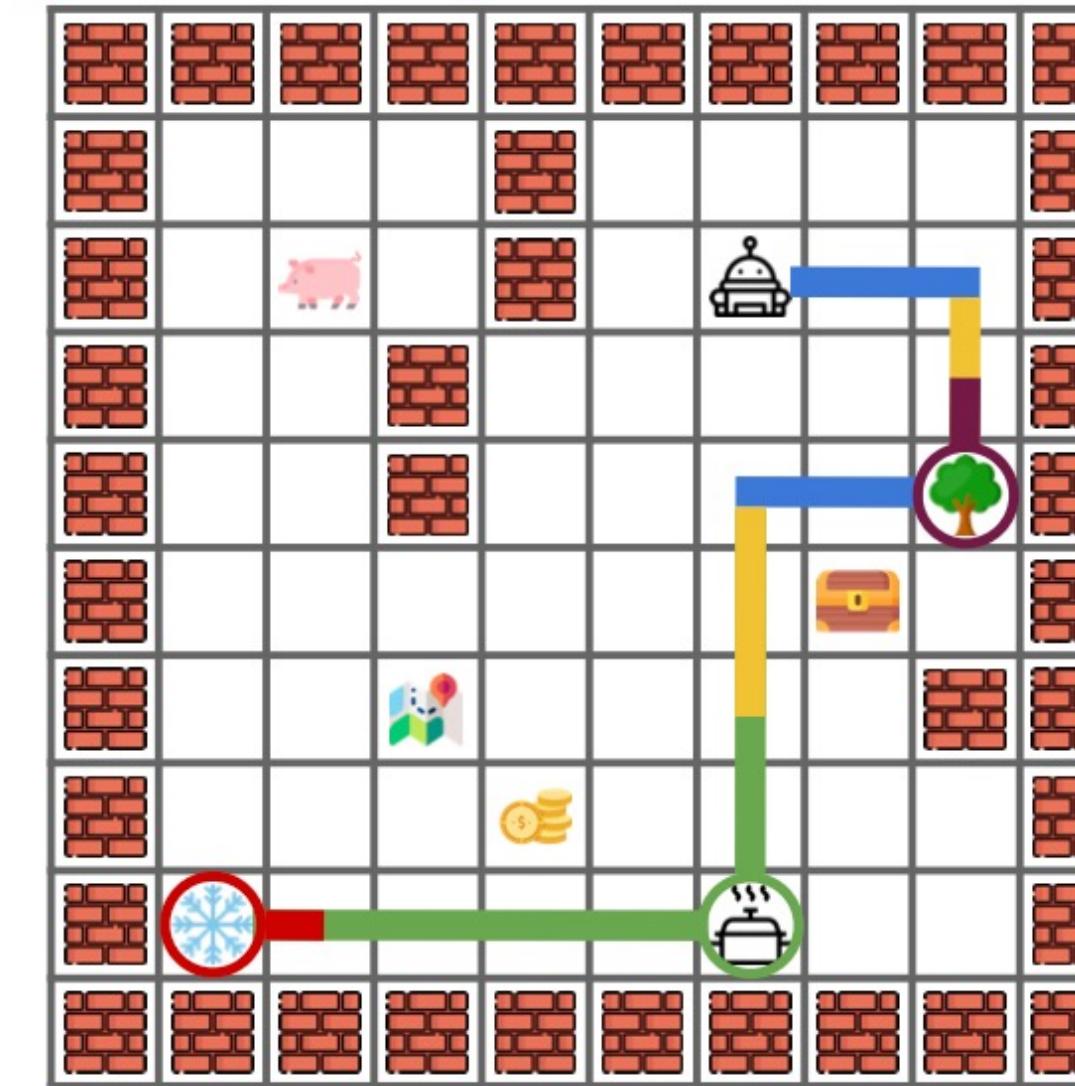
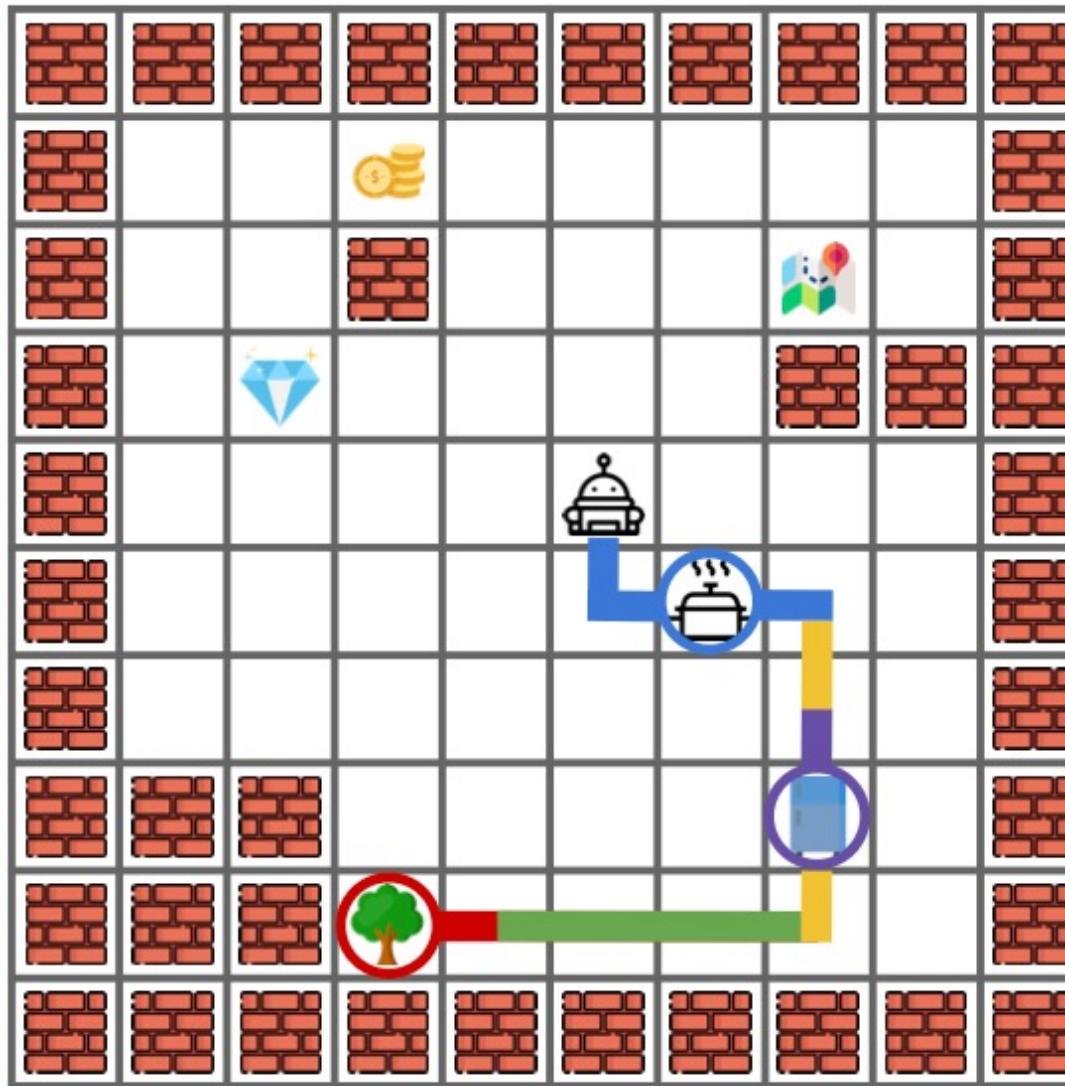
$$n_s = \mathbb{E} \left[ \sum_{t=1}^T m_t \right]$$

Avg. # of skills per trajectory

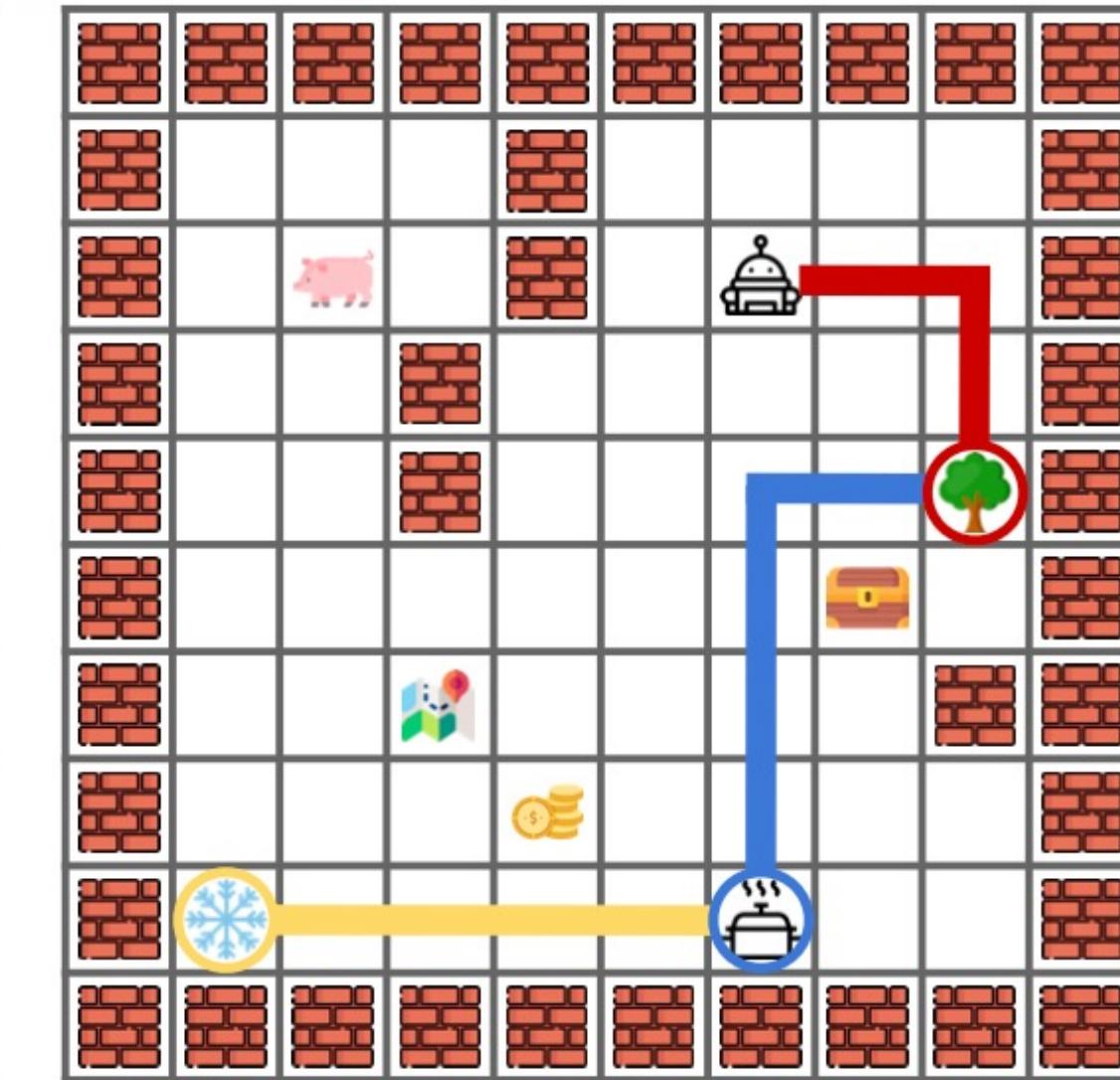
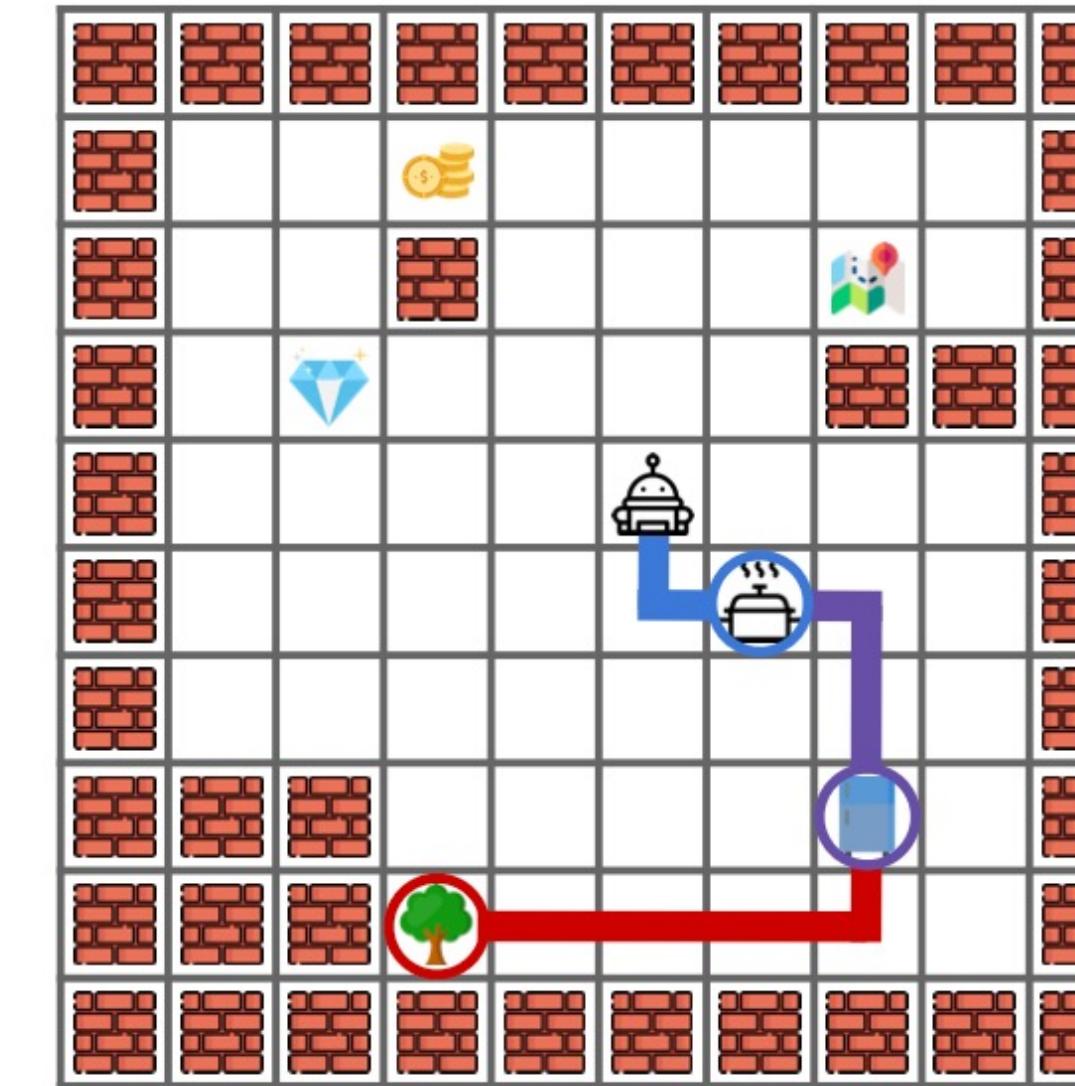
$$c_{\text{avg}} = \mathbb{E}_{p(z)} [-\log p(z)]$$

Marginal entropy

Skills that **arbitrarily segment** the trajectories

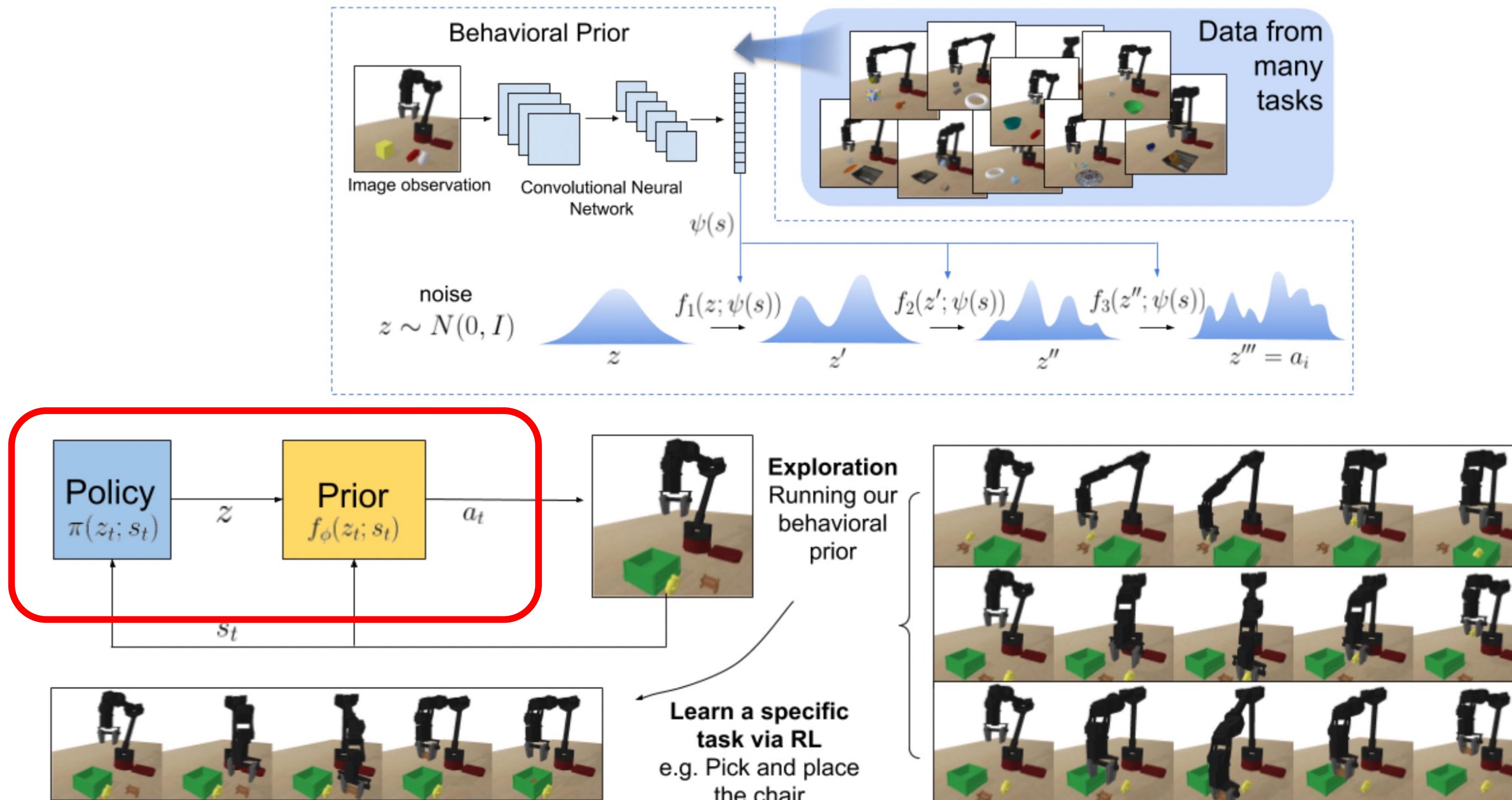


Skills that **navigate to objects**



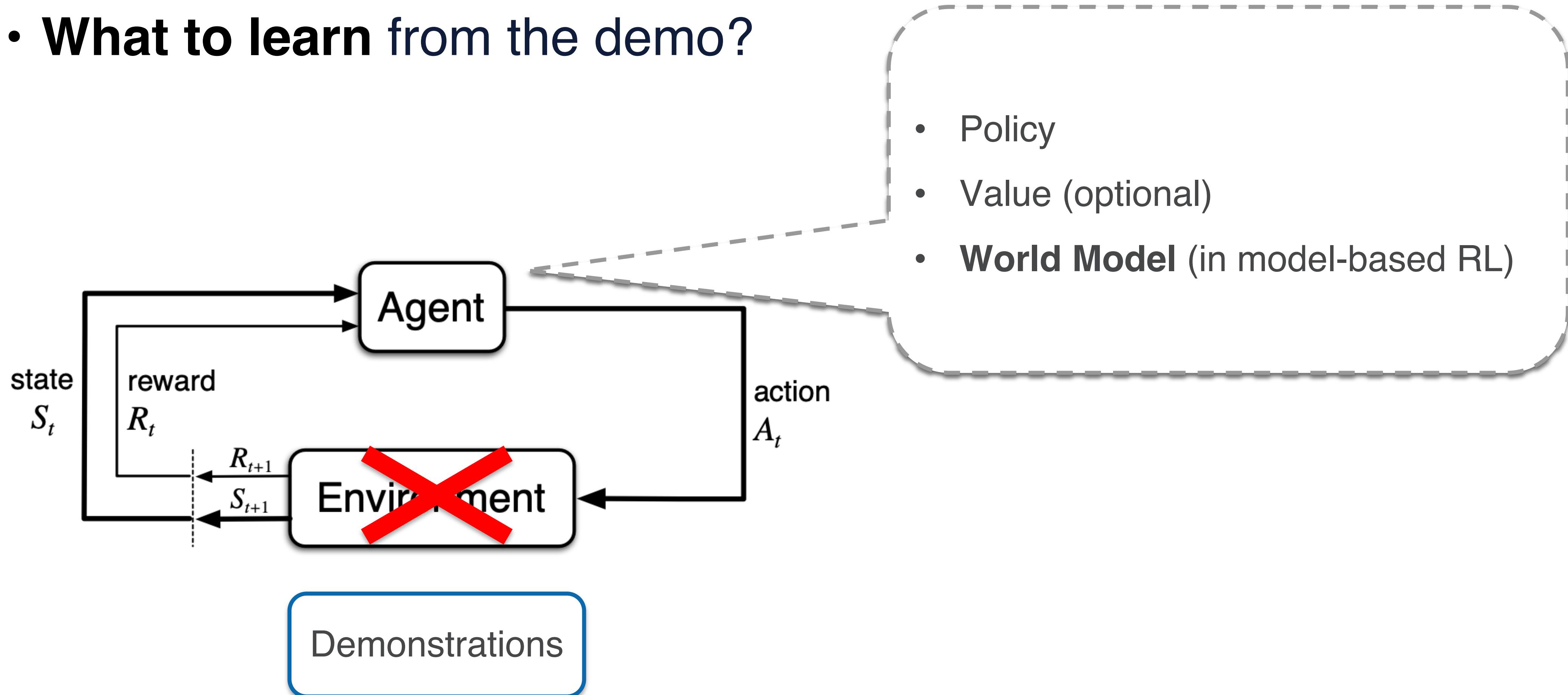
# Skill Discovery

- Utilize the learned skills by calling them in new tasks



# Use Demo Offline

- What to learn from the demo?



# Learn World Model

- Learn deterministic world model from demo

$$\min_{\theta} \mathbb{E}_{(s,a) \sim D, s' \sim T^*(\cdot|s,a)} \left[ \|s' - T_{\theta}(s, a)\|_2^2 \right]$$

↑                   ↑                   ↑  
Demo Dataset      Real Transition      Predicted Transition

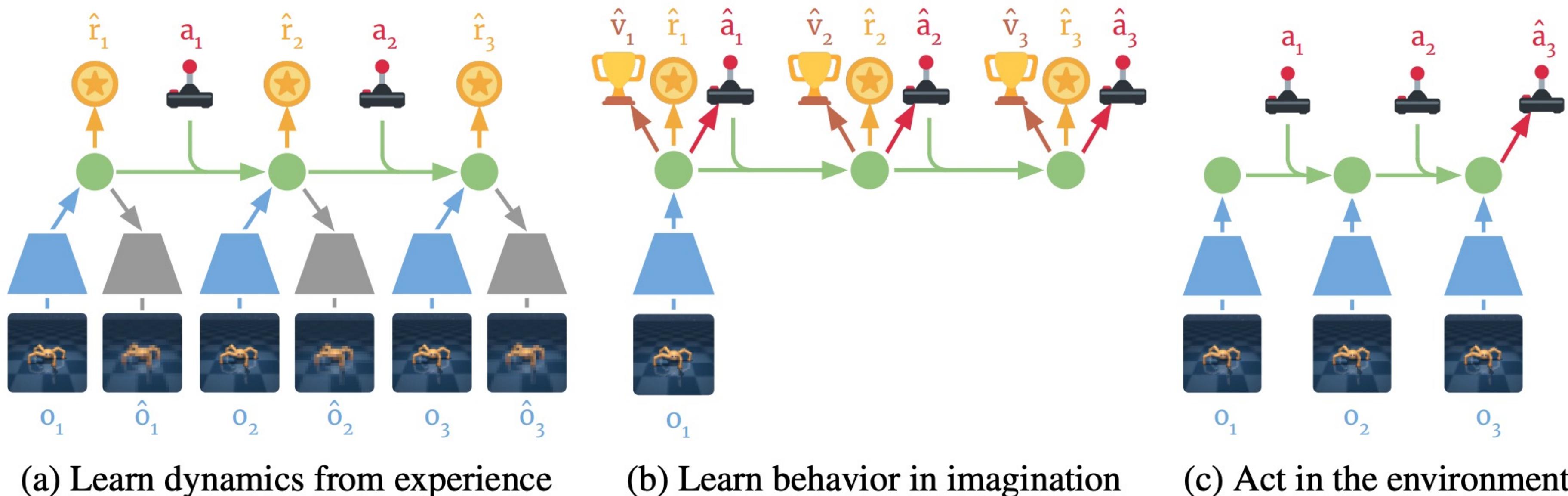
- Learn stochastic world model from demo

$$\min_{\theta} \mathbb{E}_{(s,a) \sim D} [D_{\text{KL}} (T^*(\cdot|s, a), T_{\theta}(\cdot|s, a))] := \mathbb{E}_{(s,a) \sim D, s' \sim T^*(\cdot|s,a)} \left[ \log \left( \frac{T^*(s' | s, a)}{T_{\theta}(s' | s, a)} \right) \right]$$

- Demo dataset should be **diverse**, but not necessarily be optimal

# Learn World Model – RL in World Model

- Run model-free RL in the learned world model
  - i.e., RL in imagination
  - Collected data can be used to further improve the learned world model



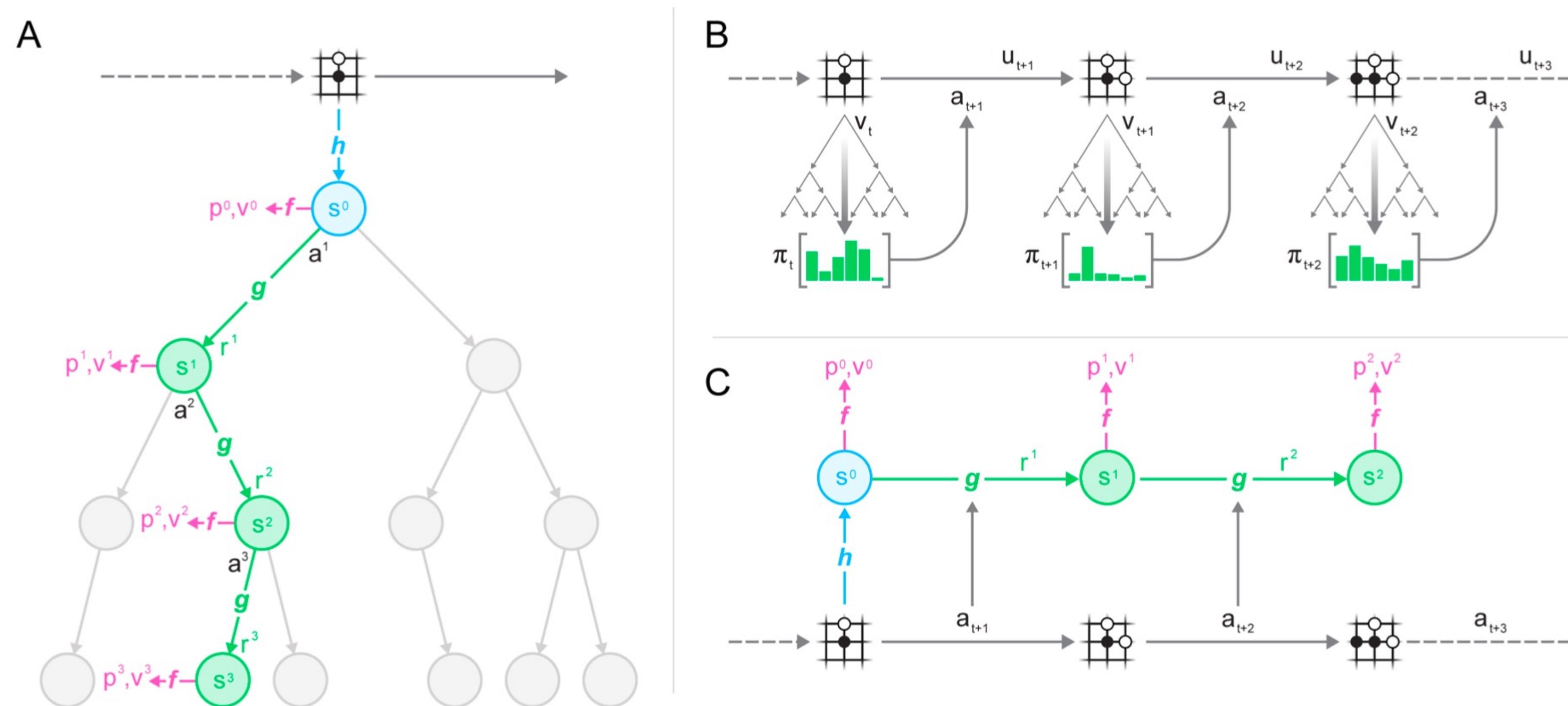
(a) Learn dynamics from experience

(b) Learn behavior in imagination

(c) Act in the environment

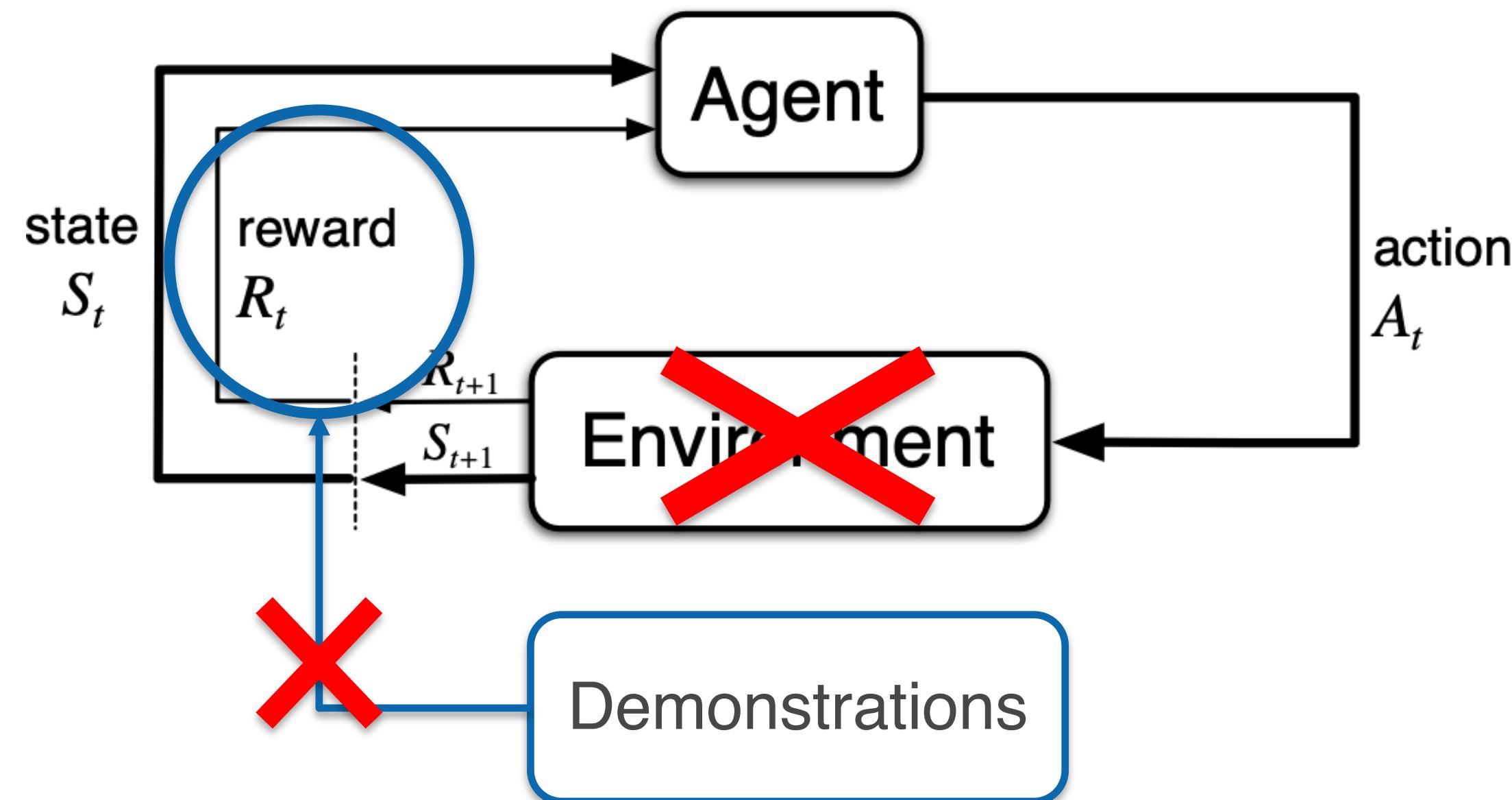
# Learn World Model – Plan in World Model

- Plan in the learned world model
  - E.g., use MCTS to search in the learned model (MuZero)



# Use Demo Offline

- What to learn from the demo?

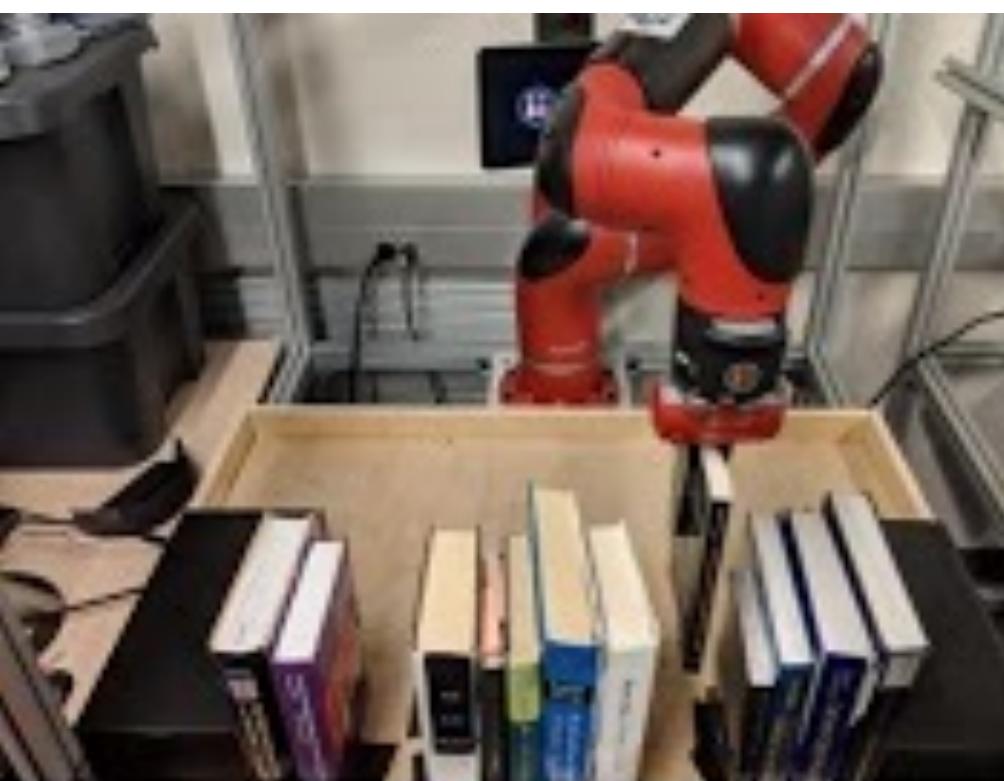


Reward is not present in demo?

We can still infer the reward offline if we have some **trajectory-level labels**,  
e.g., **success** and **failure**

# Infer Reward / Goal

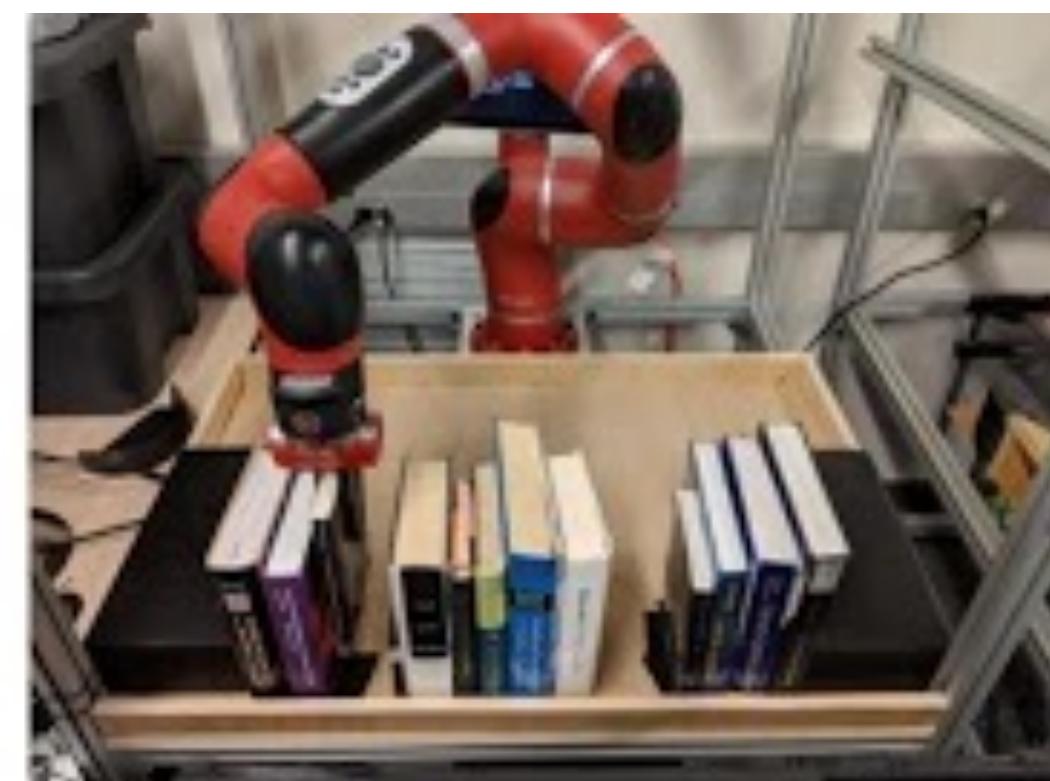
- Learn a **goal classifier** from demo, and regard it as a reward
  - Positive samples: the last observation in each success trajectory
  - Negative samples: observations from failure trajectories



Failure



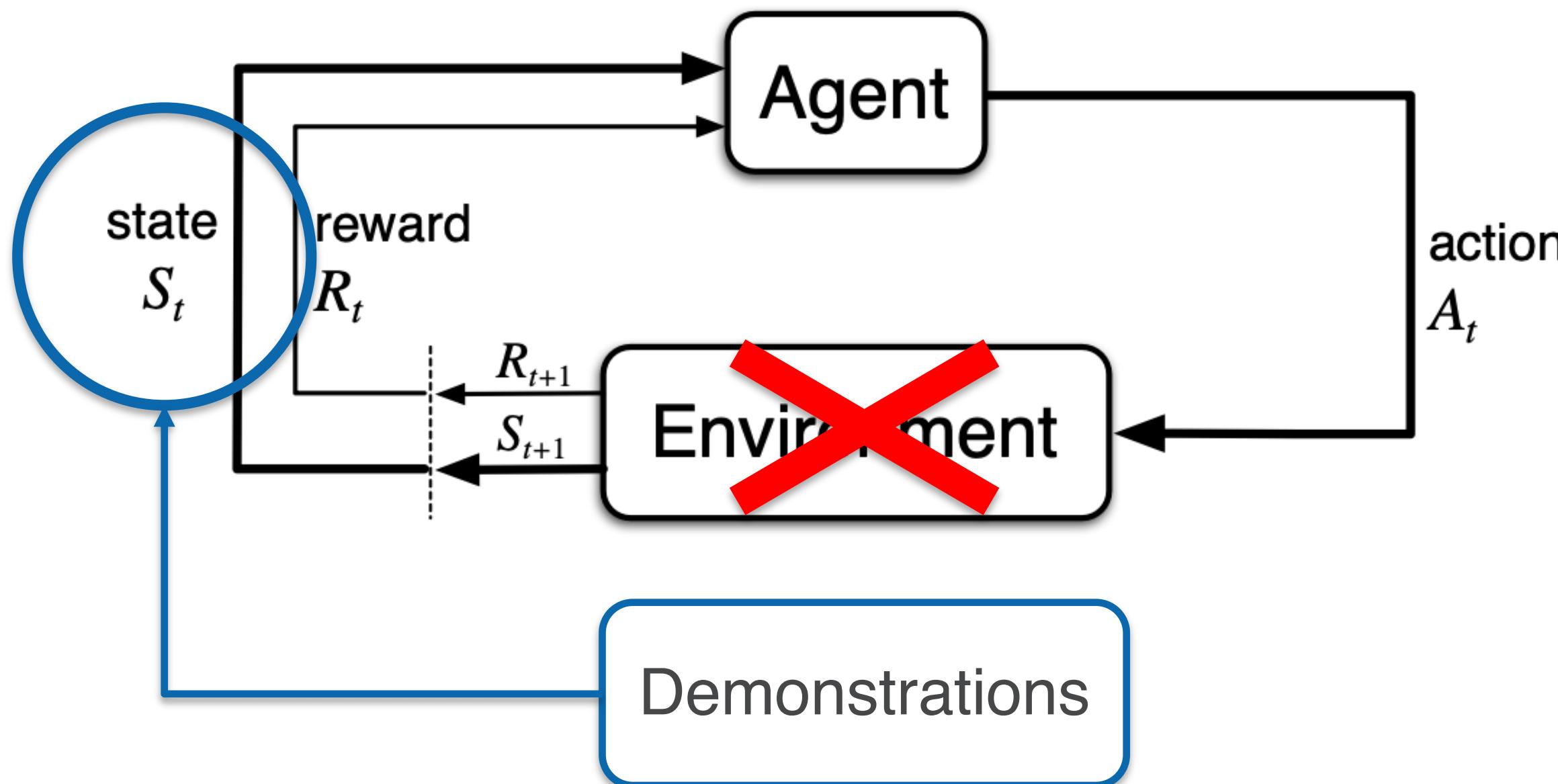
Success



Success

# Use Demo Offline

- What to learn from the demo?

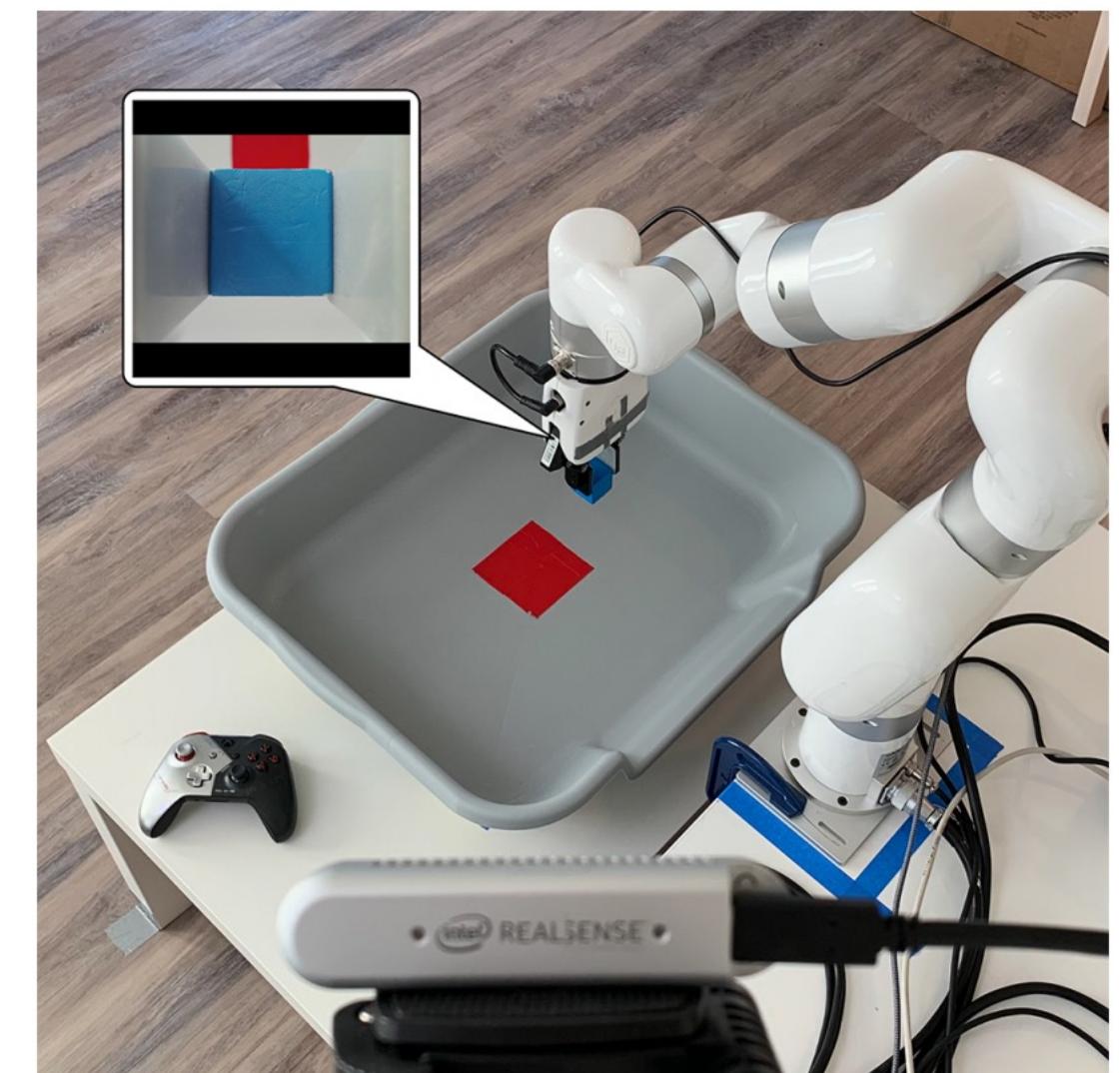
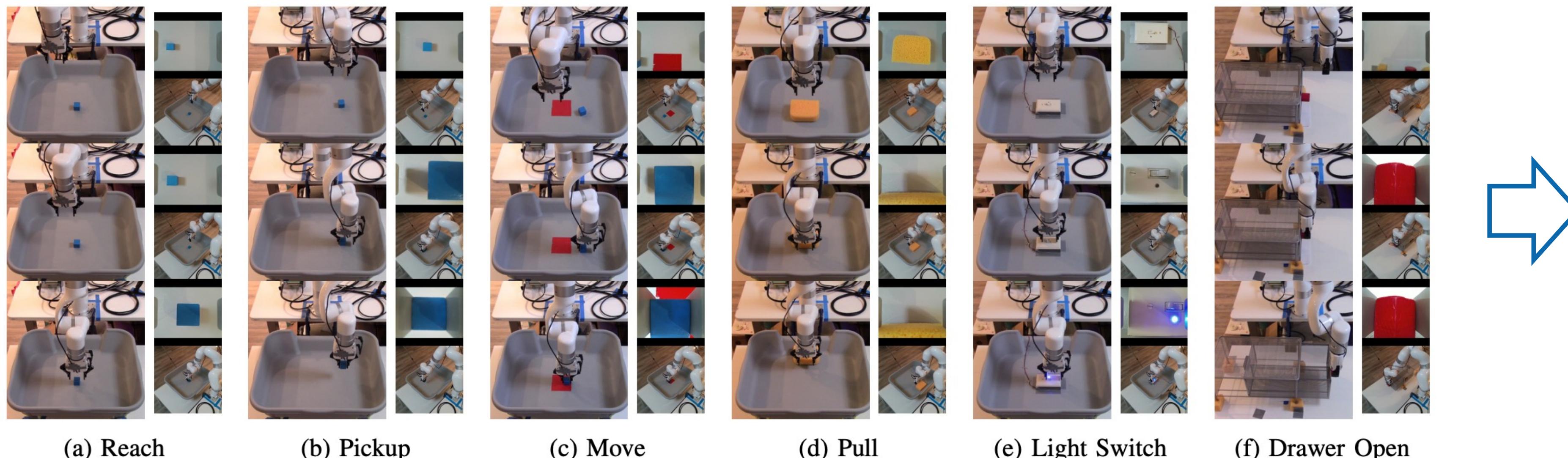


In visual RL, the observations / states usually very high-dimensional

## Representation Learning

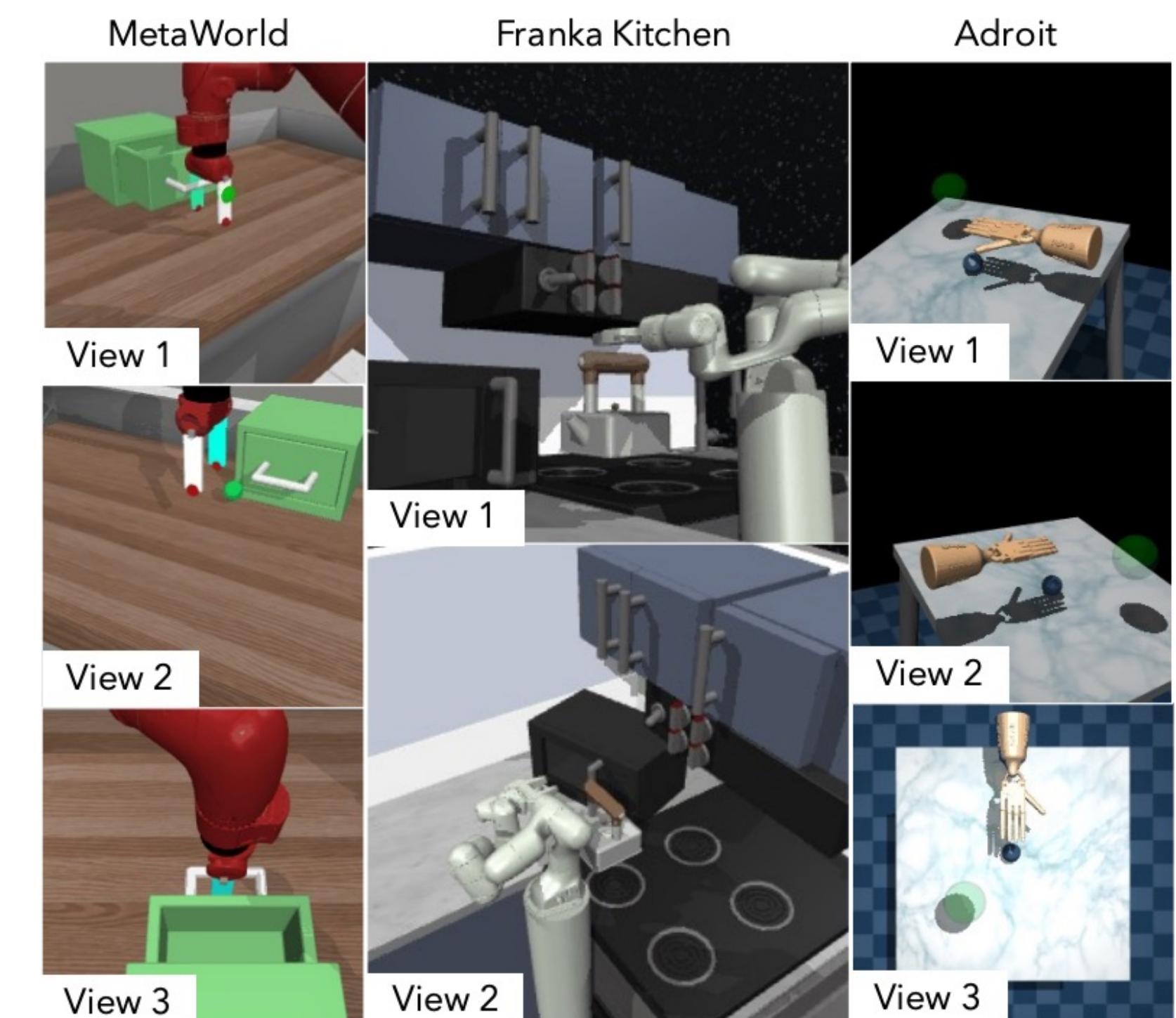
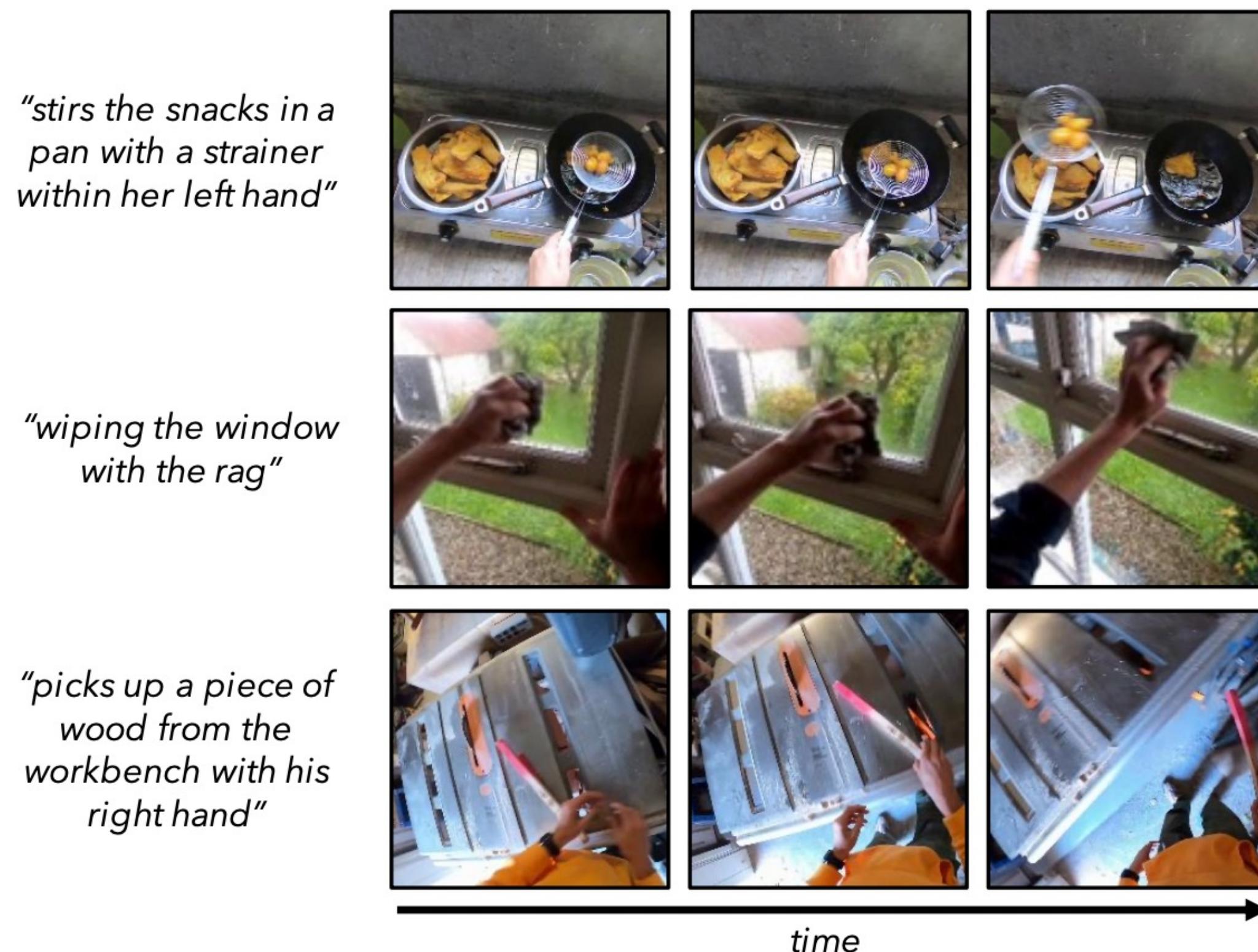
# Learn Representation

- Learn representation in in-domain dataset
  - E.g., pre-train the visual encoder by contrastive learning



# Learn Representation

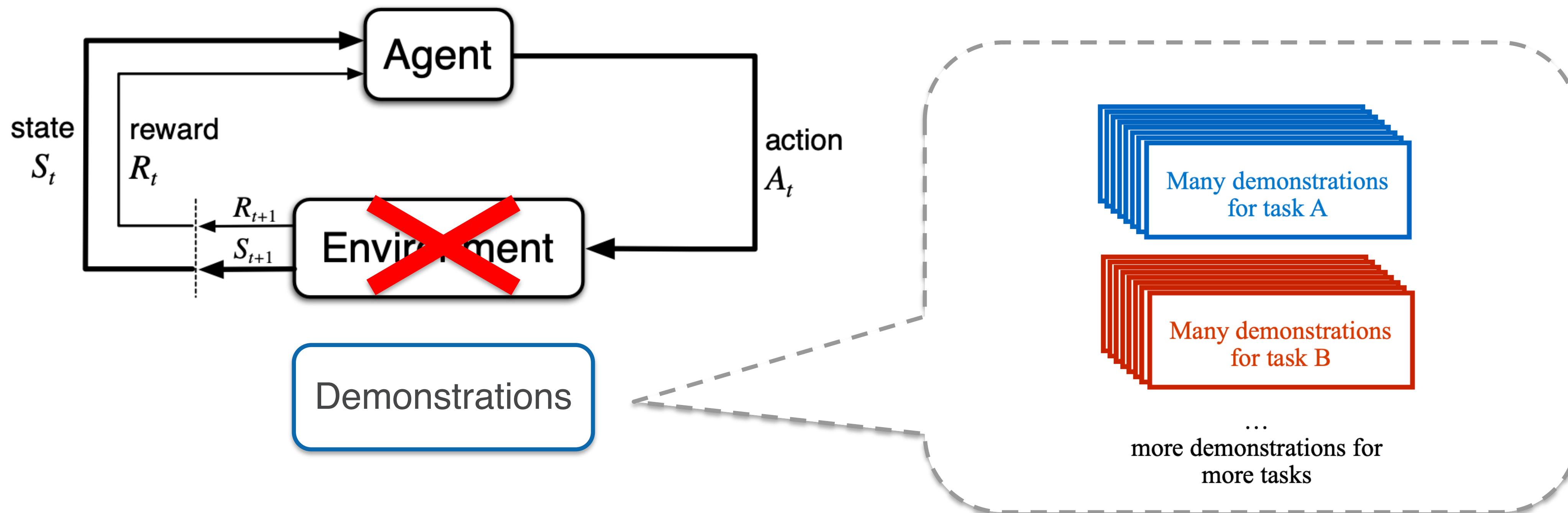
- Learn representation in large external datasets



# Use Demo Offline

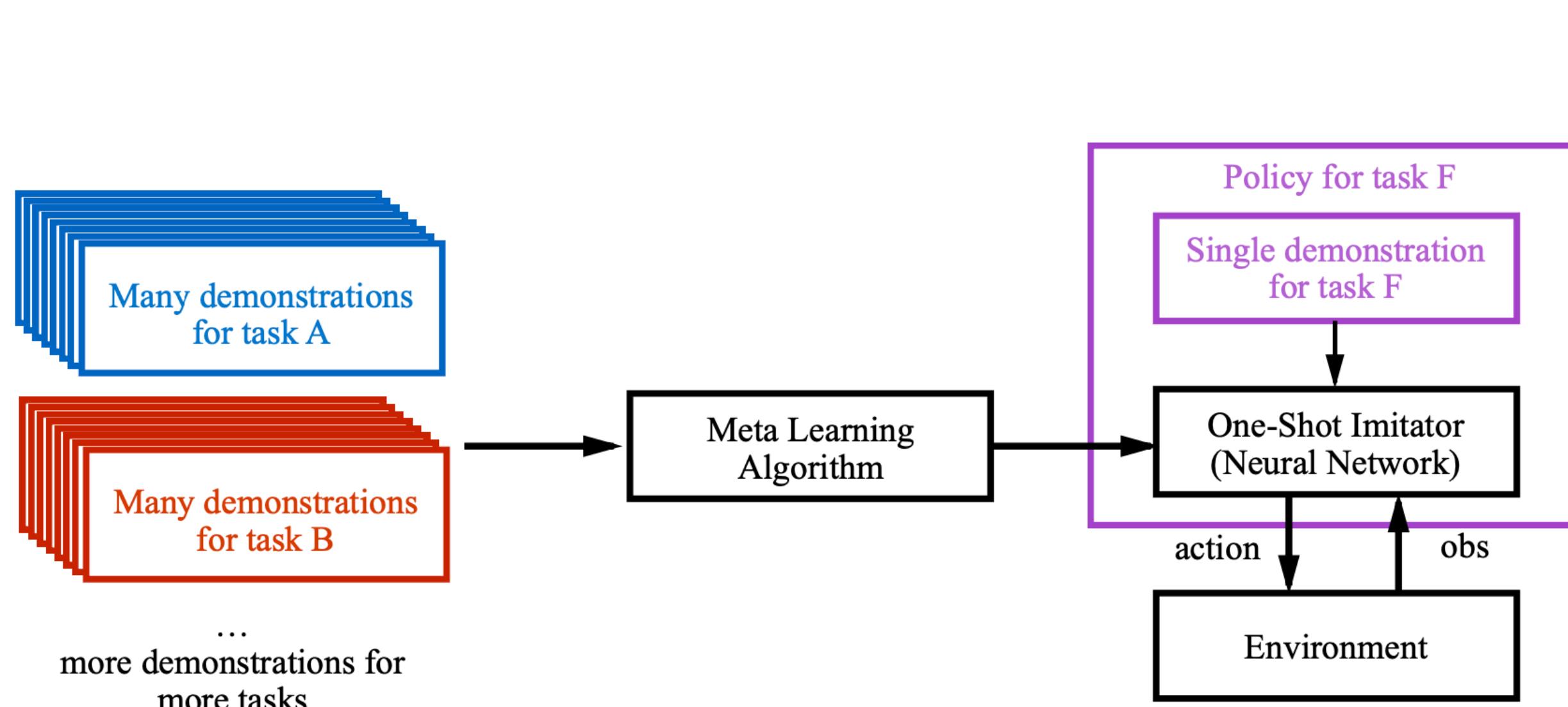
- What to learn from the demo?

What if we have demo from  
many different tasks?

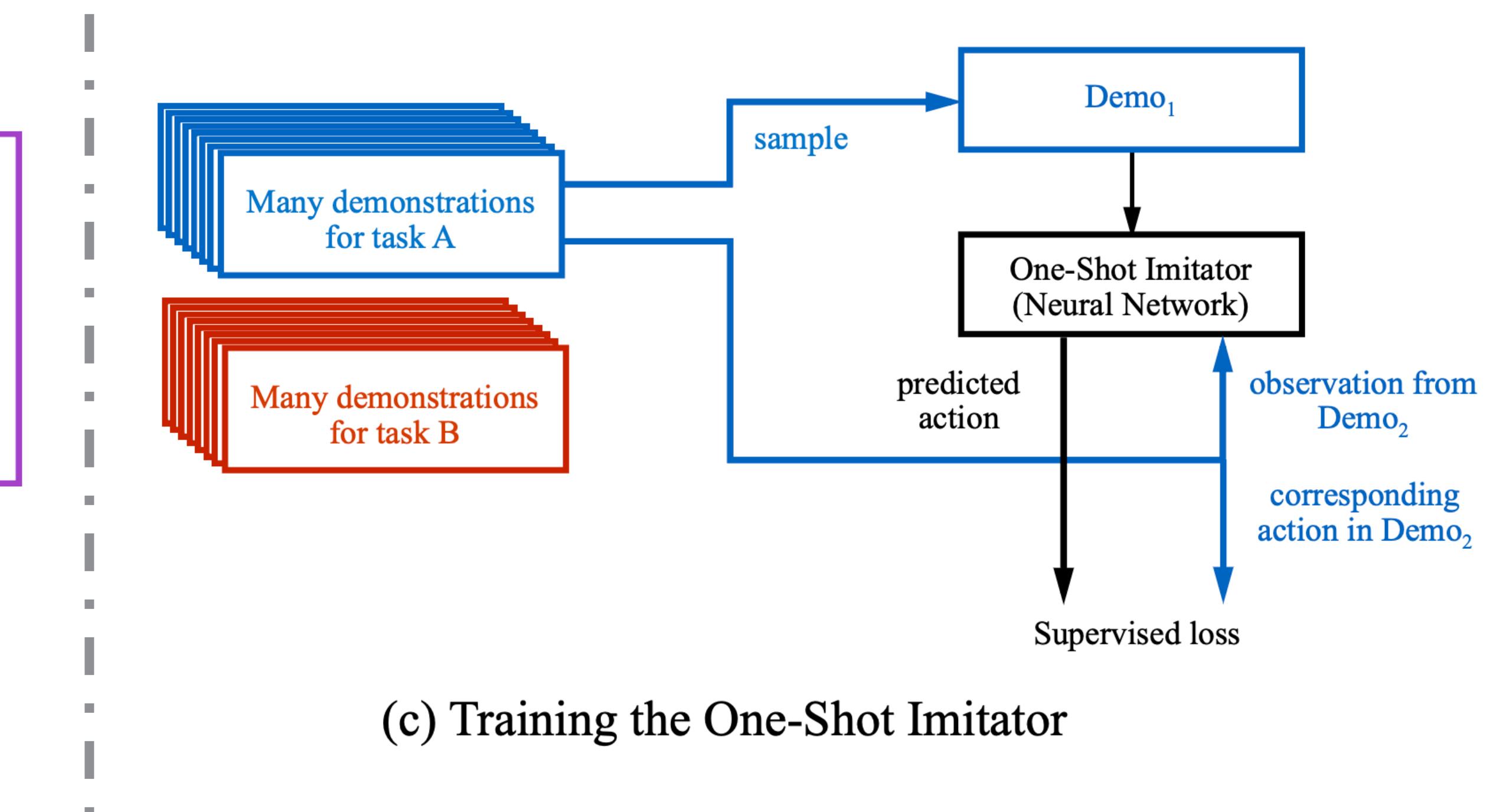


# Learn Trajectory Imitator

- One-Shot Imitation Learning
  - Given a trajectory at the test time
  - Trained to imitate the trajectory, instead of completing the tasks



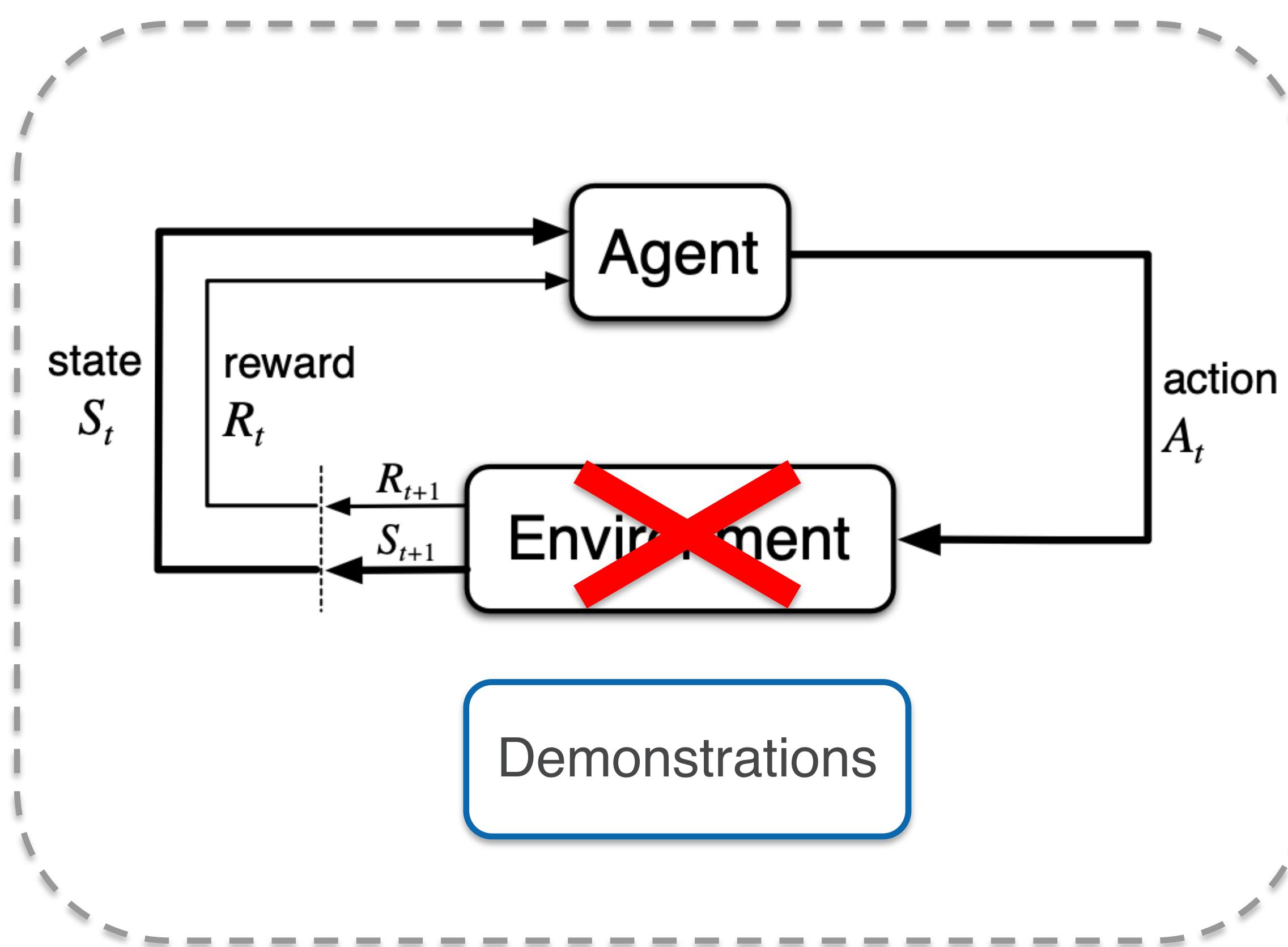
(b) One-Shot Imitation Learning



(c) Training the One-Shot Imitator

# Use Demo Offline

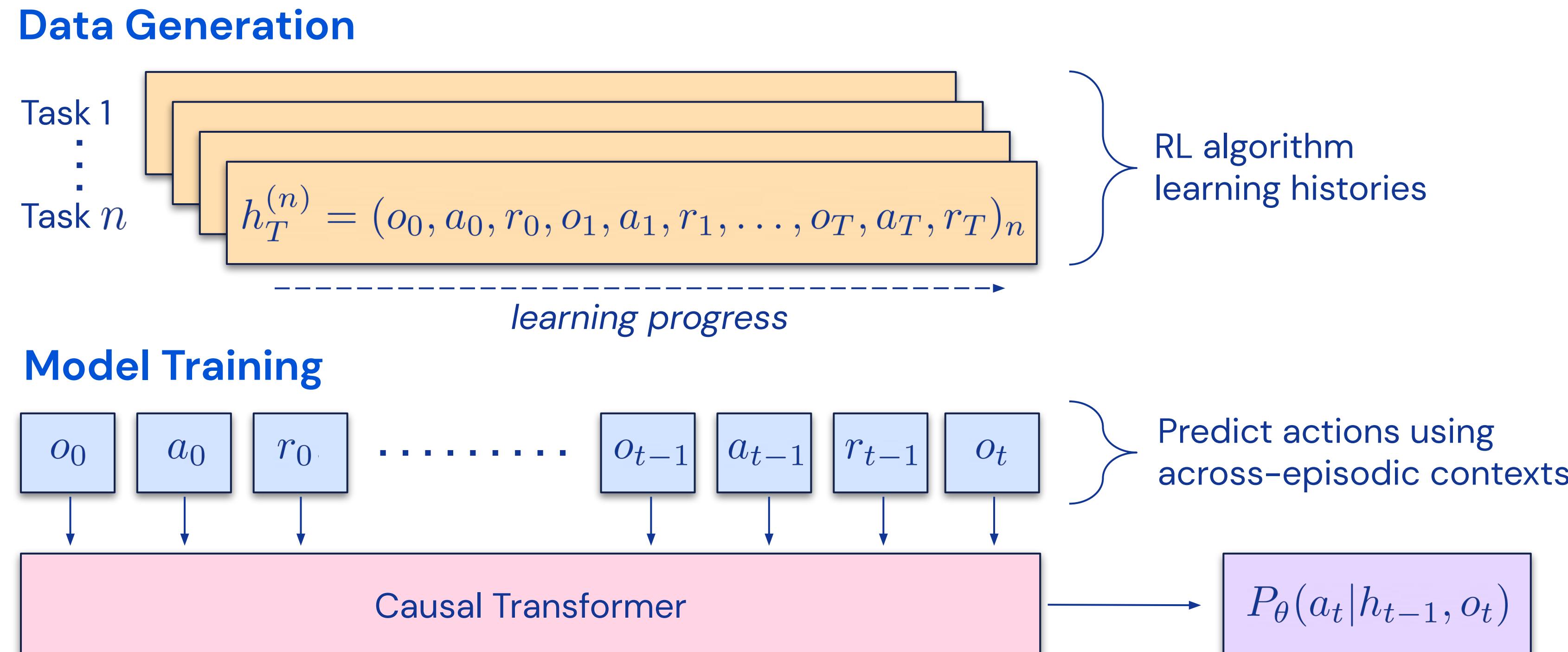
- What to learn from the demo?



Learn the RL **algorithm** itself?

# Learn Algorithm

- Algorithm Distillation
  - Assume the demo dataset contains **the whole learning history** of an agent
  - Train a transformer to predict actions given the preceding learning history



# Use Demo Offline

- **What to learn** from the demo?
  - Policy
  - Skill
  - World Model
  - Reward / Goal
  - Representation
  - Trajectory Imitator
  - Algorithm
  - ...
  - Maybe there will be more creative ways to use demo offline in the future?

# Outline

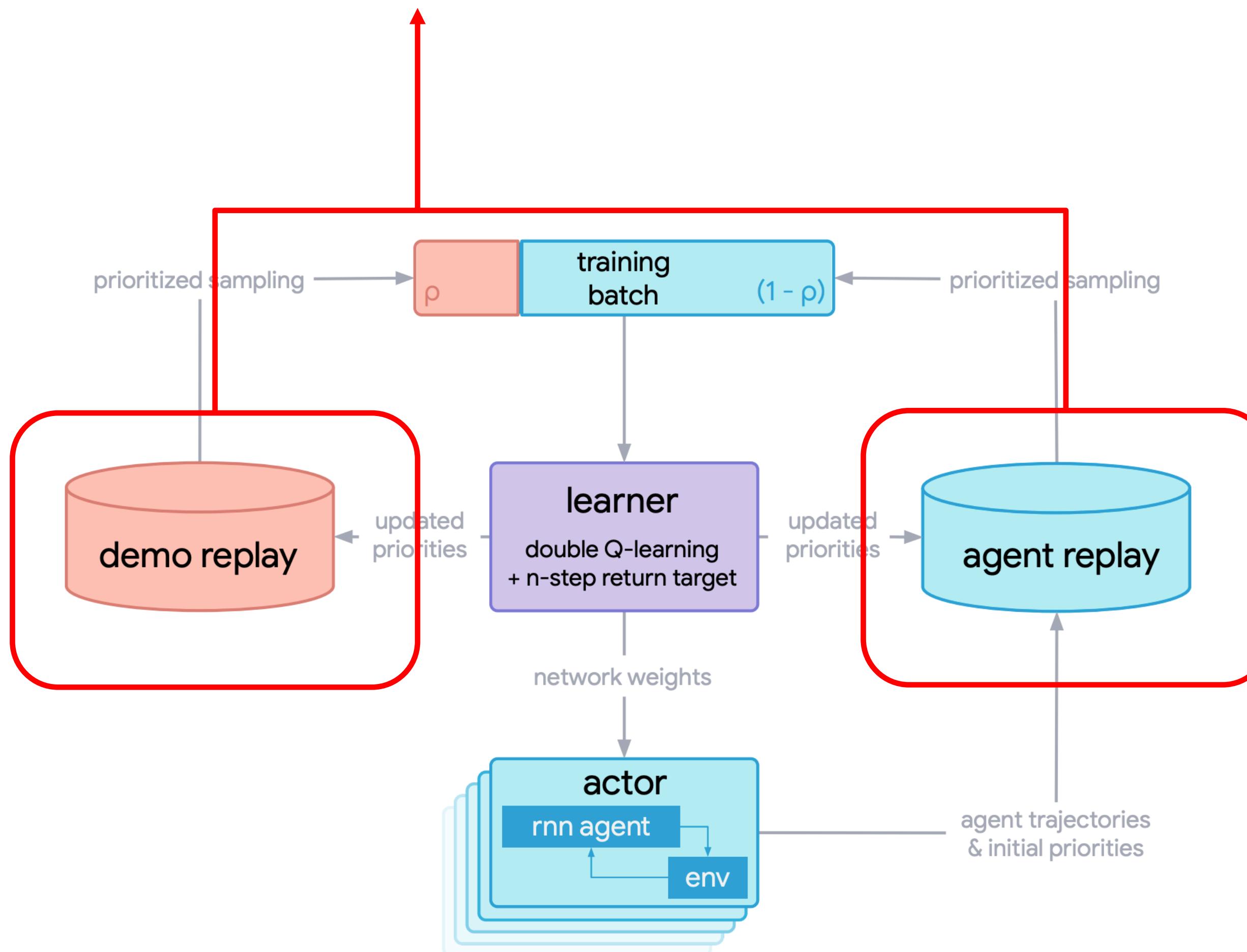
- **Background**
- **Use Demonstrations**
  - Offline – without interactions with environments
  - Online – with interactions with environments
- **Collect Demonstrations**
- **Case Study**
- **Future Directions**

# Use Demo Online

- Utilize demo during online learning directly
- A preceding offline learning stage is optional

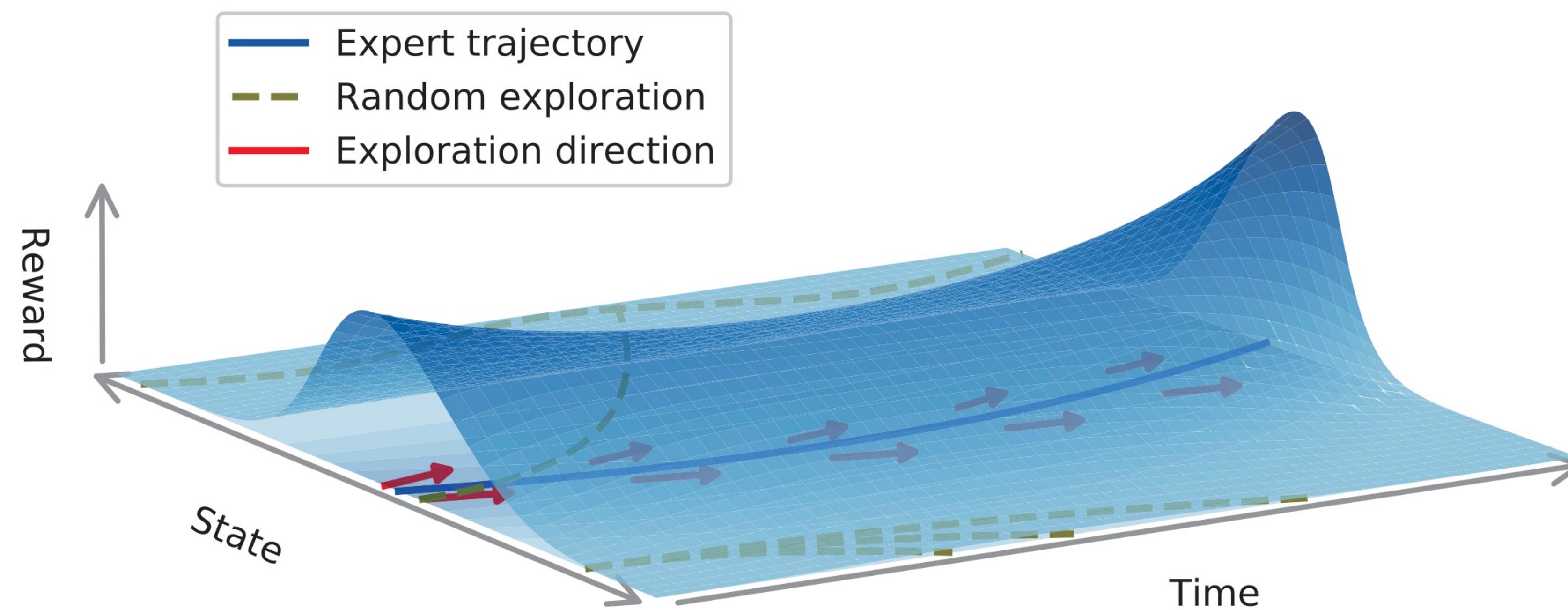
# Demo as Off-Policy Experience

- Add demo into the **replay buffer** of off-policy RL algorithms



# Demo as On-Policy Regularization

- Augment the RL objective with a regularization term
- Encourages the agent to keep the behavior close to the demo



# Demo as On-Policy Regularization

- Augment the RL objective with a regularization term
- Encourages the agent to keep the behavior close to the demo
- **POfD** (Policy Optimization from Demo)

$$\mathcal{L}(\pi_\theta) = -\eta(\pi_\theta) + \lambda_1 D_{JS}(\rho_\theta, \rho_E)$$

, where the  $D_{JS}$  denotes the Jensen-Shannon divergence, the  $\rho_\theta$  and  $\rho_E$  are the occupancy measures of agent policy and expert policy, respectively.

# Demo as On-Policy Regularization

- Augment the RL objective with a regularization term
- Encourages the agent to keep the behavior close to the demo
- **DAPG** (Demo Augmented Policy Gradient)

$$g_{aug} = \sum_{(s,a) \in \rho_\pi} \nabla_\theta \ln \pi_\theta(a|s) A^\pi(s, a) + \sum_{(s,a) \in \rho_D} \nabla_\theta \ln \pi_\theta(a|s) w(s, a)$$

, where  $w(s, a)$  is a weighting function. In practice,  $w(s, a)$  is implemented as a heuristic weighting function:

$$w(s, a) = \lambda_0 \lambda_1^k \max_{(s', a') \in \rho_\pi} A^\pi(s', a') \quad \forall (s, a) \in \rho_D$$

, where  $\lambda_0$  and  $\lambda_1$  are hyperparameters, and  $k$  is the iteration counter.

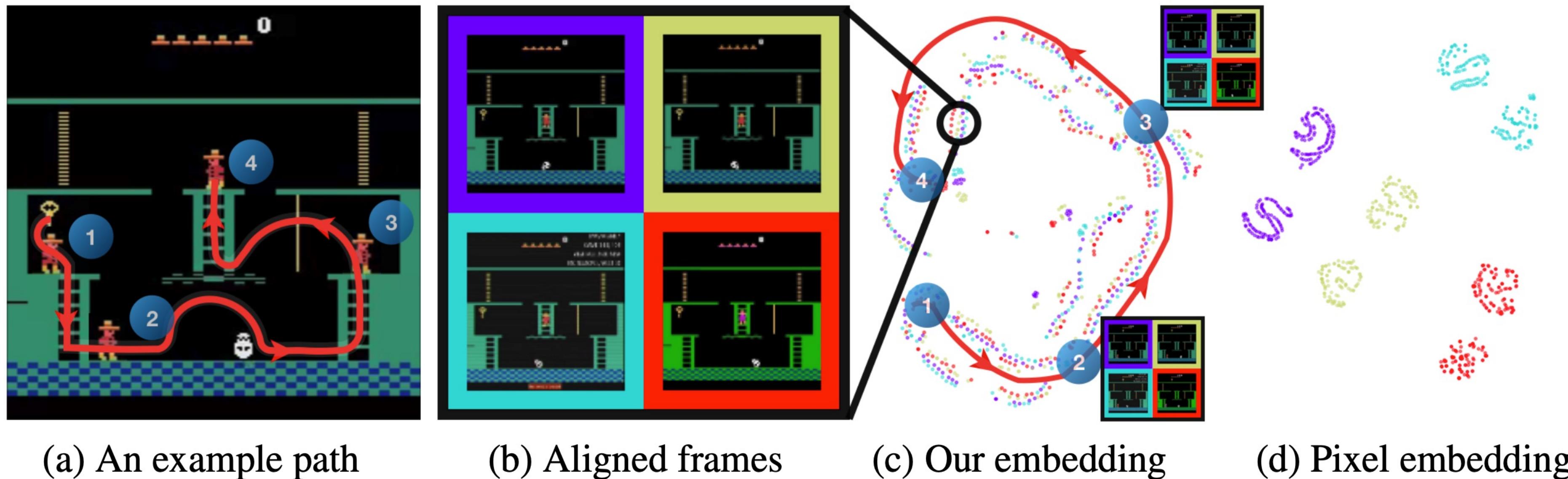
# Demo as Reference for Reward

- Regularization modifies the RL learning objective
- Another natural way is to convert demo into reward, then the demo is automatically incorporated into RL objective
- Two kinds of ideas:
  - Directly define reward based on a single demo trajectory
  - Match the distribution of demonstrations, and use the divergence as the reward

# Define Reward with a Single Demo

- Reward: the **distance to the demo trajectory** in an embedding space
- The embedding space can be **learned**

$$r_{\text{imitation}} = \begin{cases} 0.5 & \text{if } \bar{\phi}(v_{\text{agent}}) \cdot \bar{\phi}(v_{\text{checkpoint}}) > \alpha \\ 0.0 & \text{otherwise} \end{cases}$$



(a) An example path

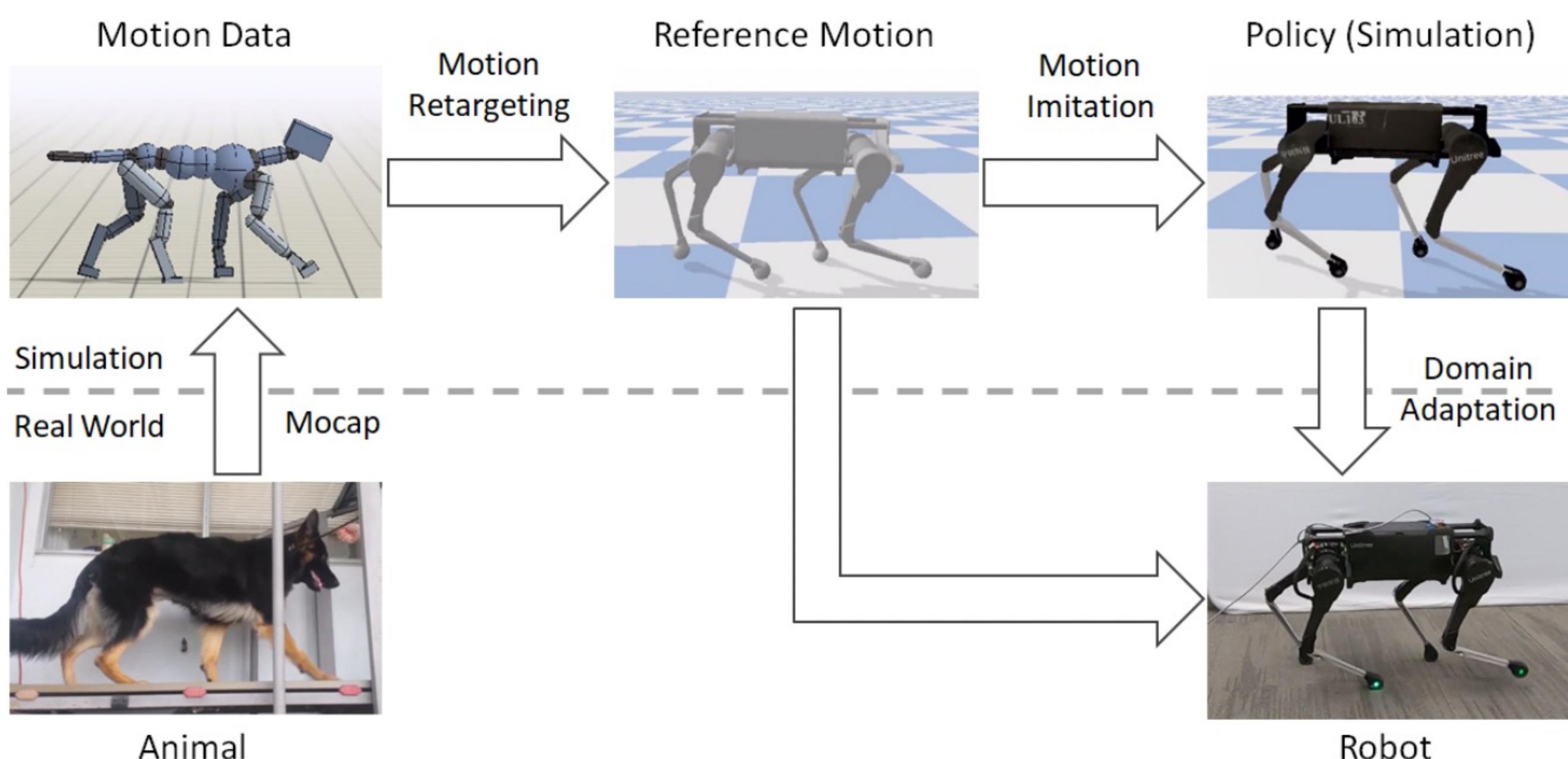
(b) Aligned frames

(c) Our embedding

(d) Pixel embedding

# Define Reward with a Single Demo

- Reward: the **distance to the demo trajectory** in an embedding space
- The embedding space can be **manually defined**



$$\begin{aligned}
 r_t^p &= \exp \left[ -5 \sum_j \|\hat{\mathbf{q}}_t^j - \mathbf{q}_t^j\|^2 \right] & r_t^v &= \exp \left[ -0.1 \sum_j \|\dot{\hat{\mathbf{q}}}_t^j - \dot{\mathbf{q}}_t^j\|^2 \right] \\
 r_t &= w^p r_t^p + w^v r_t^v + w^e r_t^e + w^{rp} r_t^{rp} + w^{rv} r_t^{rv} \\
 r_t^e &= \exp \left[ -40 \sum_e \|\hat{\mathbf{x}}_t^e - \mathbf{x}_t^e\|^2 \right] \\
 r_t^{rp} &= \exp \left[ -20 \|\hat{\mathbf{x}}_t^{\text{root}} - \mathbf{x}_t^{\text{root}}\|^2 - 10 \|\hat{\mathbf{q}}_t^{\text{root}} - \mathbf{q}_t^{\text{root}}\|^2 \right] \\
 r_t^{rv} &= \exp \left[ -2 \|\dot{\hat{\mathbf{x}}}_t^{\text{root}} - \dot{\mathbf{x}}_t^{\text{root}}\|^2 - 0.2 \|\dot{\hat{\mathbf{q}}}_t^{\text{root}} - \dot{\mathbf{q}}_t^{\text{root}}\|^2 \right]
 \end{aligned}$$

# Match the Distribution of Demo

- Reward: the **divergence** between agent trajectories and demo
  - Usually needs to be approximated due to the computation cost
- This is actually the idea behind most **inverse RL** methods

---

**Algorithm 1:** Template for IRL

---

**Input:**  $\mathcal{M}_{RE} = \langle S, A, T, \gamma \rangle$ ,

Set of trajectories demonstrating desired behavior:

$\mathcal{D} = \{(s_0, a_0), (s_1, a_1), \dots, (s_t, a_t)\}, \dots\}$ ,  $s_t \in S$ ,  $a_t \in A$ ,  $t \in \mathbb{N}$ ,

or expert's policy:  $\pi_E$ , and reward function features

**Output:**  $\hat{R}_E$

- 1 Model the expert's observed behavior as the solution of an MDP whose reward function is not known;
- 2 Initialize the parameterized form of the reward function using any given features (linearly weighted sum of feature values, distribution over rewards, or other);
- 3 Solve the MDP with current reward function to generate the learned behavior or policy;
- 4 Update the optimization parameters to minimize the divergence between the observed behavior (or policy) and the learned behavior (policy);
- 5 Repeat the previous two steps till the divergence is reduced to a desired level.

# Match the Distribution of Demo

- **GAIL** (Generative adversarial imitation learning)
  - Reward is JS divergence, implemented in a way similar to GANs
  - Generator: policy  $\pi(a|s)$
  - Discriminator: predicts  $(s, a)$  from agent or demo

---

**Algorithm 1** Generative adversarial imitation learning

---

- 1: **Input:** Expert trajectories  $\tau_E \sim \pi_E$ , initial policy and discriminator parameters  $\theta_0, w_0$
- 2: **for**  $i = 0, 1, 2, \dots$  **do**
- 3:   Sample trajectories  $\tau_i \sim \pi_{\theta_i}$
- 4:   Update the discriminator parameters from  $w_i$  to  $w_{i+1}$  with the gradient

$$\hat{\mathbb{E}}_{\tau_i} [\nabla_w \log(D_w(s, a))] + \hat{\mathbb{E}}_{\tau_E} [\nabla_w \log(1 - D_w(s, a))] \quad (17)$$

- 5:   Take a policy step from  $\theta_i$  to  $\theta_{i+1}$ , using the TRPO rule with cost function  $\log(D_{w_{i+1}}(s, a))$ . Specifically, take a KL-constrained natural gradient step with

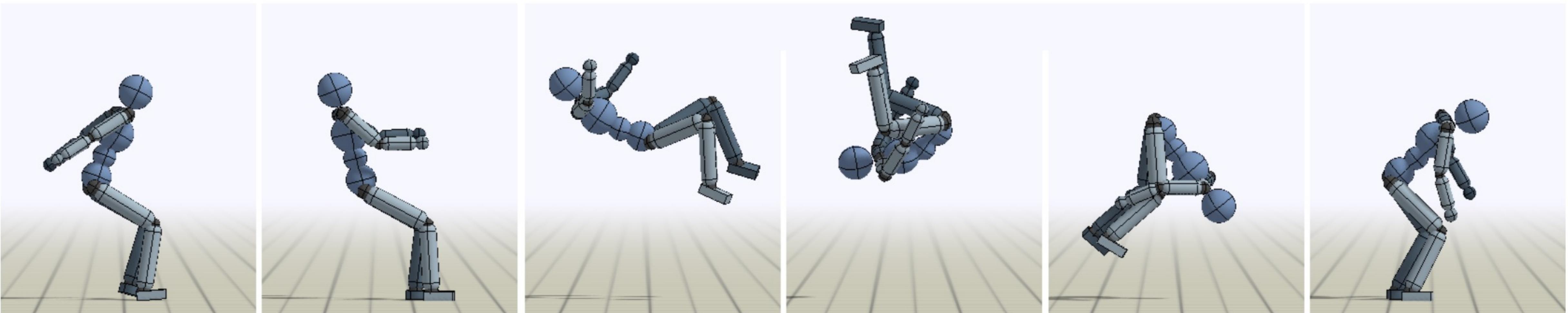
$$\hat{\mathbb{E}}_{\tau_i} [\nabla_\theta \log \pi_\theta(a|s) Q(s, a)] - \lambda \nabla_\theta H(\pi_\theta), \quad (18)$$

where  $Q(\bar{s}, \bar{a}) = \hat{\mathbb{E}}_{\tau_i} [\log(D_{w_{i+1}}(s, a)) \mid s_0 = \bar{s}, a_0 = \bar{a}]$

- 6: **end for**
-

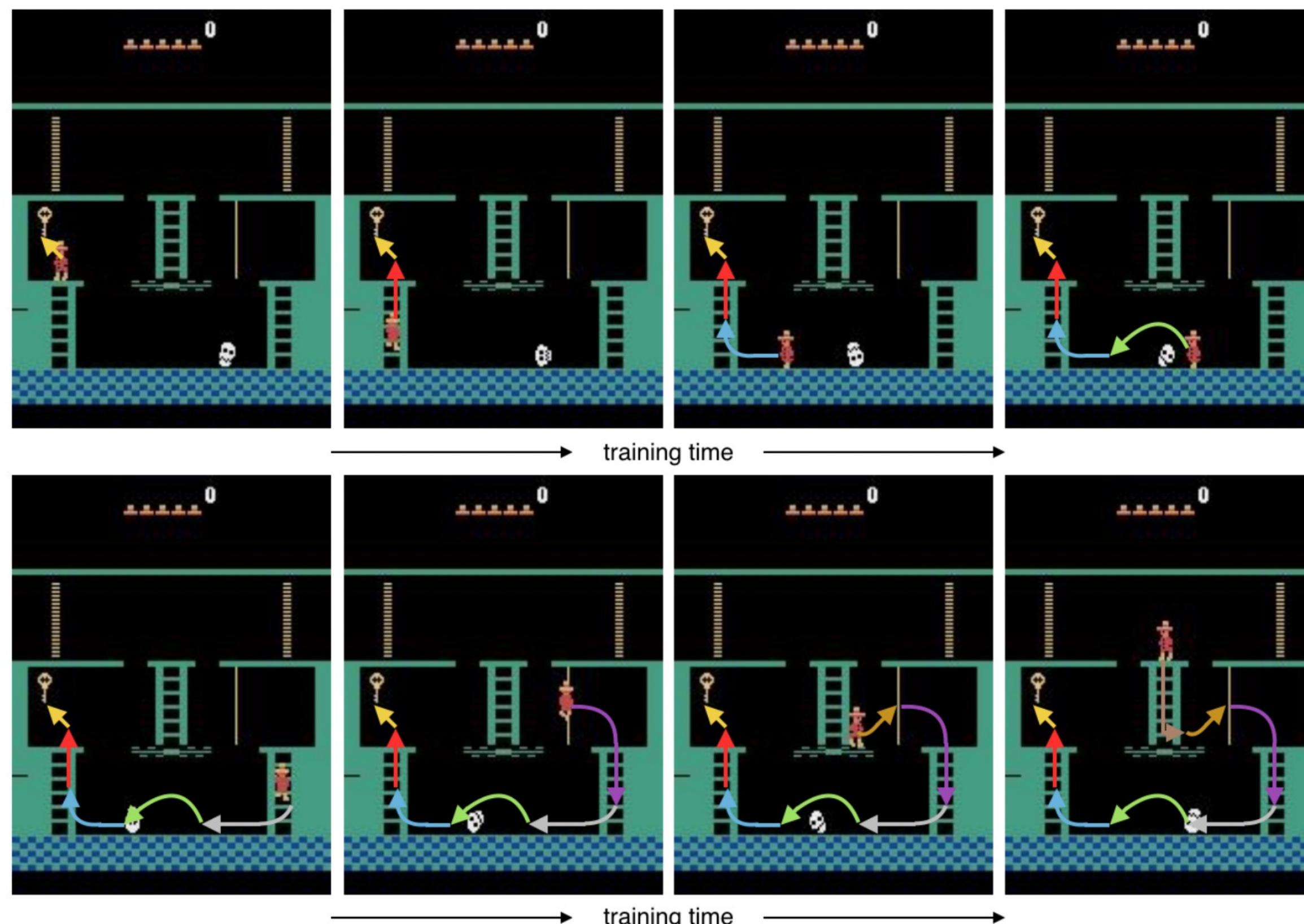
# Demo as Curriculum of Start States

- Assume the environment is a simulator that we can fully control
- **Reset to the states in demo, and start to explore from there**
  - Uniformly sample states in demo as the start states



# Demo as Curriculum of Start States

- Assume the environment is a simulator that we can fully control
- **Reset to the states in demo, and start to explore from there**
  - Or start from the end of the demo, and gradually go backwards



# Outline

- **Background**
- **Use Demonstrations**
  - Offline – without interactions with environments
  - Online – with interactions with environments
- **Collect Demonstrations**
- **Case Study**
- **Future Directions**

# Collect Demo

- Use Embodied AI as an example domain
  - Demo can be collected in various ways
    - By human or by robots
    - In simulators or in the real world
  - One of the most popular domains to use demonstrations
- Focus on acquiring expert demonstrations
  - Non-expert demo are relatively easy to get



# Collect Demo

- **Embodiment:** the “physical body” of the agent
- **Operator:** the “brain” to control the agent

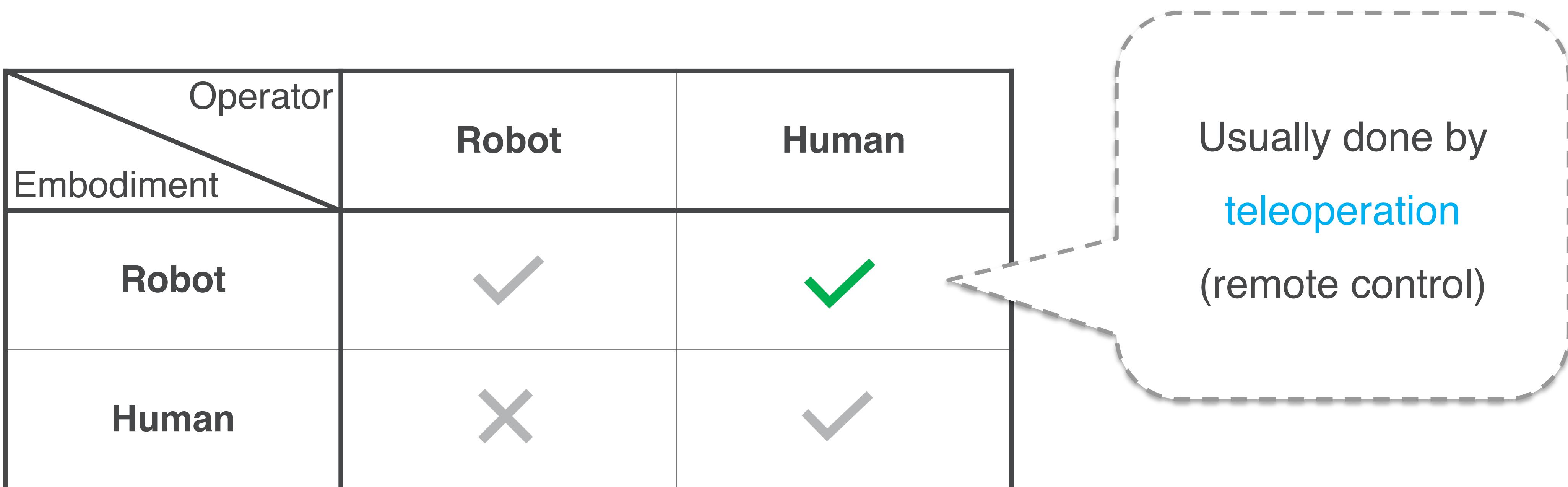
Operator		
Embodiment		
Robot		
Human		

A 3x3 matrix diagram illustrating the relationship between Operator, Embodiment, Robot, and Human. The rows are labeled from top to bottom: Operator, Embodiment, Robot, and Human. The columns are labeled from left to right: Operator, Embodiment, Robot, and Human. The diagonal from top-left to bottom-right is labeled "Robot". The cell at the intersection of "Robot" (row) and "Robot" (column) contains a green checkmark. The cells at the intersections of "Robot" (row) and "Human" (column), and "Human" (row) and "Robot" (column) also contain green checkmarks. The cell at the intersection of "Human" (row) and "Human" (column) contains a large red X.

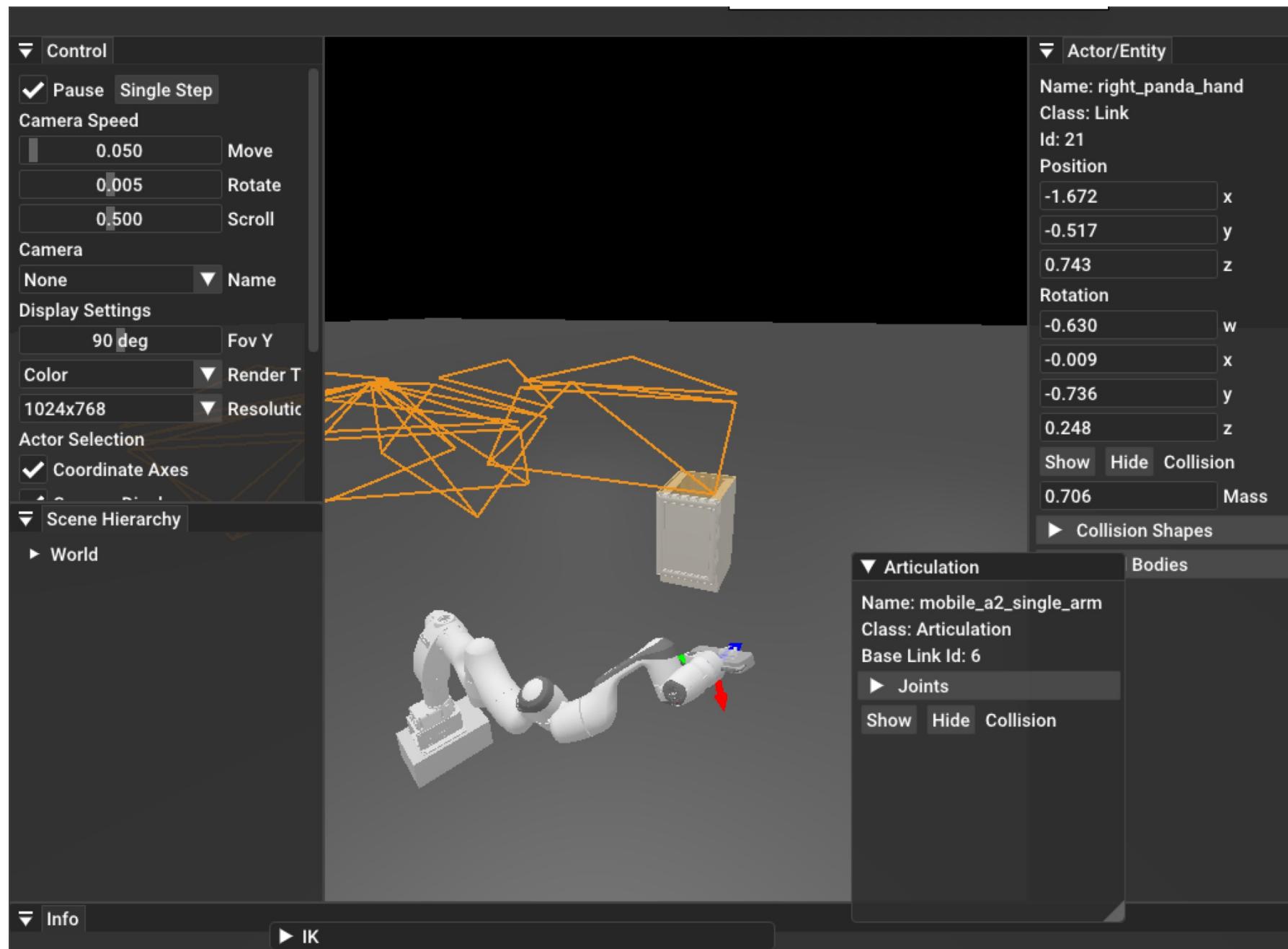


# Collect Demo

- **Embodiment:** the “physical body” of the agent
- **Operator:** the “brain” to control the agent

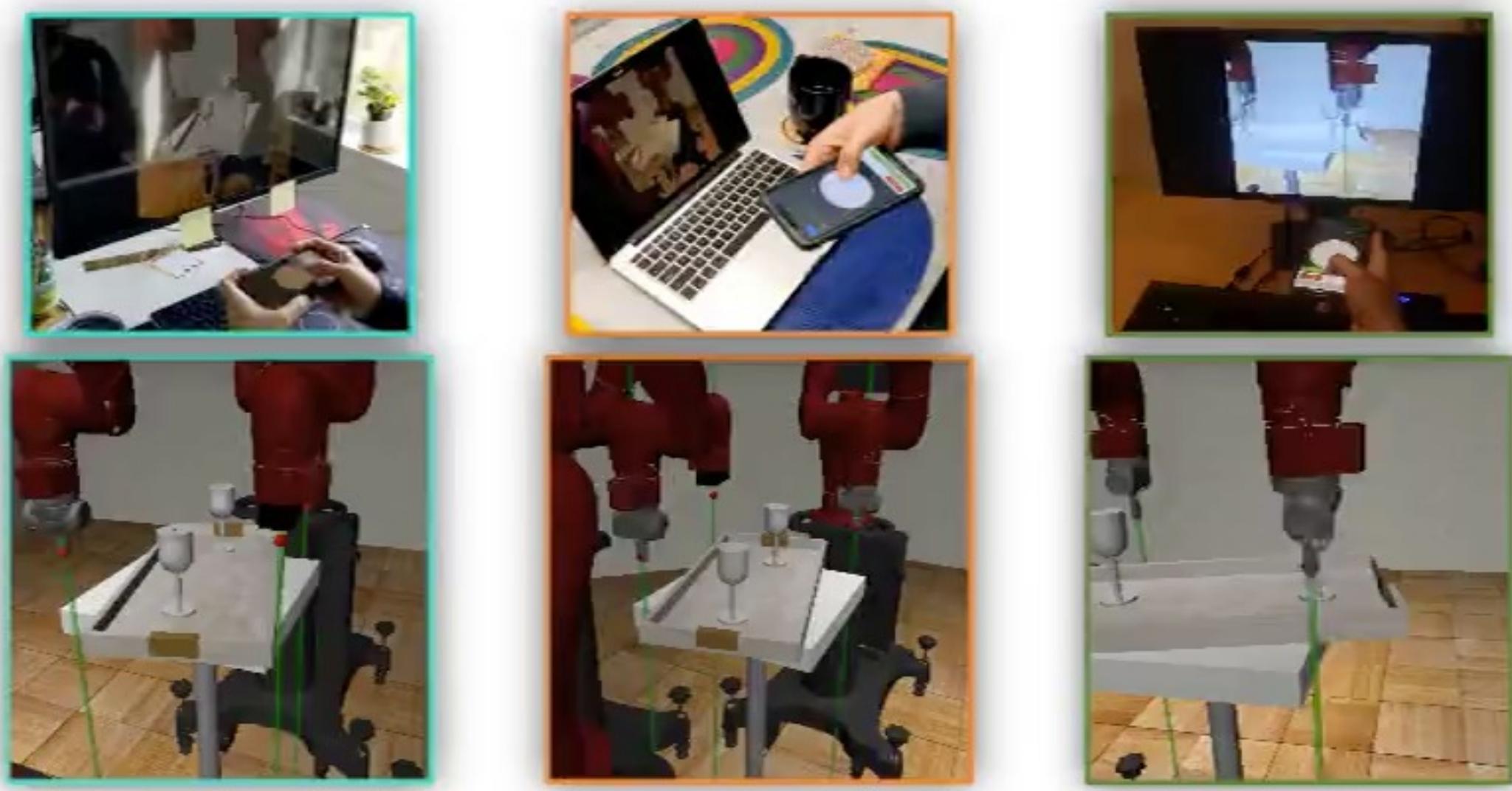


# Teleoperation – Basic Devices



Keyboard + Mouse  
(Many Simulation Environments)

Multi-Arm RoboTurk (MART): Collaborative Teleoperation



Smartphone (RoboTurk)

# Teleoperation – Virtual Reality

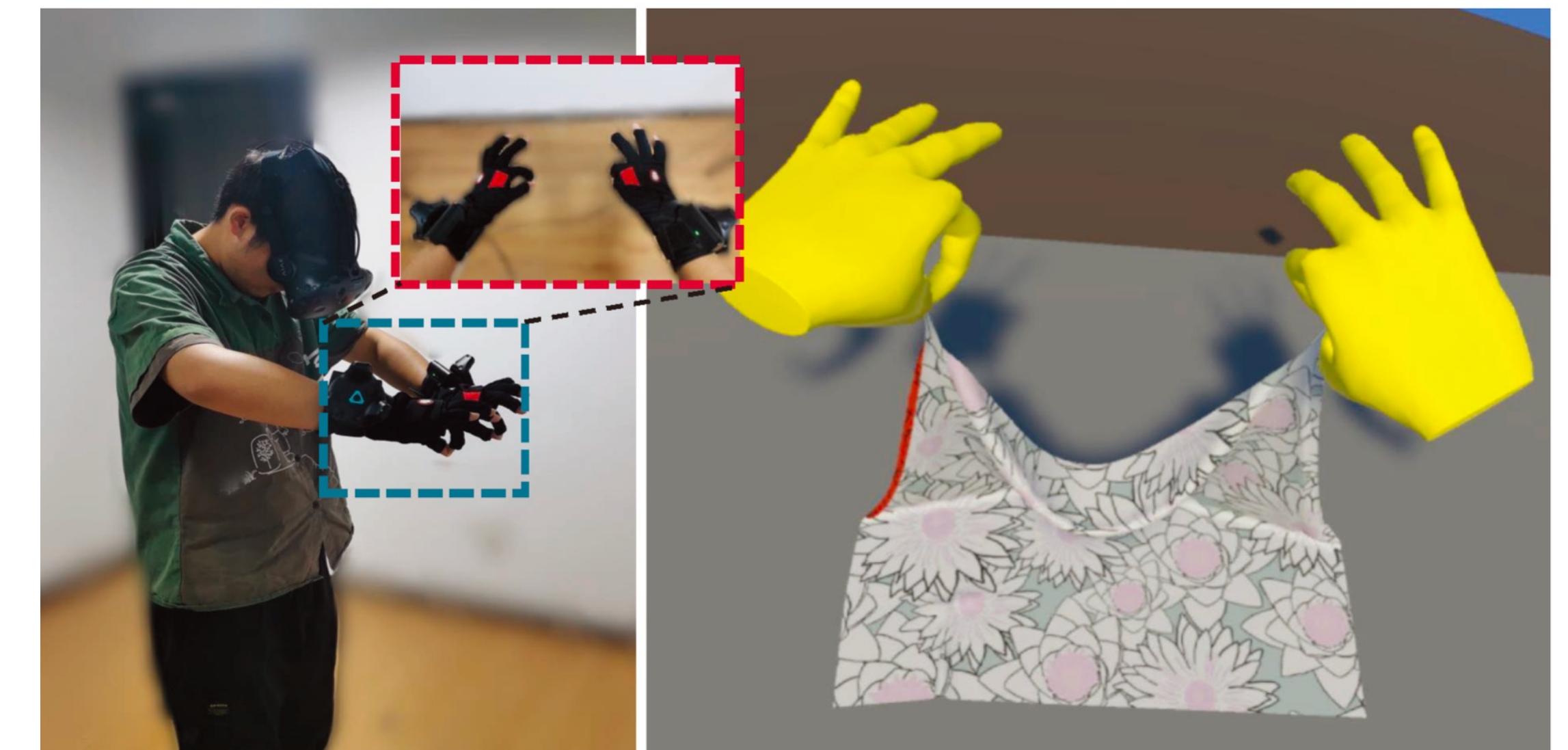
- VR is widely adopted in many simulation environments

## VR Interface

Cooking Onion



Headset + Regular Hand Controllers  
(iGibson 2)



Headset + Motion Capture Gloves  
(RFUniverse)

# Teleoperation – Virtual Reality

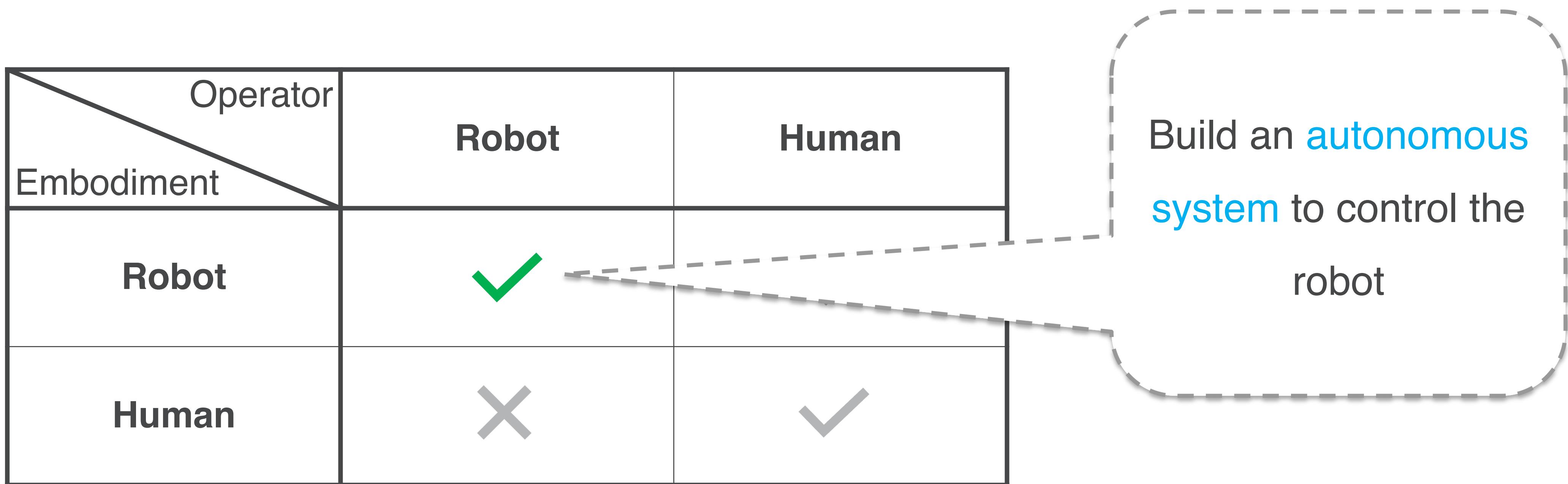
- VR is also used when collecting robot demo in the real world



VR remotes + joystick  
(RT-1)

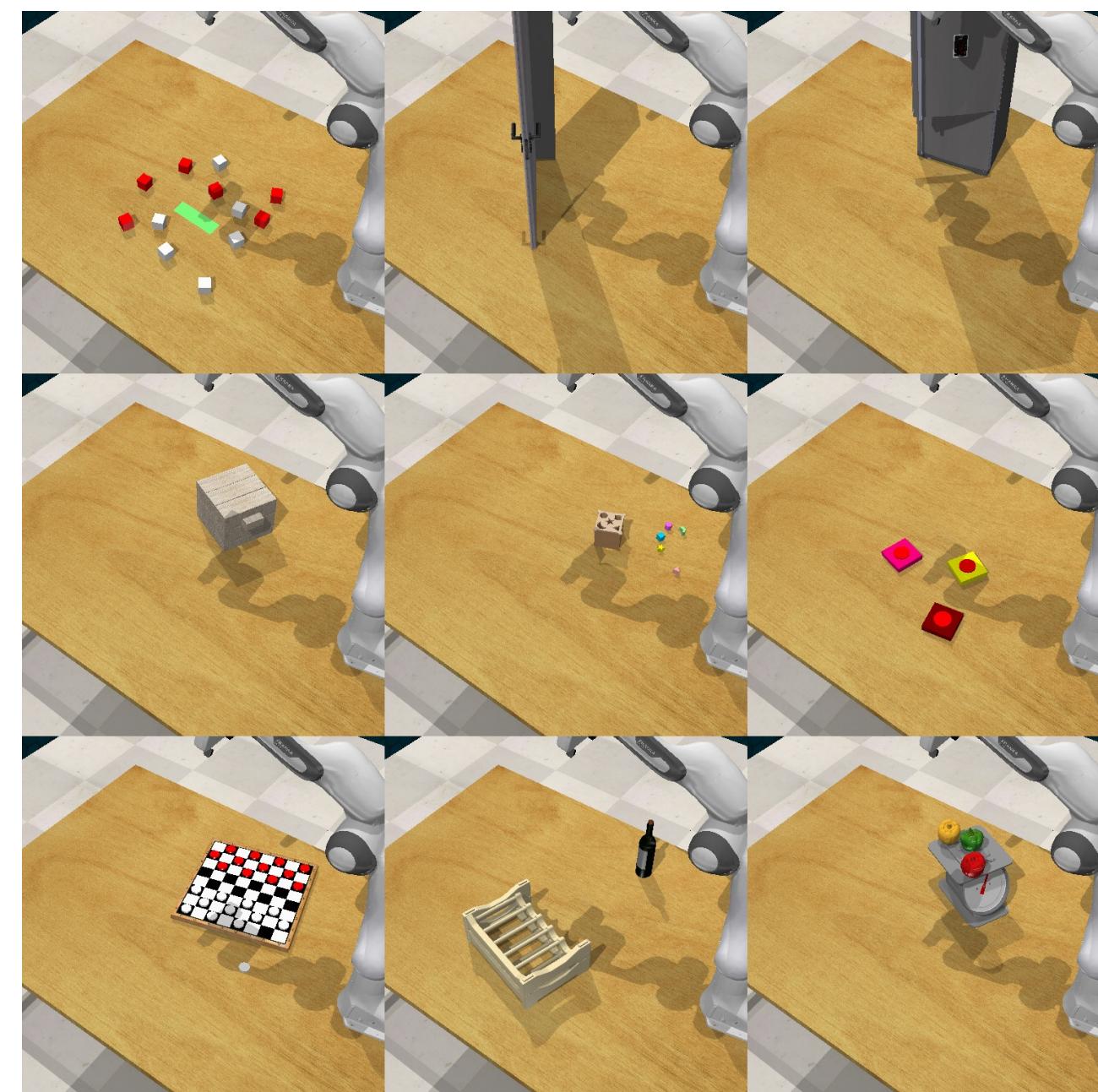
# Collect Demo

- **Embodiment:** the “physical body” of the agent
- **Operator:** the “brain” to control the agent



# Autonomous System – Planning

- Use **planners** to generate demo (assume world model is known)

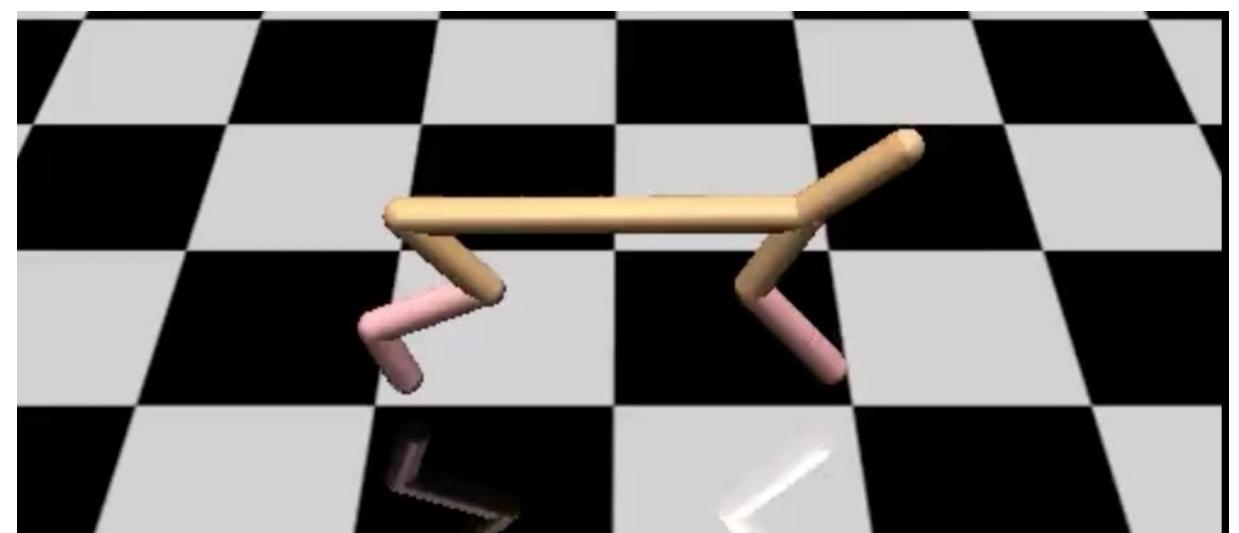


Symbolic Planner for High-level Tasks  
(ALFRED)

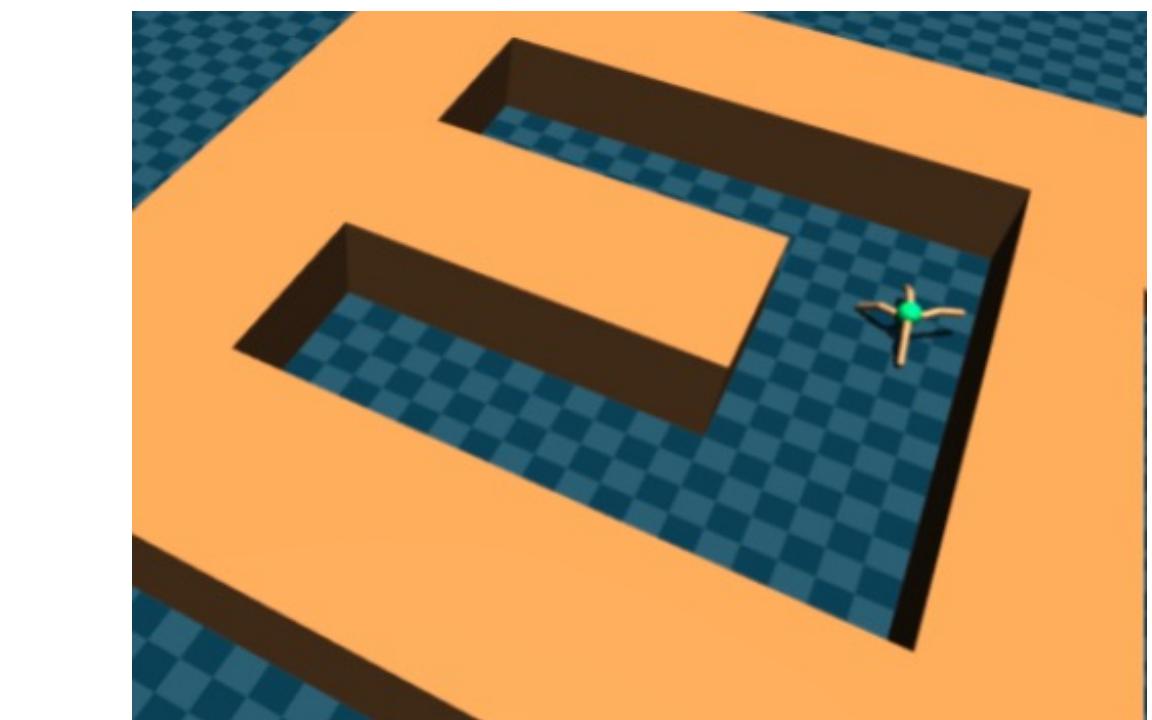
Motion Planner for Low-level Tasks  
(RLBench)

# Autonomous System – Learning

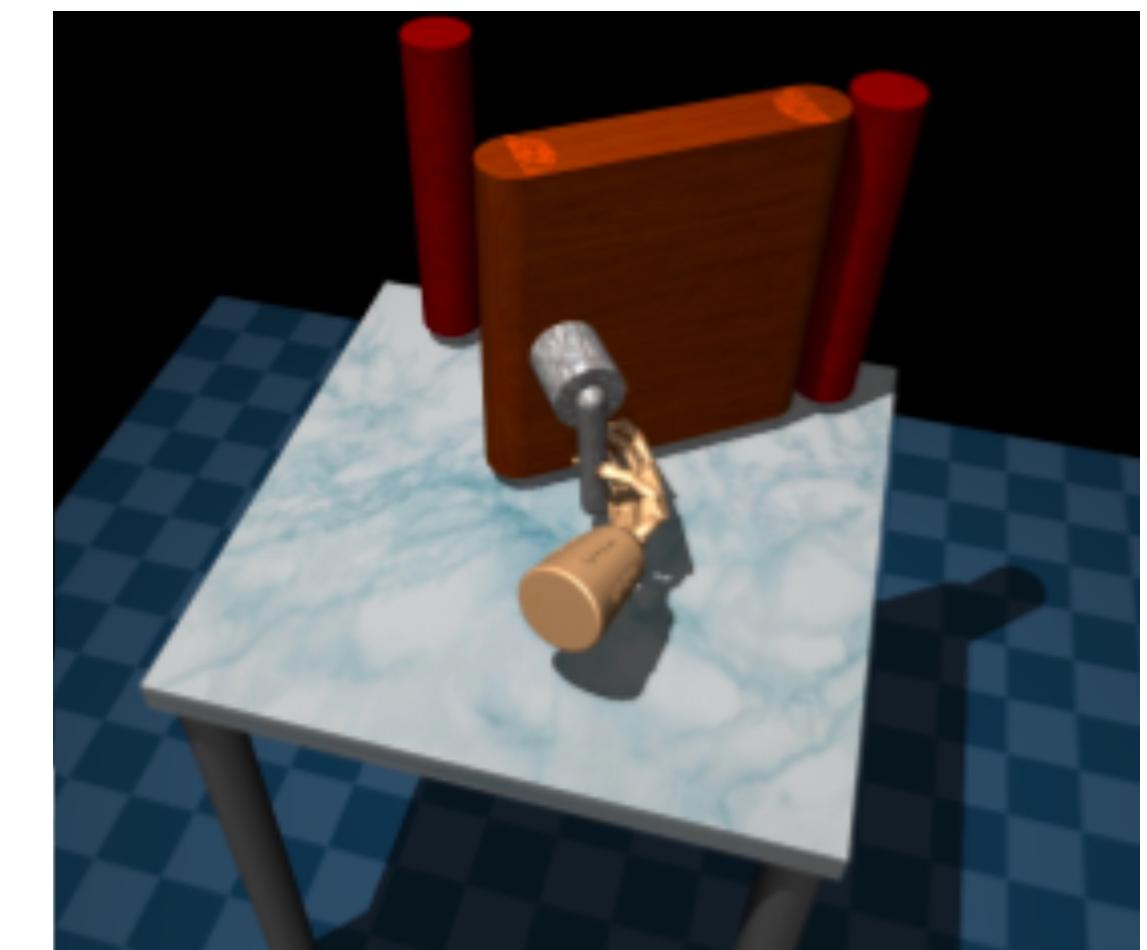
- Learning-based methods
- Design different methods for different tasks



RL alone is already enough  
for many tasks



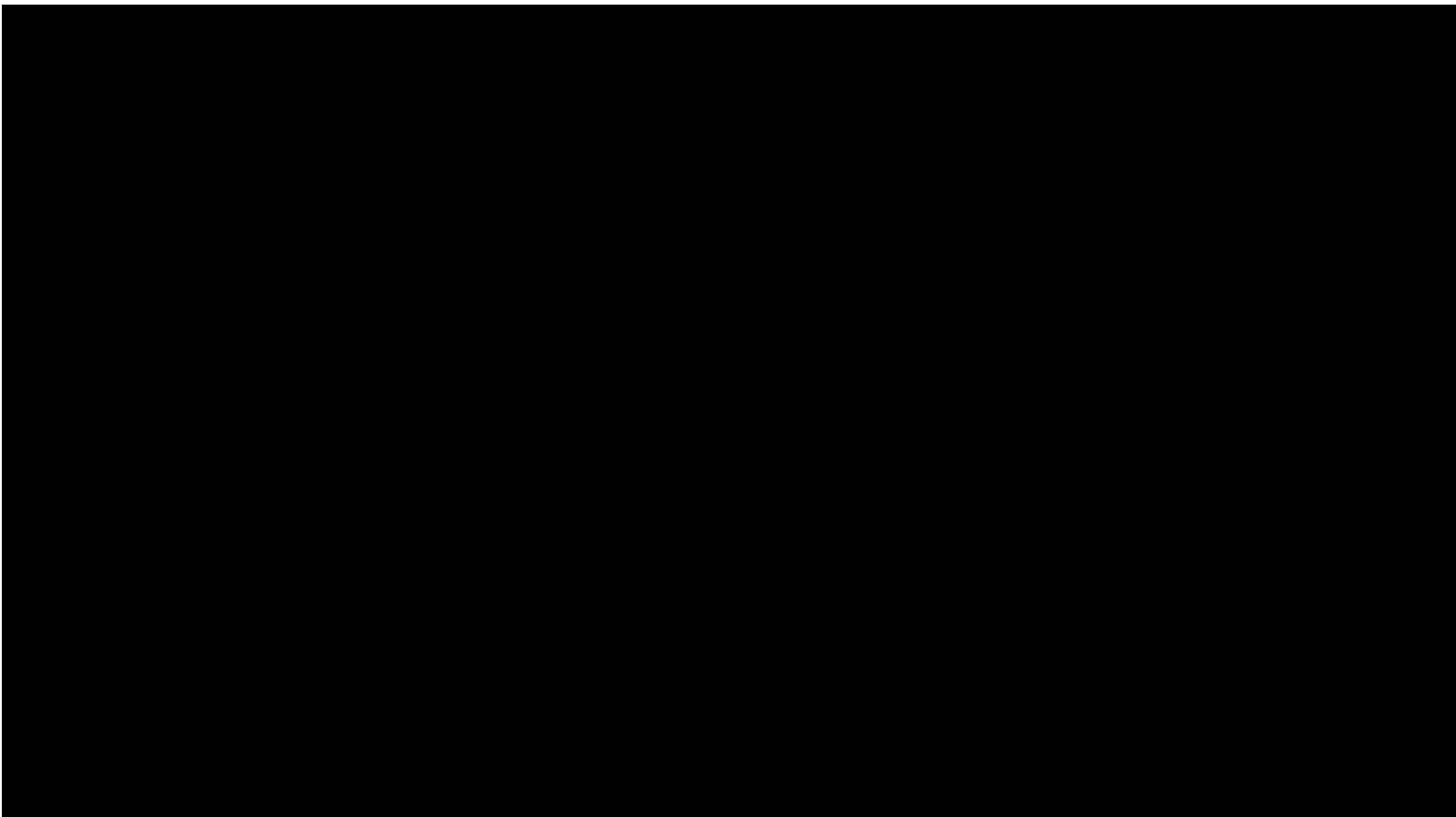
Waypoint generator + Goal-reaching policy from RL



DAPG + a few human demo

# Autonomous System – Self-supervised

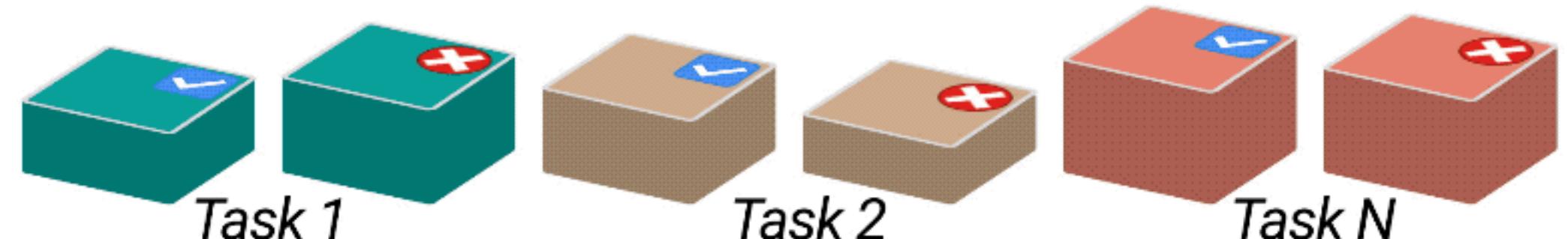
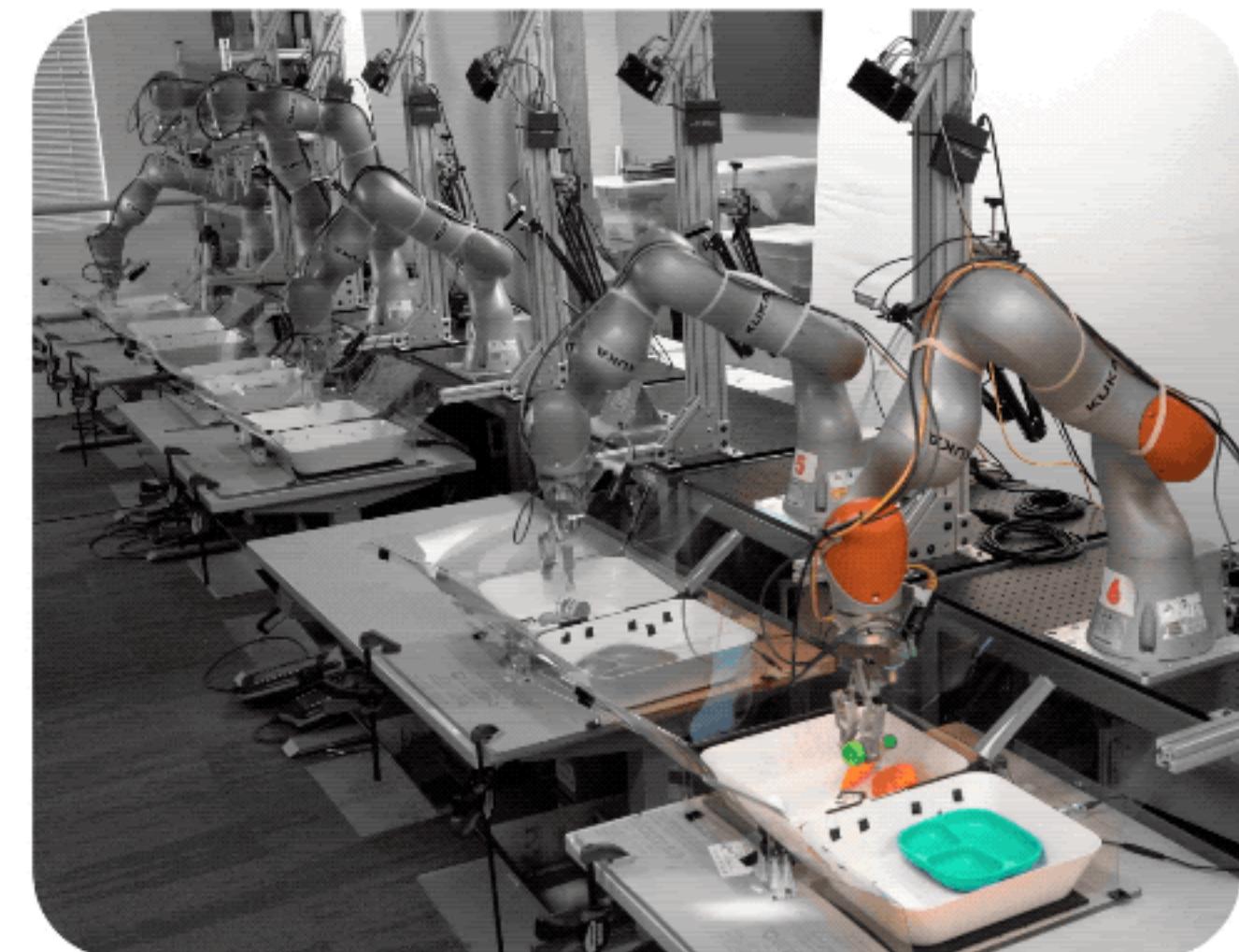
- **Self-supervised system** in the real world



- QT-opt
- Run RL in the real world
- Vision system to get a sparse reward

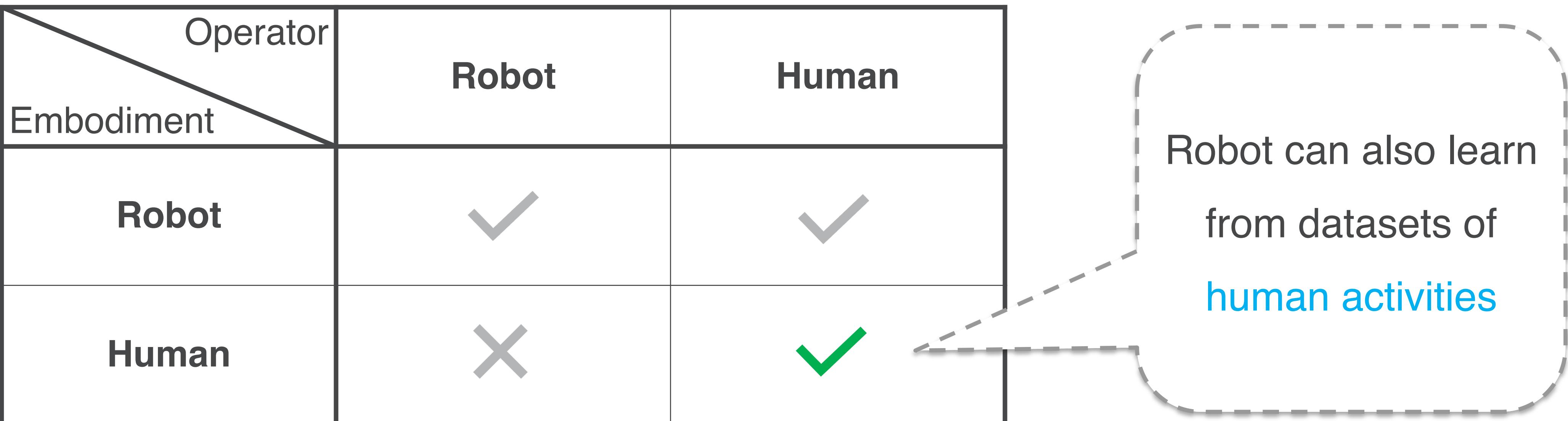
# Autonomous System – Self-supervised

- MT-Opt
  - Reset by special boxes
  - Success detectors trained on data from all tasks
  - Use the solutions to easier tasks to bootstrap learning of more complex tasks



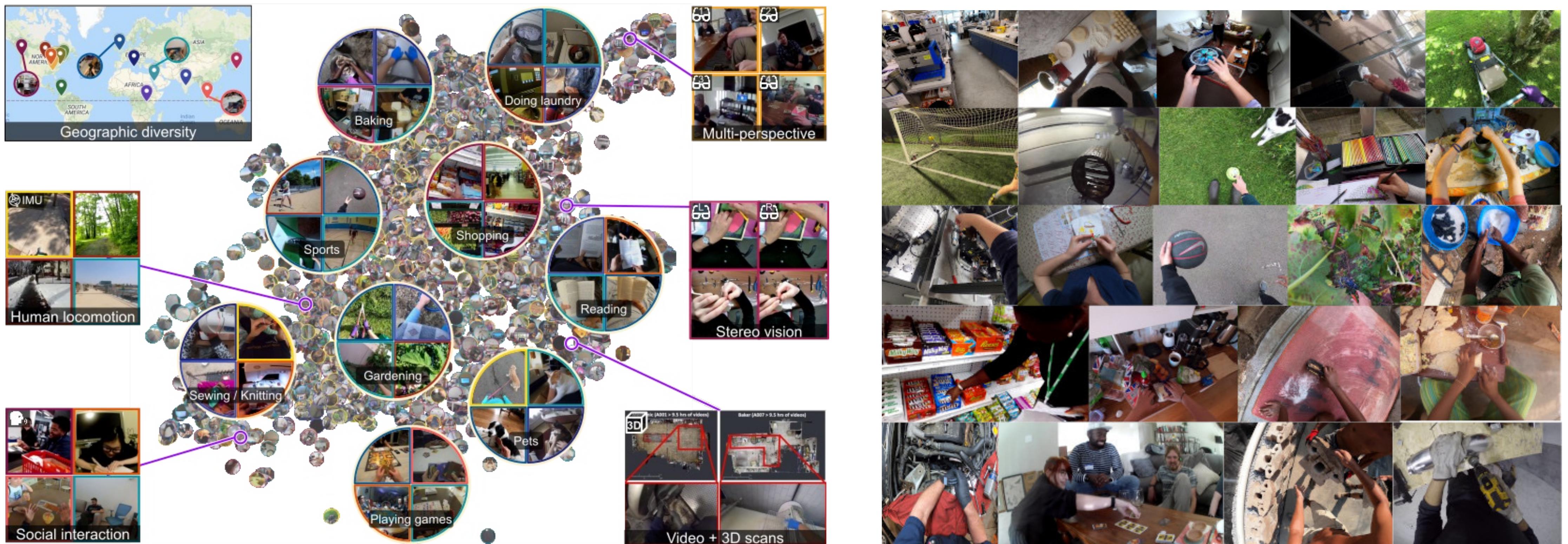
# Collect Demo

- **Embodiment:** the “physical body” of the agent
- **Operator:** the “brain” to control the agent



# Human Demo Datasets

- Ego4D
  - Ego-centric, multi-modal dataset of human activities



# Human Demo Datasets

- RoboTube
  - Human video dataset + its digital twin in simulation environment



# Outline

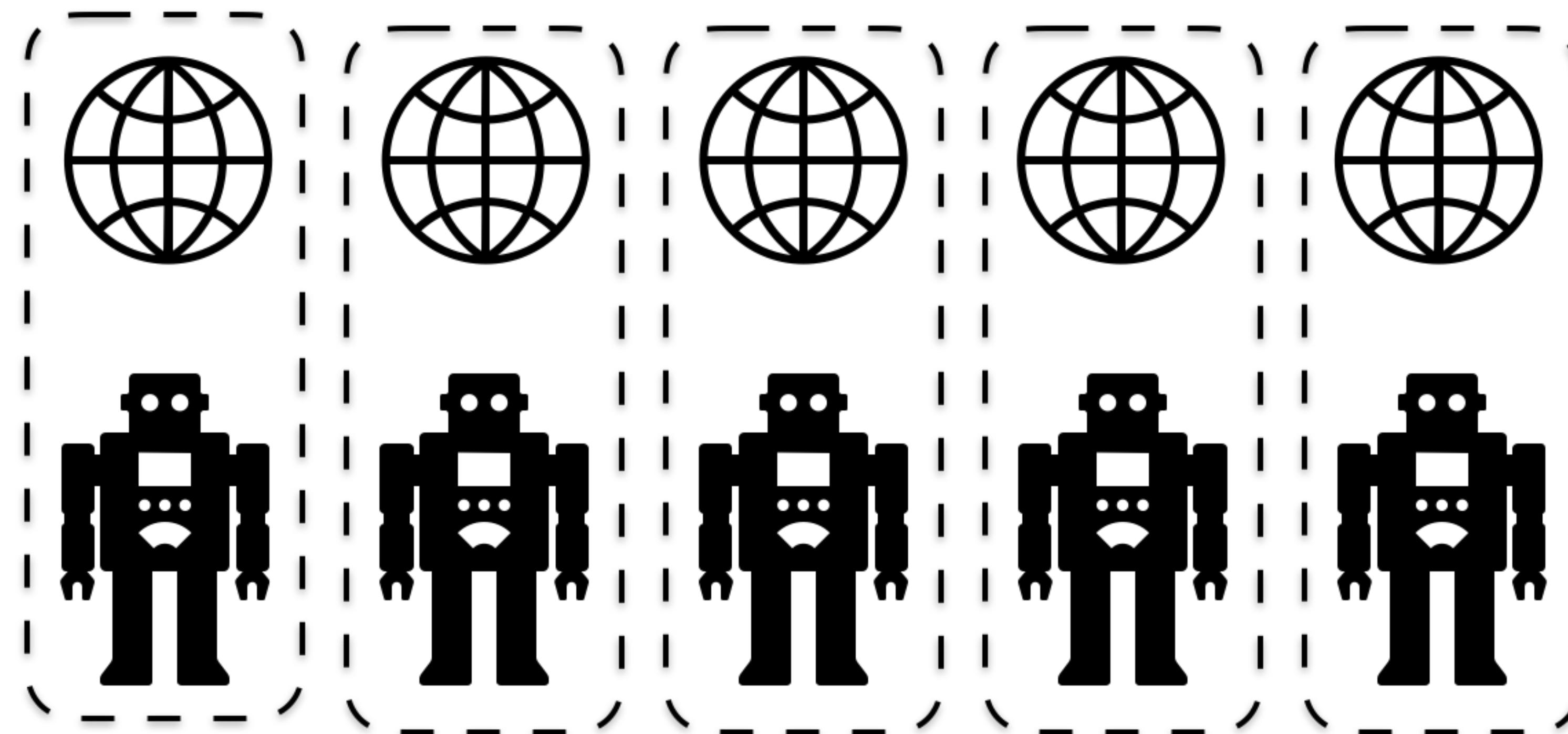
- **Background**
- **Use Demonstrations**
  - Offline – without interactions with environments
  - Online – with interactions with environments
- **Collect Demonstrations**
- **Case Study**
- **Future Directions**

# ManiSkill Benchmark



# Collect Demo

- Divide and Conquer RL
  - For each task, we design a reward template ([shared across all objects in the same task](#))
  - Train an RL agent on each object individually



# Use Demo

- Imitation Learning: Behavior Cloning (BC)
- Offline RL: Batch-Constrained Q-Learning (BCQ)
- Demo as Regularization: TD3 with Behavior Cloning (TD3+BC)

Algorithm	BC				BCQ		TD3+BC	
Architecture	PointNet		PointNet + Transformer		PointNet + Transformer		PointNet + Transformer	
Split	Training	Test	Training	Test	Training	Test	Training	Test
OpenCabinetDoor	0.18±0.02	0.04±0.03	0.30±0.06	0.11±0.02	0.16±0.02	0.04±0.02	0.13±0.03	0.04±0.02
OpenCabinetDrawer	0.24±0.03	0.11±0.03	0.37±0.06	0.12±0.02	0.22±0.04	0.11±0.03	0.18±0.02	0.10±0.02
PushChair	0.11±0.02	0.09±0.02	0.18±0.02	0.08±0.01	0.11±0.01	0.08±0.01	0.12±0.02	0.08±0.01
MoveBucket	0.03±0.01	0.02±0.01	0.15±0.01	0.08±0.01	0.08±0.01	0.06±0.01	0.05±0.01	0.03±0.01

# Outline

- **Background**
- **Use Demonstrations**
  - Offline – without interactions with environments
  - Online – with interactions with environments
- **Collect Demonstrations**
- **Case Study**
- **Future Directions**

# Future Directions

- Scale up demo collection
  - Teleoperation-based approaches
    - Pros: provides high-quality and diverse demo
    - Cons: very costly, hard to scale up
  - Learning-based autonomous data collection pipelines
    - Pros: generates unlimited data, easier to scale up
    - Cons: not strong enough to solve some complex tasks, quality of demo is an issue
  - Combine them together?
    - Human-in-the-loop autonomous system which improves itself over time?
    - Still a long way to go...

# Future Directions

- What kinds of demo do we really need?
  - Quality: always need near-optimal demo?
  - Modality: learn robot policies from videos and language descriptions?
  - Embodiment: learn robot policies from human?
  - Content: always need actions? rewards?
  - ...

# Future Directions

- Combine offline learning from demo with online learning
  - Though there have been several preliminary attempts, the problem is still not solved yet
  - Demo can come in different forms and different qualities
  - Solutions might need to be designed for each different scenario
  - An interesting problem to study
    - Low to learn from non-optimal, cross-domain, partially observed demonstrations
    - This kind of demo is what we usually get in the real world

**Thank you!**