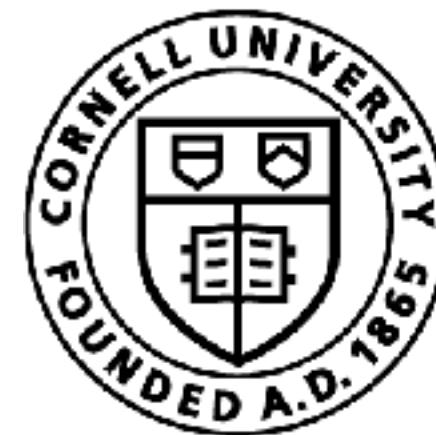


Generative World Models v/s Sim2Real

Sanjiban Choudhury



Cornell Bowers CIS
Computer Science

Today's class

- Practical MBRL
- Learning Models: The DREAMER algorithm
- Leveraging Physics: Sim2Real
Case study: OpenAI Dactyl Hand

Modelling complex tasks from video input



Challenges with learning complex models

Challenge 1: Can't see state, only get high-dimensional observations

Challenge 2: Planning with complex dynamics

Two strategies

“Physics the sh*t out it!”

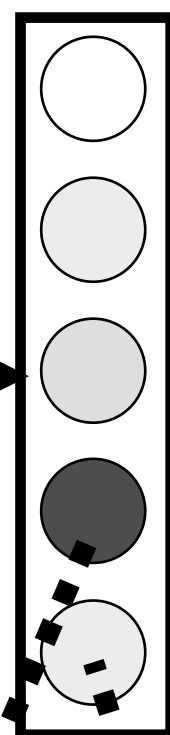
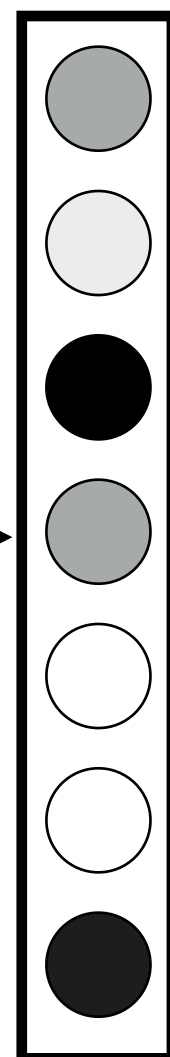
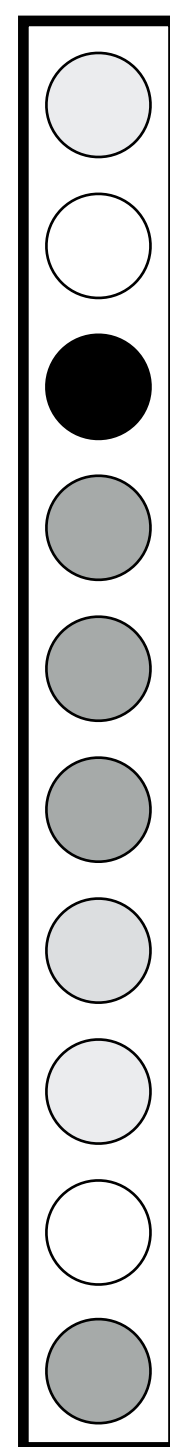
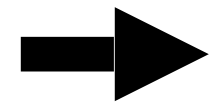
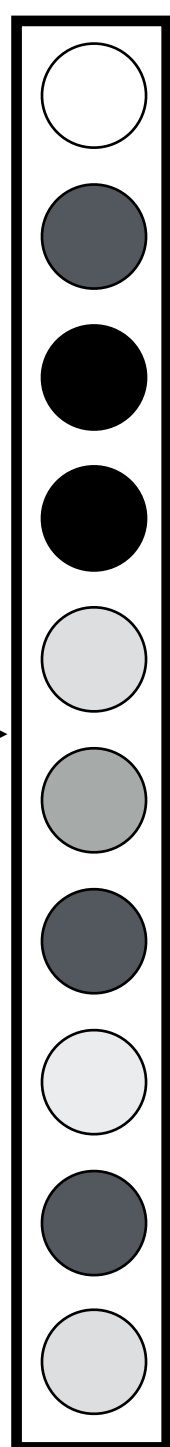
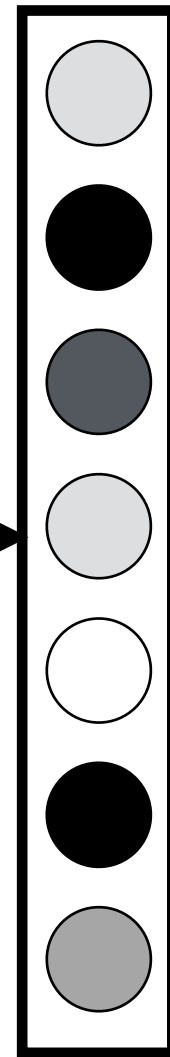
“Be lazy and use ML”

How can we learn latent low-dimensional state from high-dimensional observations?

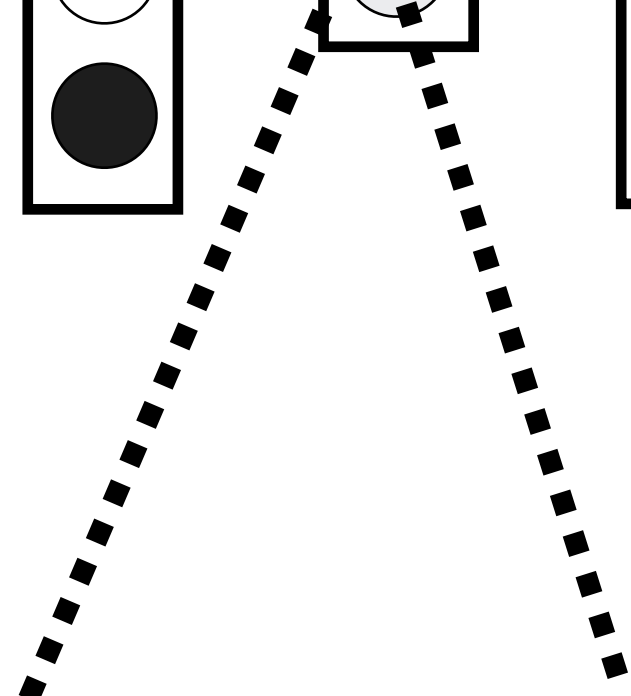
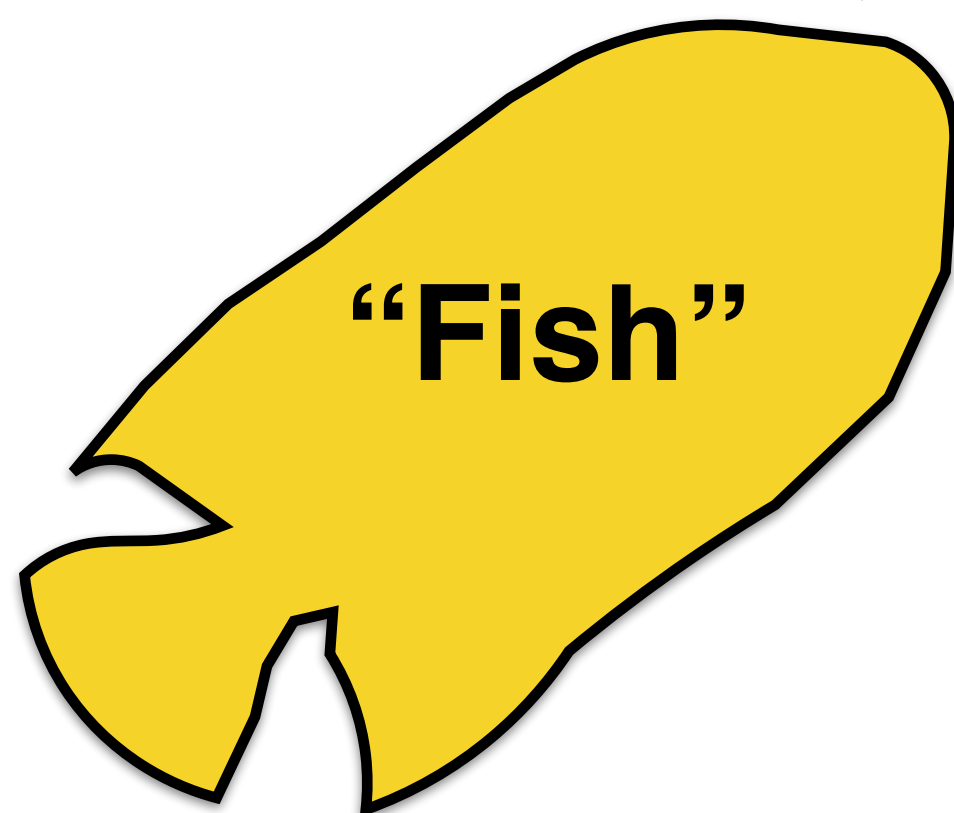
Idea: Use “auto-encoder” trick from
computer vision

\mathbf{X} 

Image

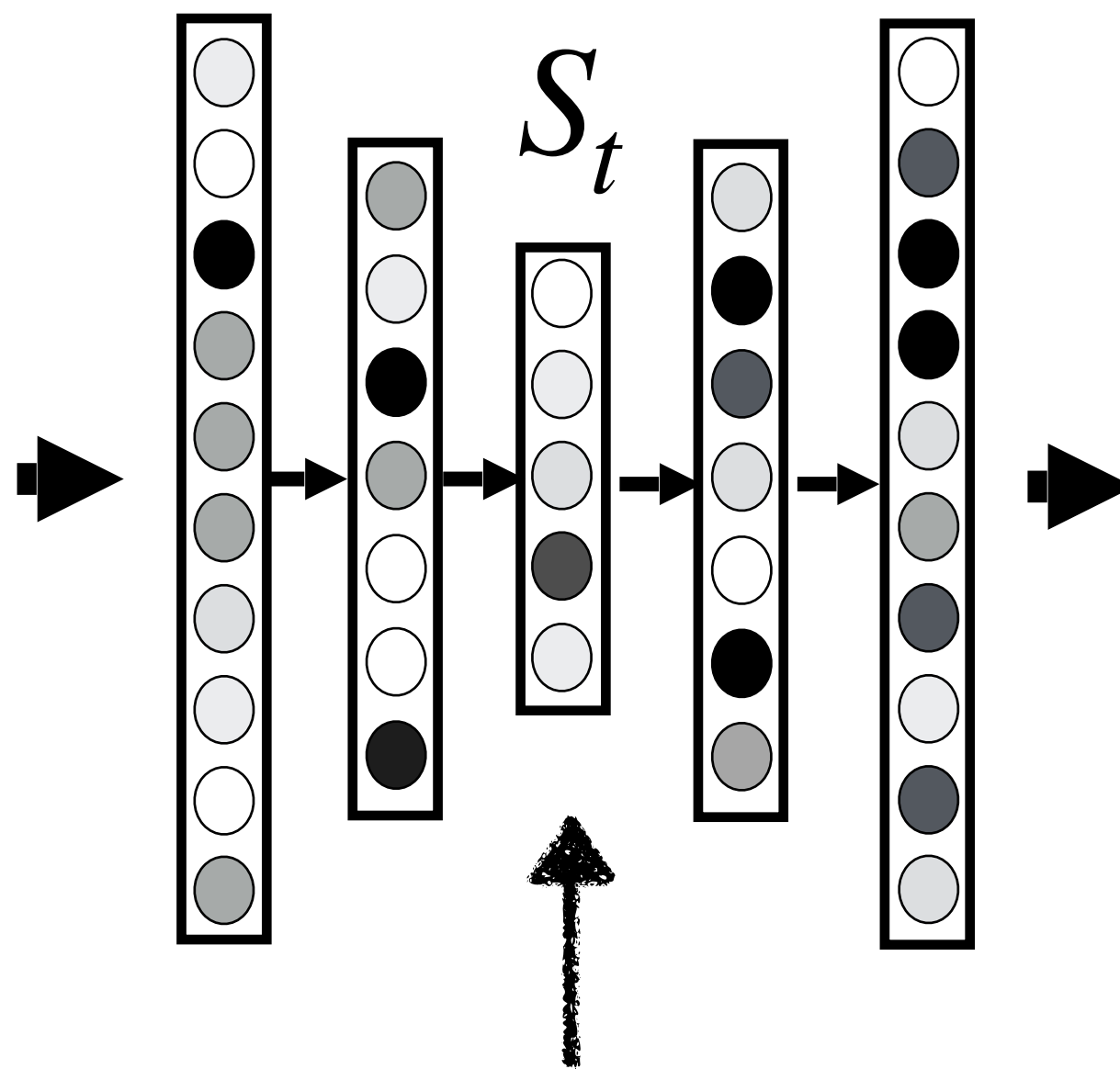
 \mathcal{F}  $\hat{\mathbf{X}} = \mathcal{F}(\mathbf{X})$ 

Reconstructed image

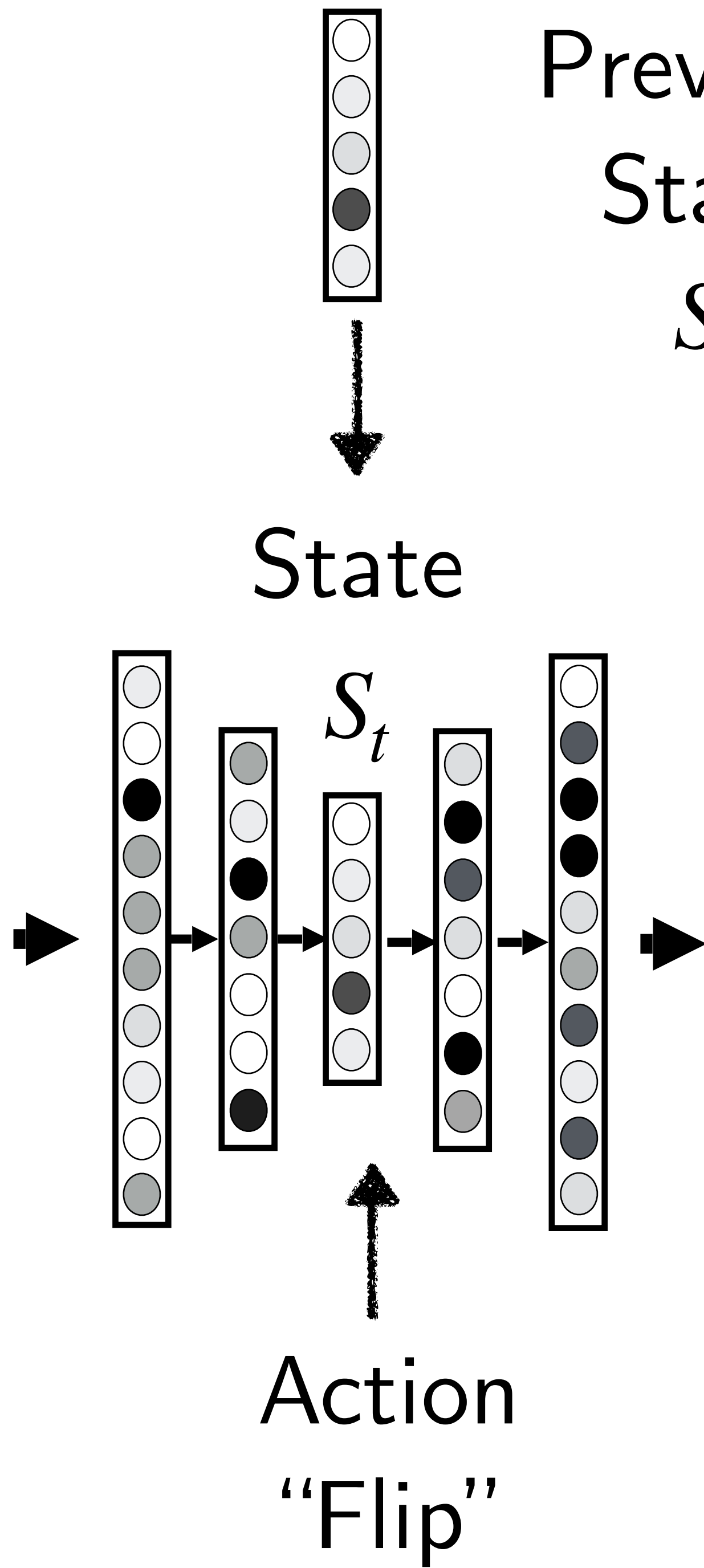




State



Action
"Flip"

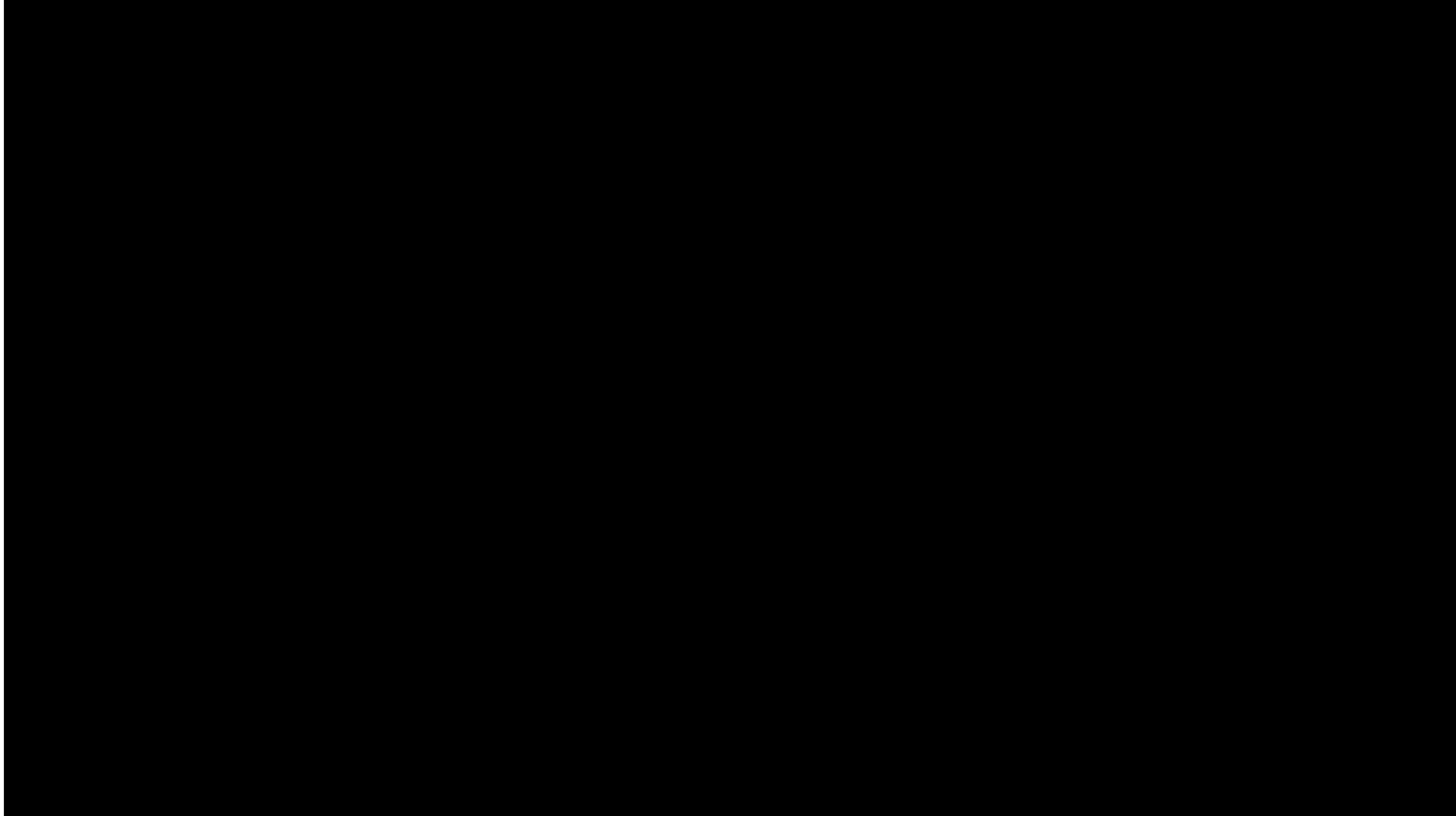


Today's class

- ☑ Practical MBRL
(Only observations, complex dynamics)
- ☐ Learning Models: The DREAMER algorithm
- ☐ Leveraging Physics: Sim2Real
Case study: OpenAI Dactyl Hand

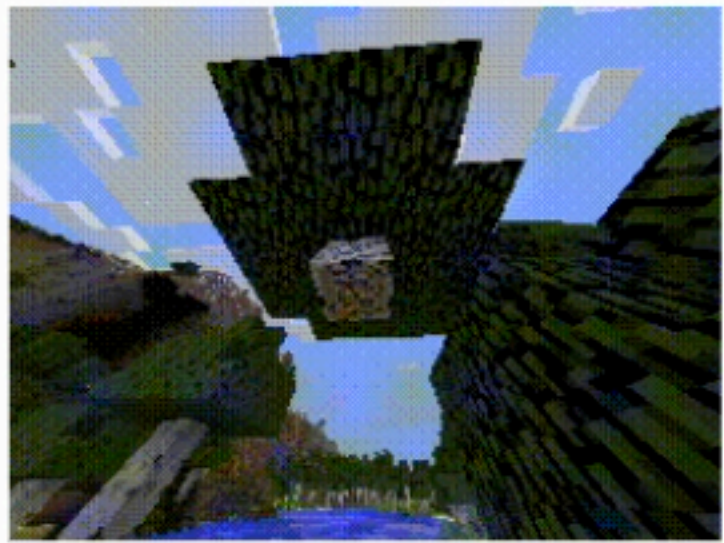
The DREAMER Algorithms

MineRL Diamond Challenge



MineRL Diamond Challenge

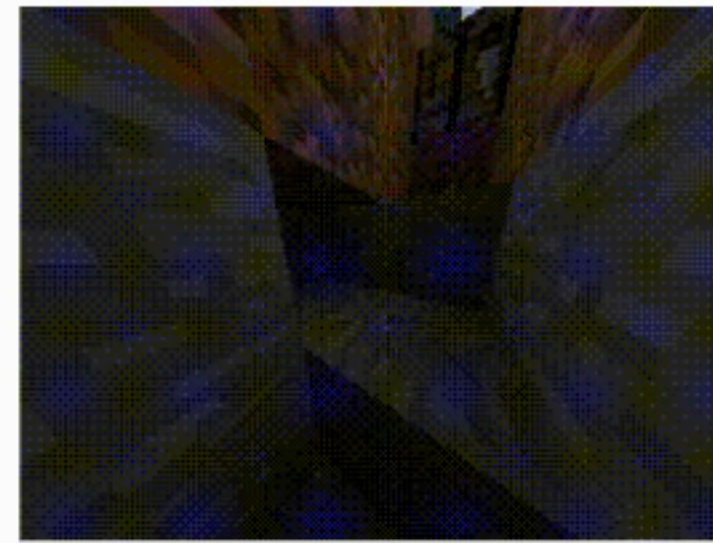
Gather Wood



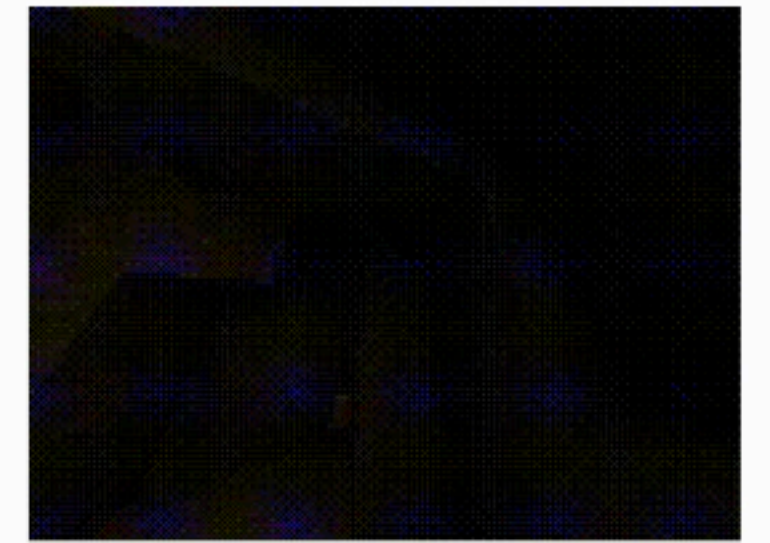
Create Wood Pickaxe



Mine Stone and Create Stone Pickaxe



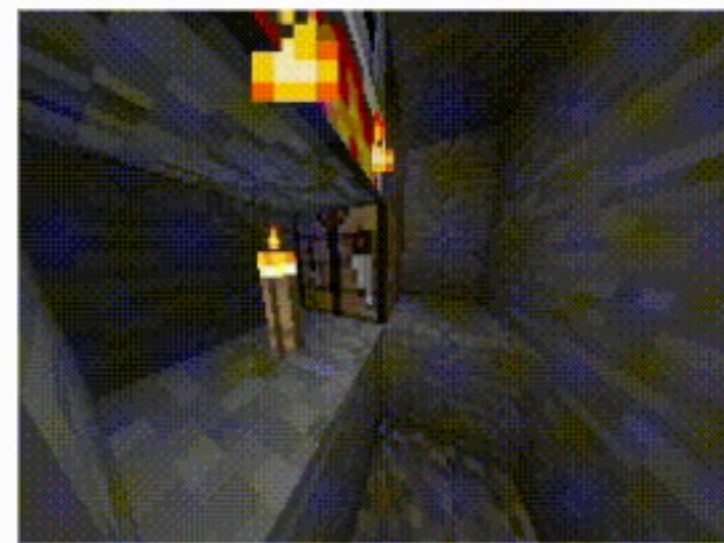
Mine Iron Ore



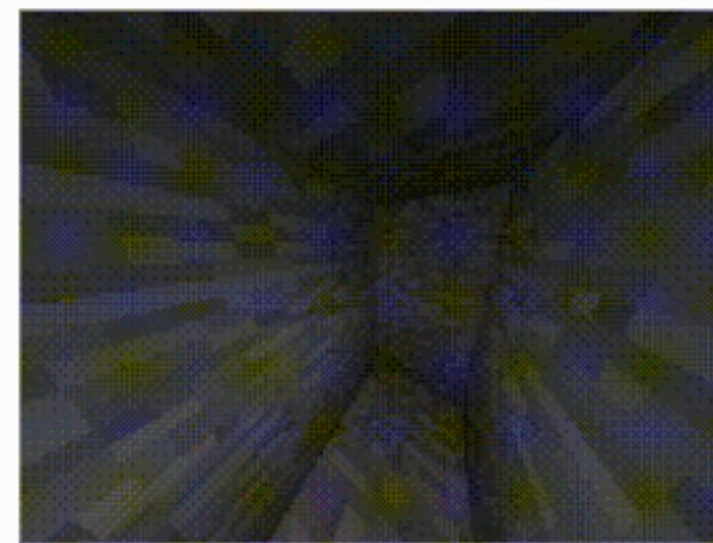
Create Furnace



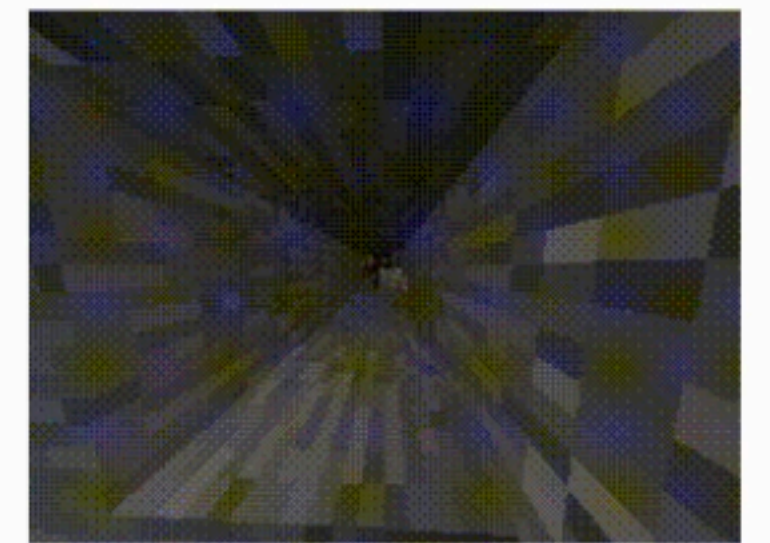
Smelt Iron and Create Iron Pickaxe



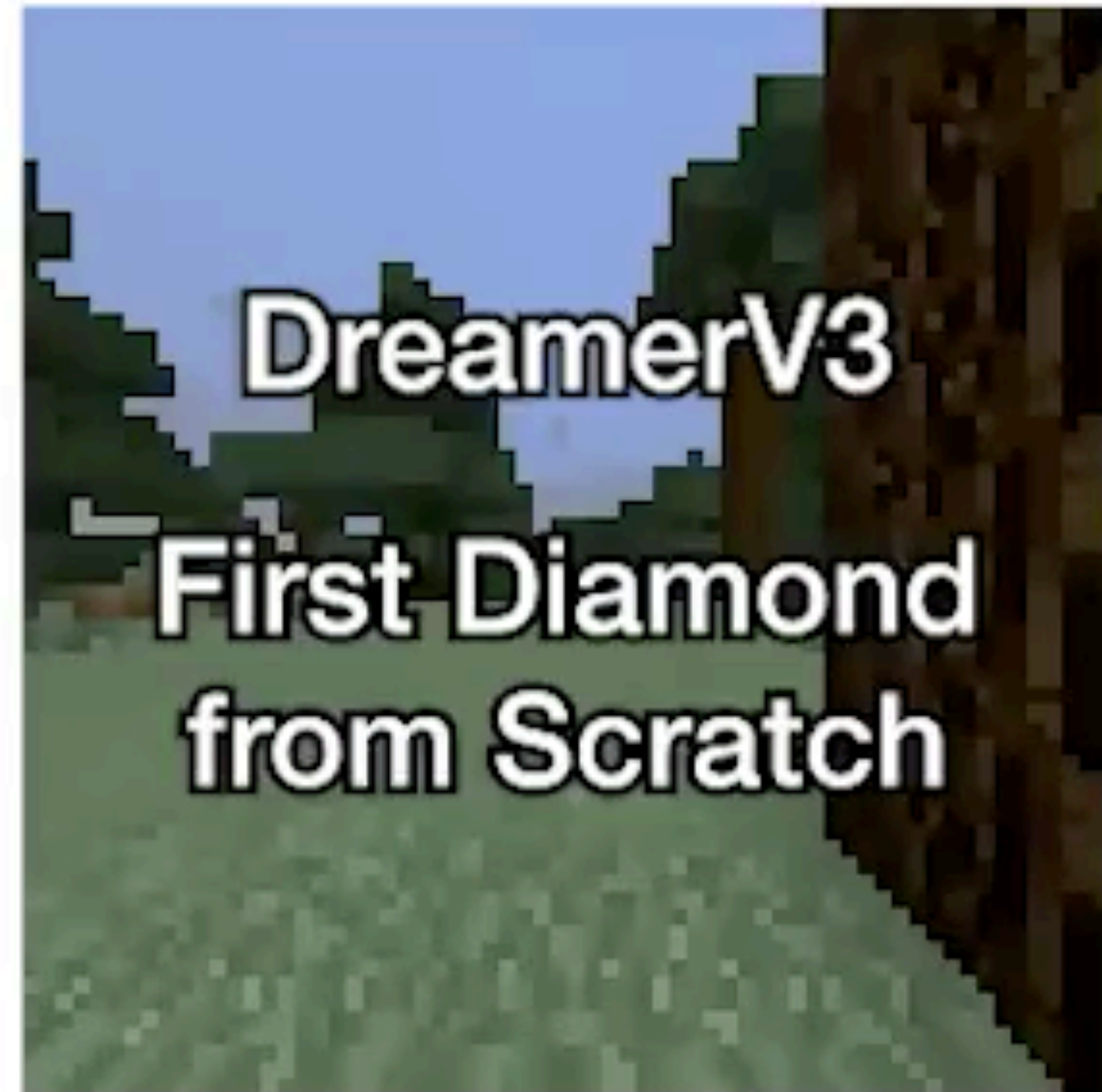
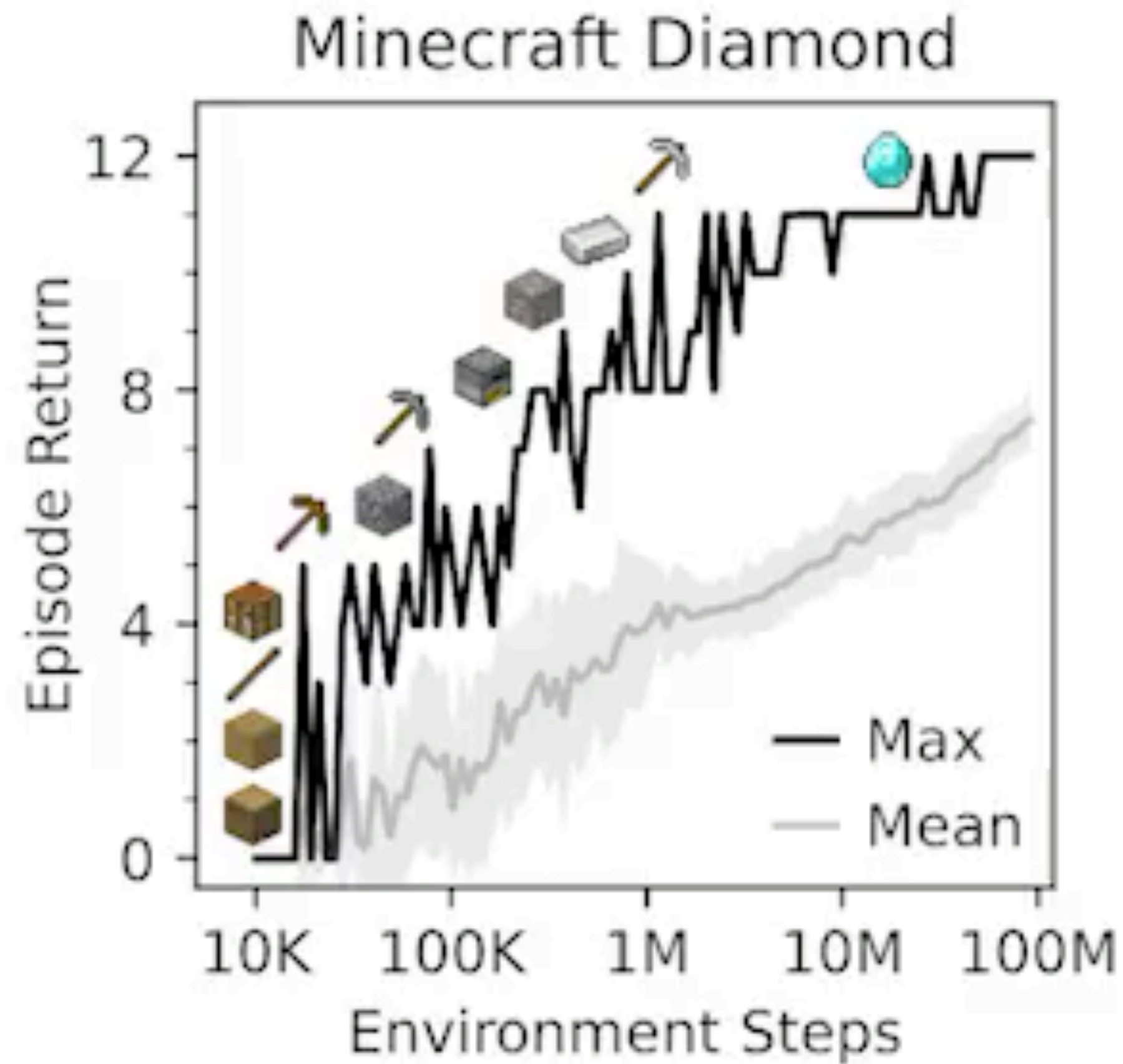
Search



Mine Diamond



DreamerV3 solved this task!





The
DREAMER
Algorithm



DREAM TO CONTROL: LEARNING BEHAVIORS BY LATENT IMAGINATION

Danijar Hafner *

University of Toronto

Google Brain

Timothy Lillicrap

DeepMind

Jimmy Ba

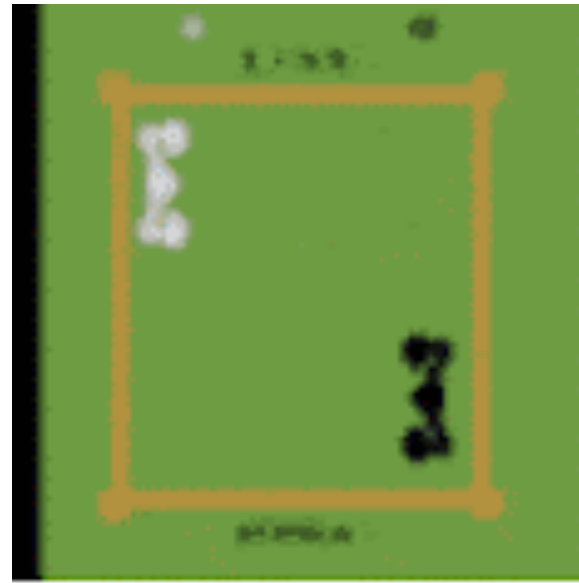
University of Toronto

Mohammad Norouzi

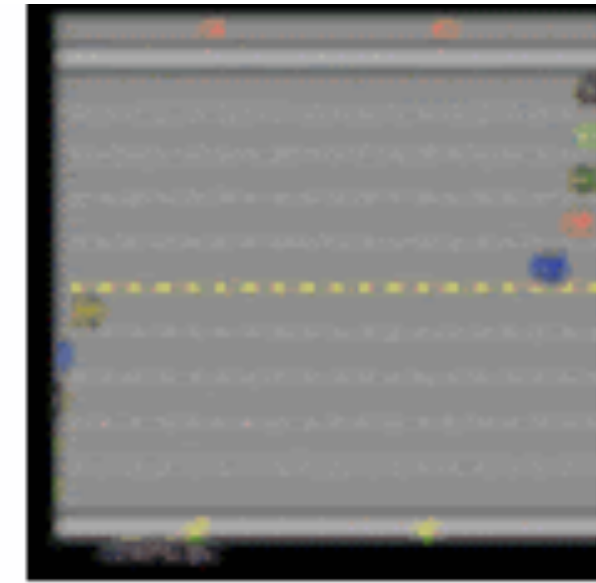
Google Brain

2020

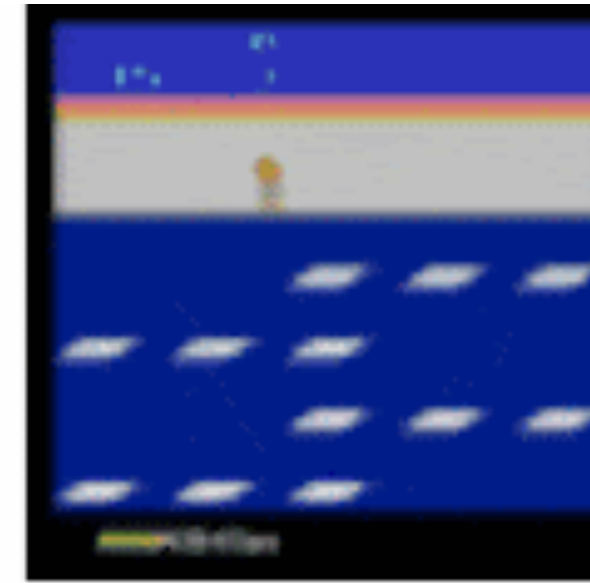
Look at the videos below



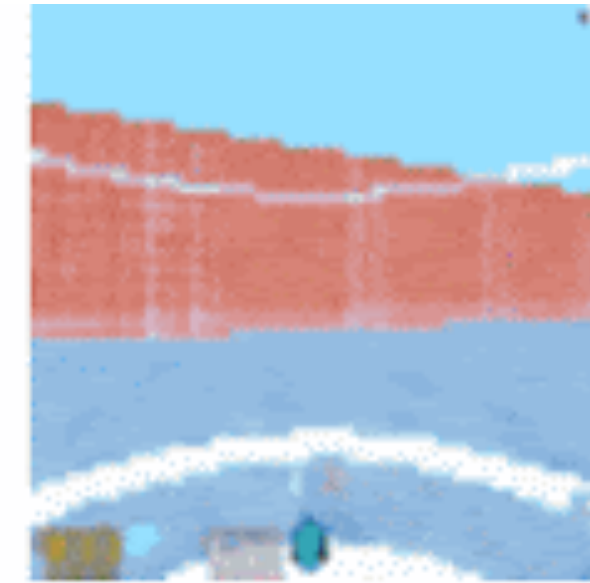
Boxing



Freeway



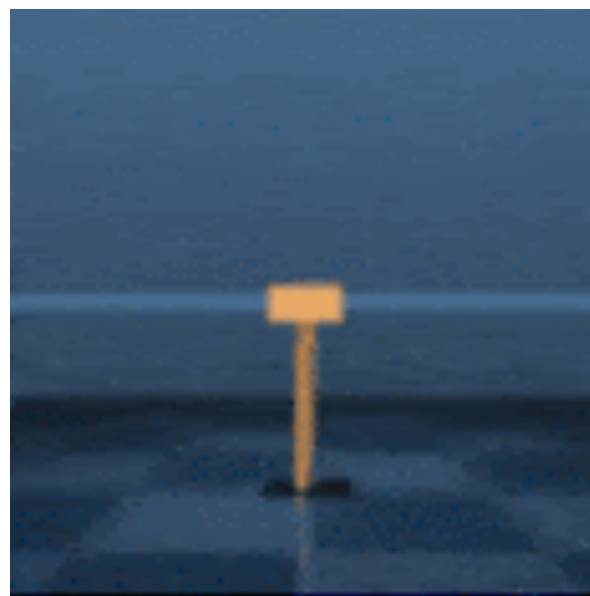
Frostbite



Collect Objects



Watermaze



Sparse Cartpole



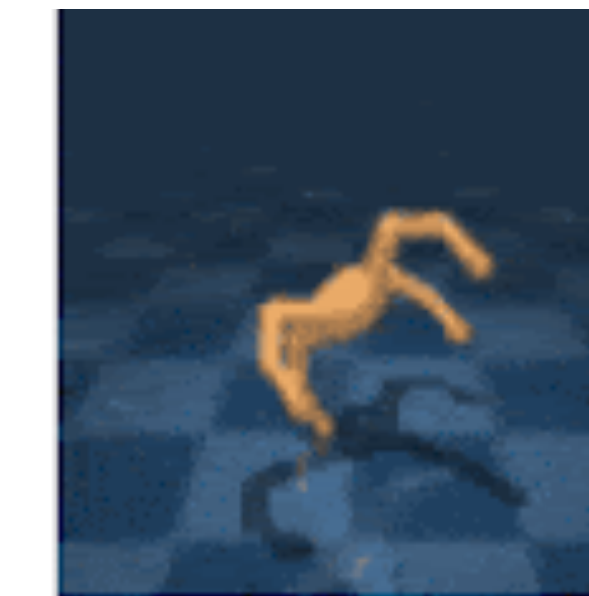
Acrobot Swingup



Hopper Hop



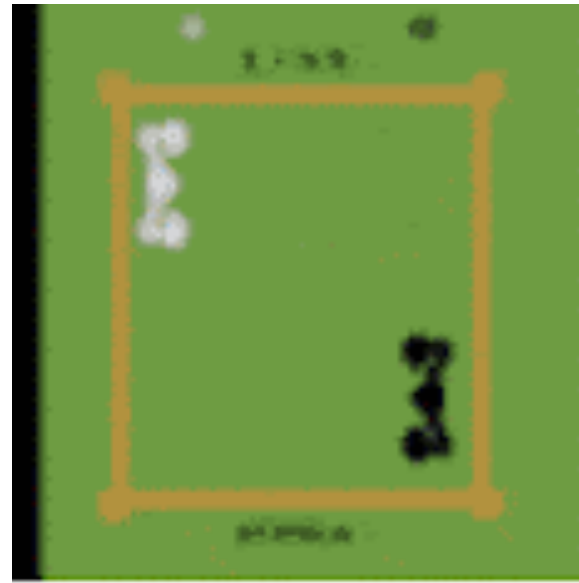
Walker Run



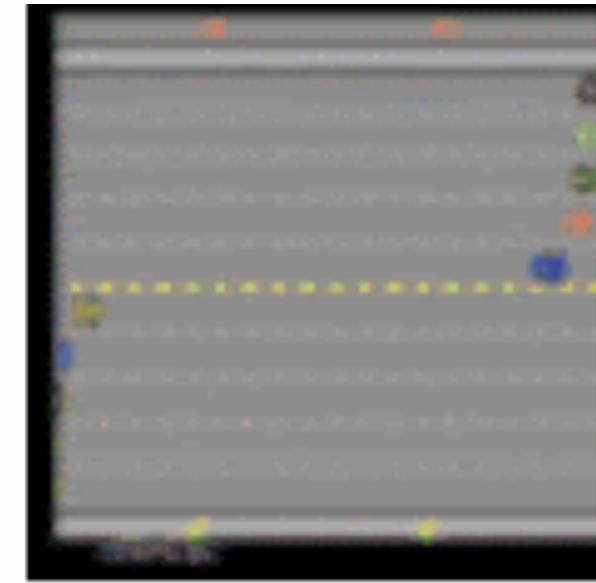
Quadruped Run

Which of these are real vs model?

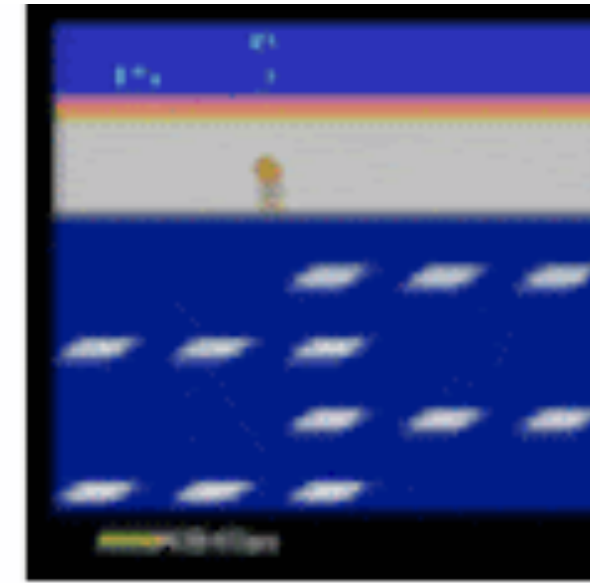
Look at the videos below



Boxing



Freeway



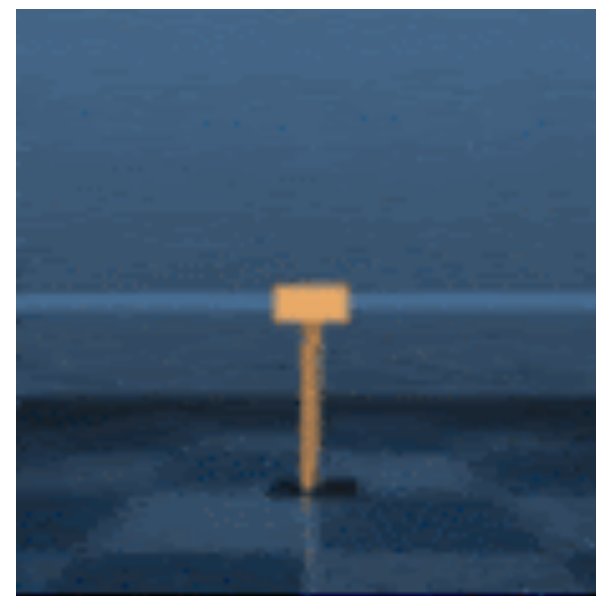
Frostbite



Collect Objects



Watermaze



Sparse Cartpole



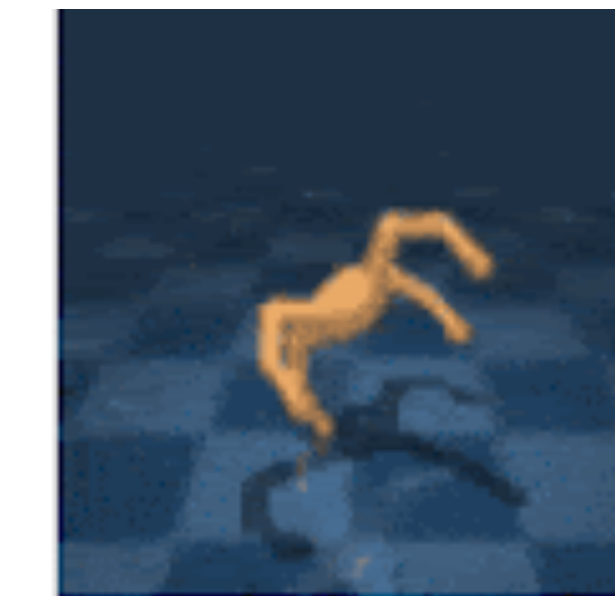
Acrobot Swingup



Hopper Hop



Walker Run

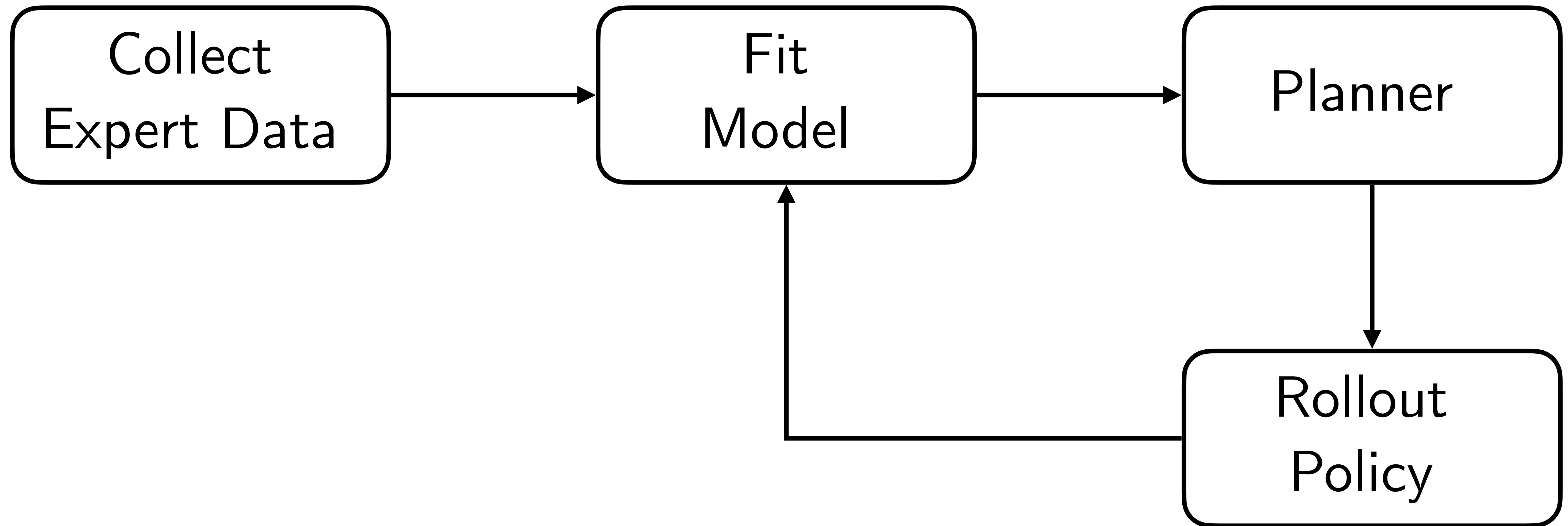


Quadruped Run

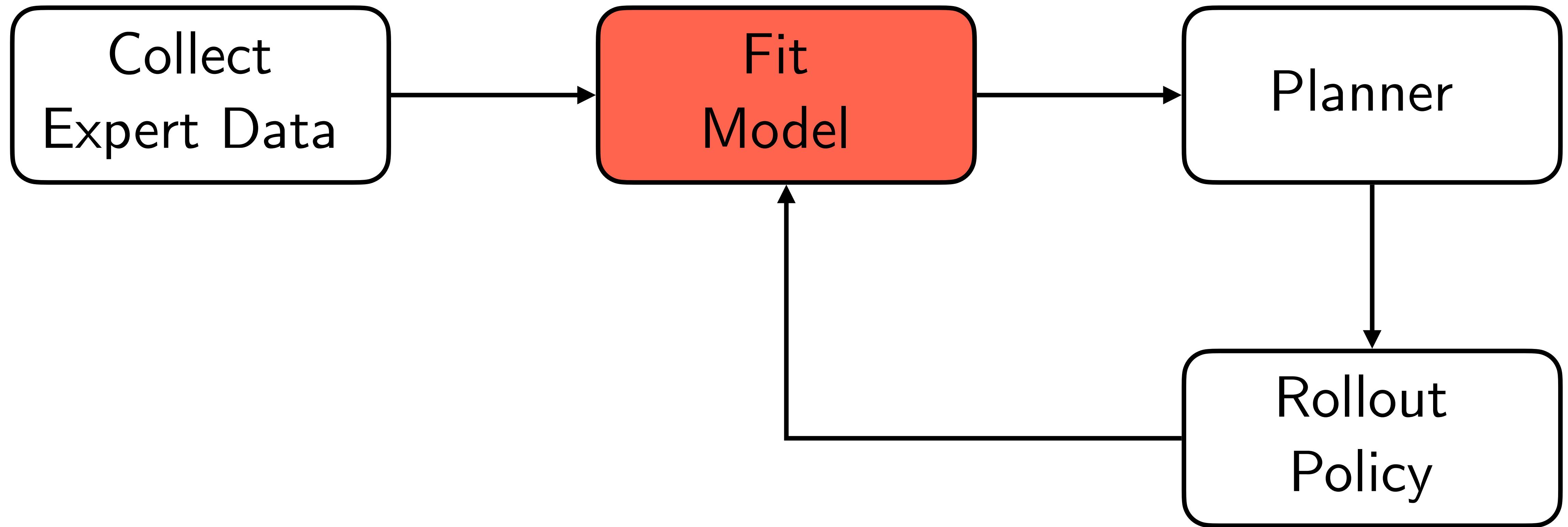
They are all from a model!

Recap: Model-based RL

(Ross & Bagnell, 2012)



How does DREAMER fit a model?



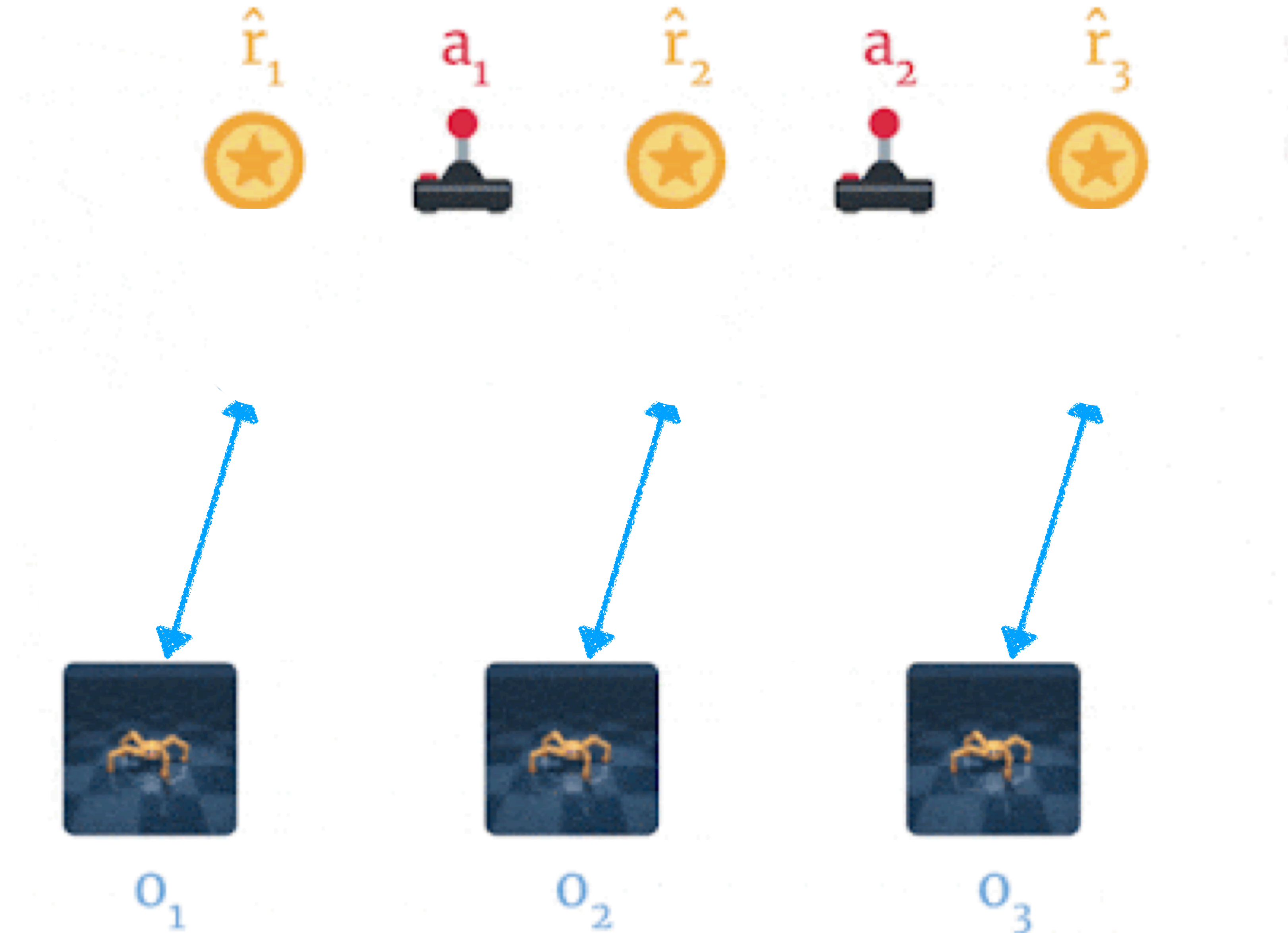
Goal: Fit a Model given data

Given Data:

Observations, rewards,
actions

Goal: Fit a Model given data

Given Data:
Observations, rewards,
actions



Predict:
States,
Dynamics Function,
Reward Function

Actions



Observations



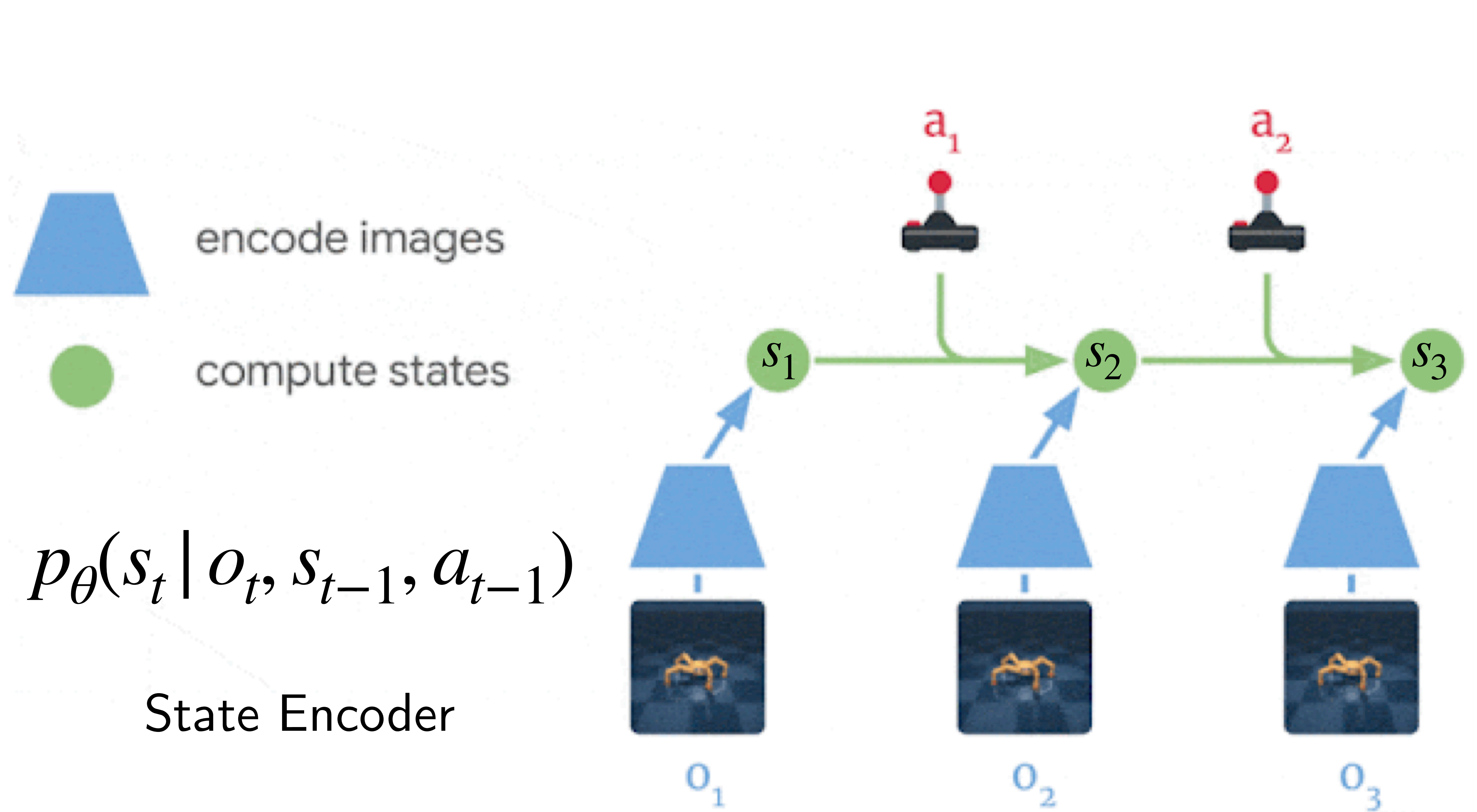
o_1



o_2



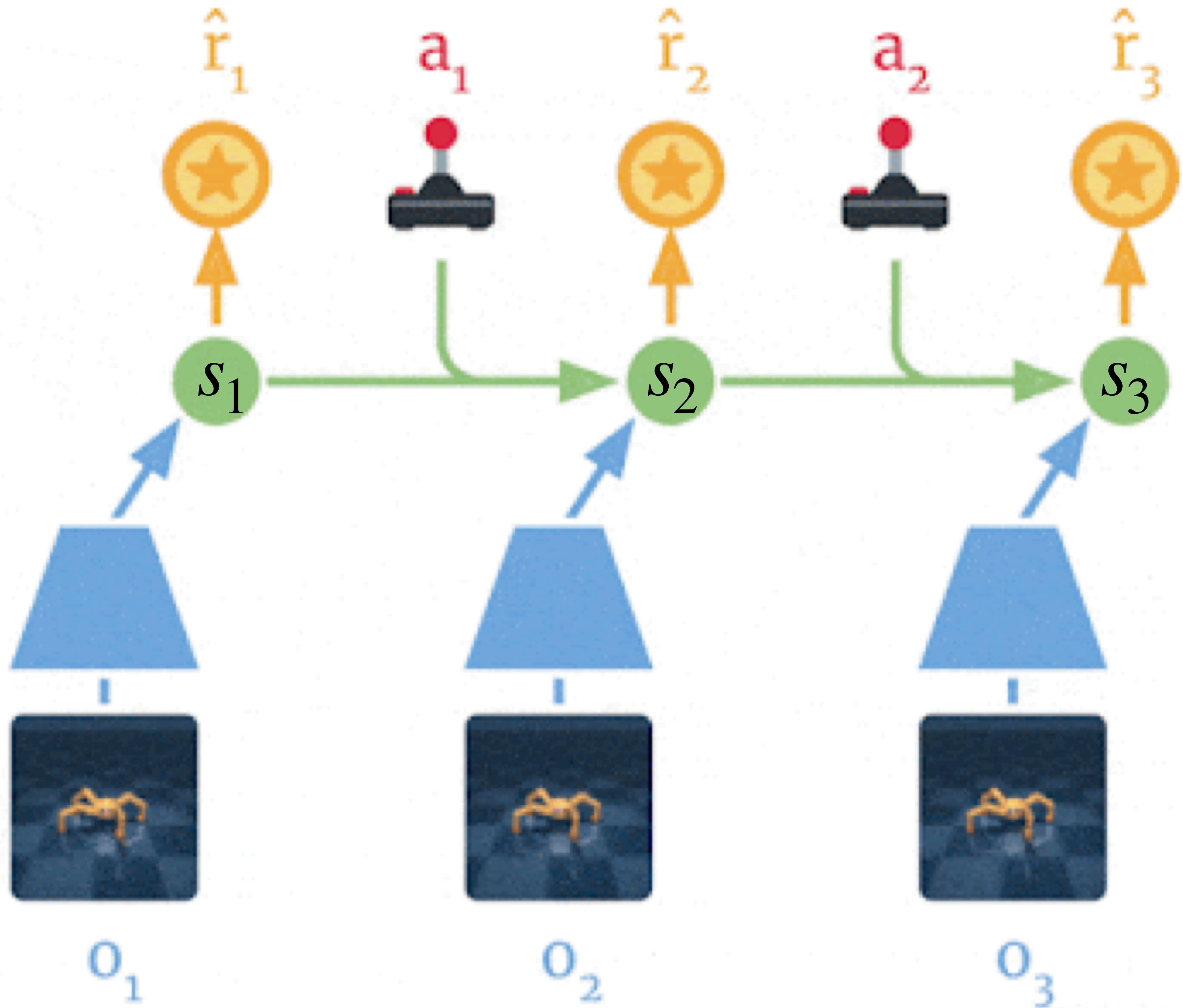
o_3



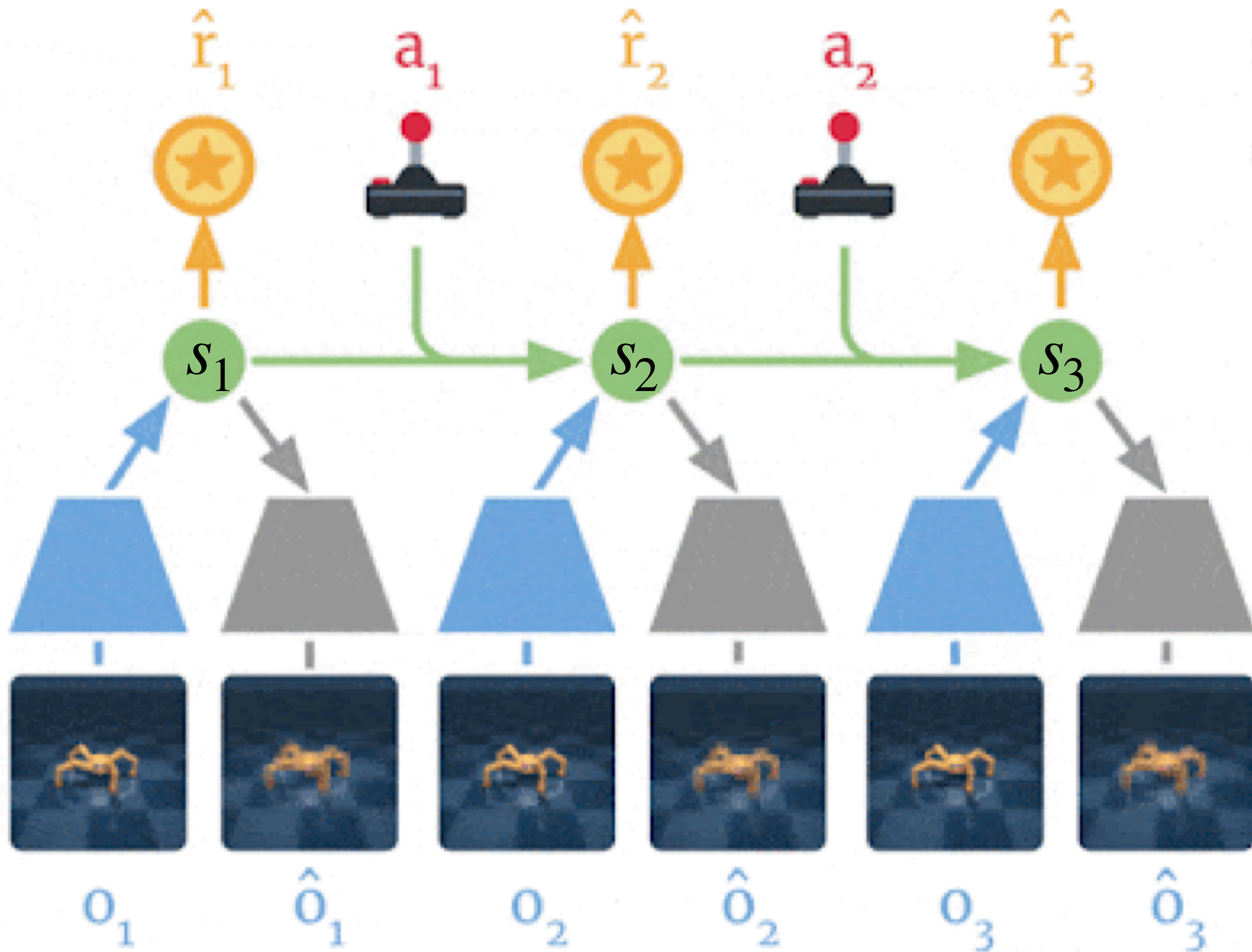
$$\mathcal{L} = (r_t - \hat{r}_t)^2$$

$$q_{\theta}(r_t | s_t)$$

Reward Decoder



$$\ell = (o_t - \hat{o}_t)^2$$

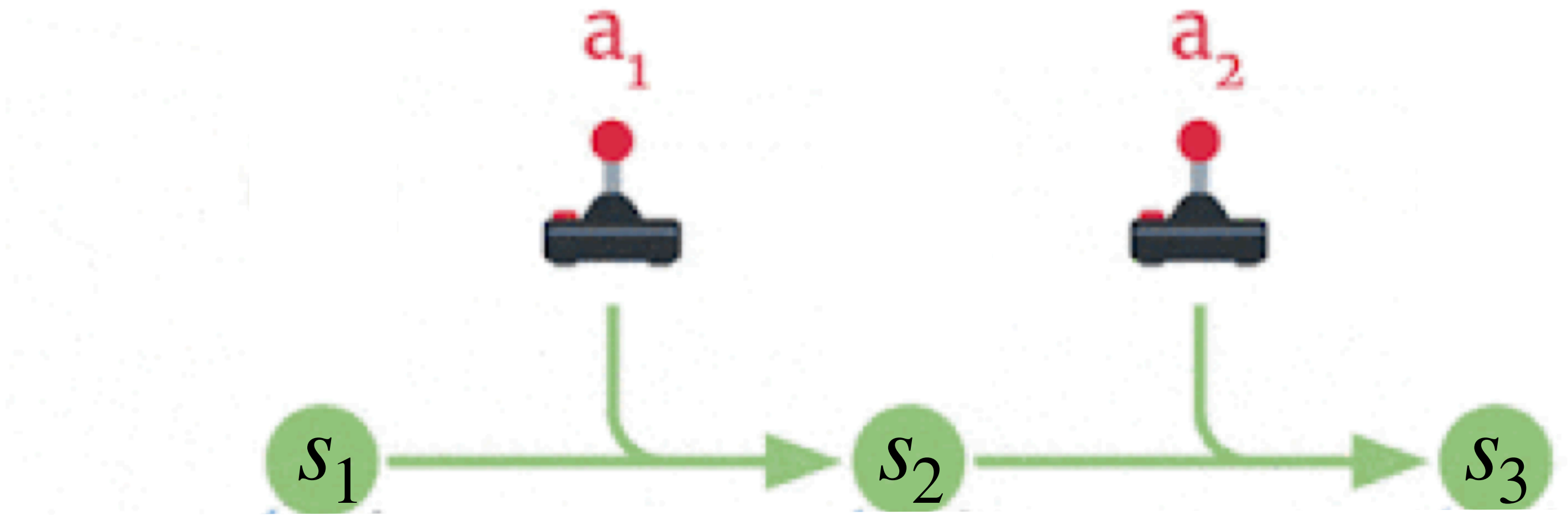


$$q_{\theta}(o_t | s_t)$$

Observation Decoder

$$q_{\theta}(s_{t+1} | s_t, a_t)$$

Dynamics
Function



$$\ell = KL(p_{\theta}(s_t | o_t, s_{t-1}, a_{t-1}) || q_{\theta}(s_{t+1} | s_t, a_t))$$

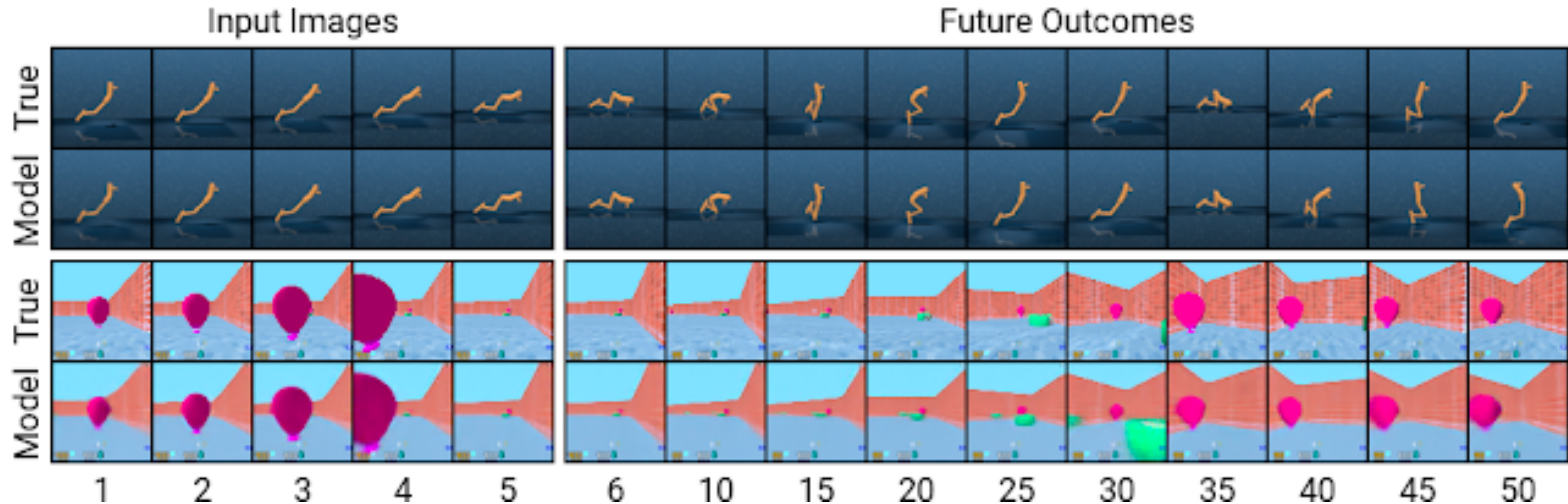


Freeze gradients

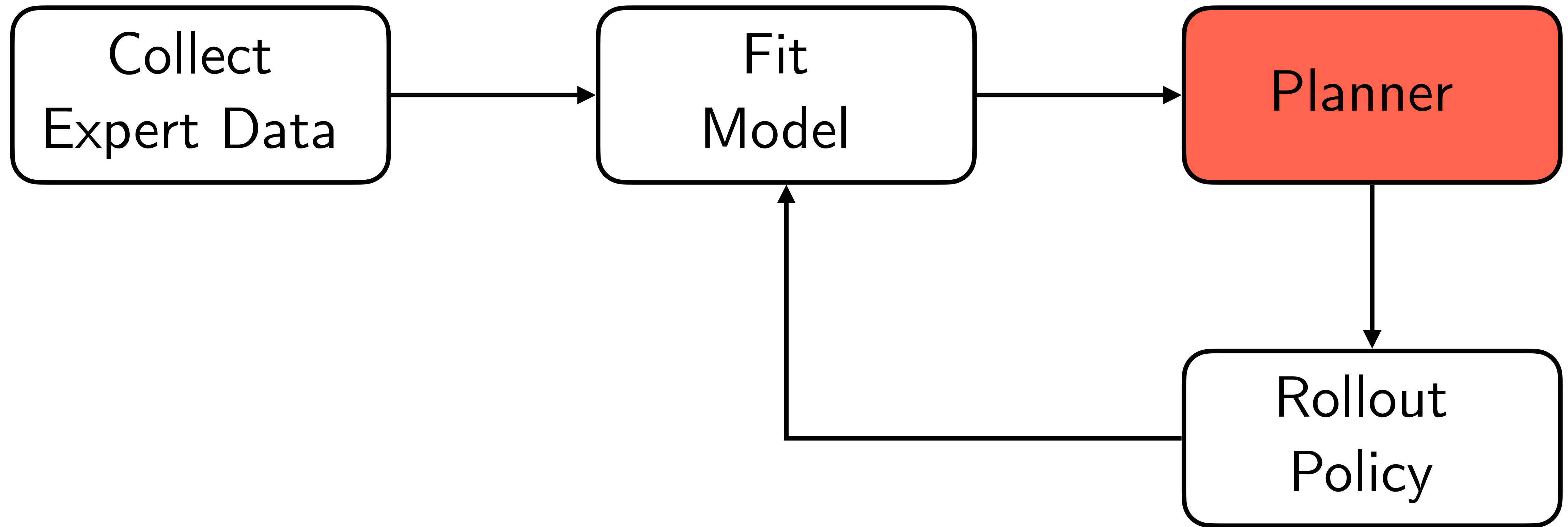
Results: Learning World Model



Results: Learning World Model



How does DREAMER do planning?



Goal: Learn a Policy using Actor-Critic

$$\pi_{\phi}(a_t | s_t)$$

Actor

$$V_{\psi}(s_t)$$

Critic

From rollouts in the model

$$q_{\theta}(s_t | s_{t-1}, a_{t-1})$$

Recall: Actor-Critic for model-free RL

Start with an arbitrary initial policy $\pi_\phi(a | s)$

while *not converged* **do**

Roll-out $\pi_\phi(a | s)$ to collect trajectories $D = \{s^i, a^i, r^i, s_+^i\}_{i=1}^N$

Fit value function $V_\psi(s^i)$ using TD, i.e. minimize $(r^i + \gamma V_\psi(s_+^i) - V_\psi(s^i))^2$

Compute advantage $\hat{A}(s^i, a^i) = r(s^i, a^i) + \gamma V_\psi(s_+^i) - V_\psi(s^i)$

Compute gradient

$$\nabla_\phi J(\phi) = \frac{1}{N} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\phi(a_t^i | s_t^i) \hat{A}(s^i, a^i) \right]$$

Update parameters

$$\phi \leftarrow \phi + \alpha \nabla_\phi J(\phi)$$

Actor-Critic in Model-based RL

Start with an arbitrary initial policy $\pi_\phi(a | s)$

while *not converged* **do**

Roll-out $\pi_\phi(a | s)$ in the model $q_\theta(s' | s, a)$ to collect trajectories $D = \{s^i, a^i, r^i, s_+^i\}_{i=1}^N$

Fit value function $V_\psi(s^i)$ using TD, i.e. minimize $(r^i + \gamma V_\psi(s_+^i) - V_\psi(s^i))^2$

Compute advantage $\hat{A}(s^i, a^i) = r(s^i, a^i) + \gamma V_\psi(s_+^i) - V_\psi(s^i)$

Compute gradient

$$\nabla_\phi J(\phi) = \frac{1}{N} \left[\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\phi(a_t^i | s_t^i) \hat{A}(s^i, a^i) \right]$$

Update parameters

$$\phi \leftarrow \phi + \alpha \nabla_\phi J(\phi)$$

Actor-Critic in Model-based RL

Start with an arbitrary initial policy $\pi_\phi(a | s)$

while *not converged* **do**

Roll-out $\pi_\phi(a | s)$ in the model $q_\theta(s' | s, a)$ to collect trajectories $D = \{s^i, a^i, r^i, s_+^i\}_{i=1}^N$

Fit value function $V_\psi(s^i)$ using TD, i.e. minimize $(r^i + \gamma V_\psi(s_+^i) - V_\psi(s^i))^2$

Directly backprop gradients through model to update policy!

$$\phi \leftarrow \phi + \alpha \nabla_\phi V_\psi(s^i)$$

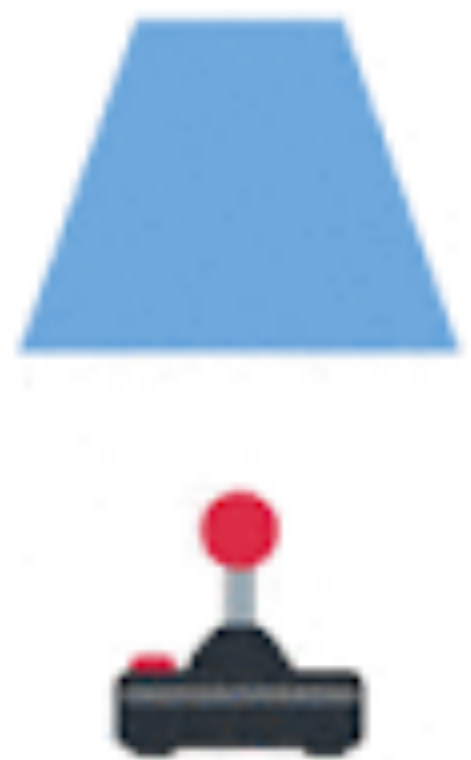


O_1



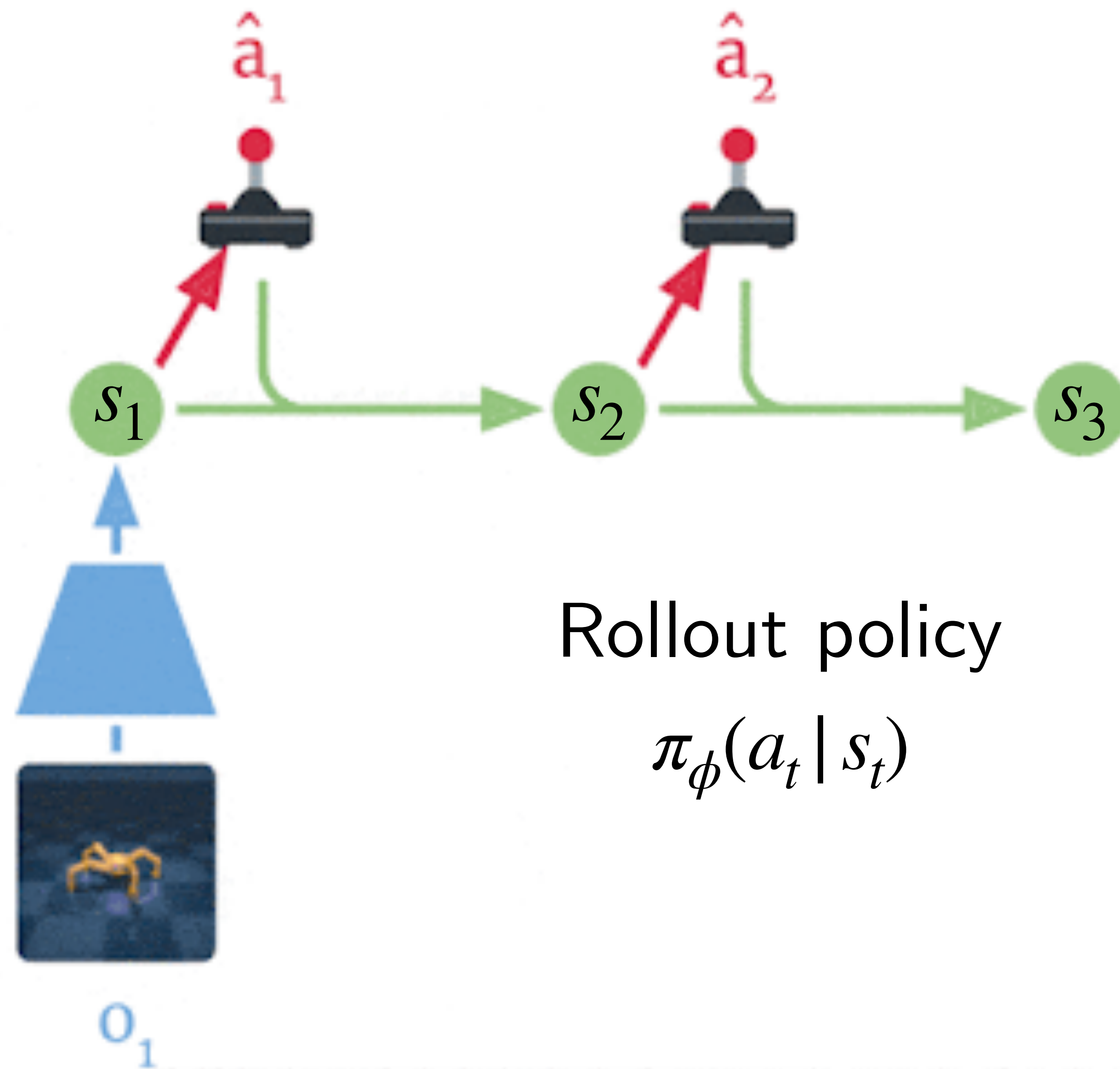
encode images





encode images

imagine ahead





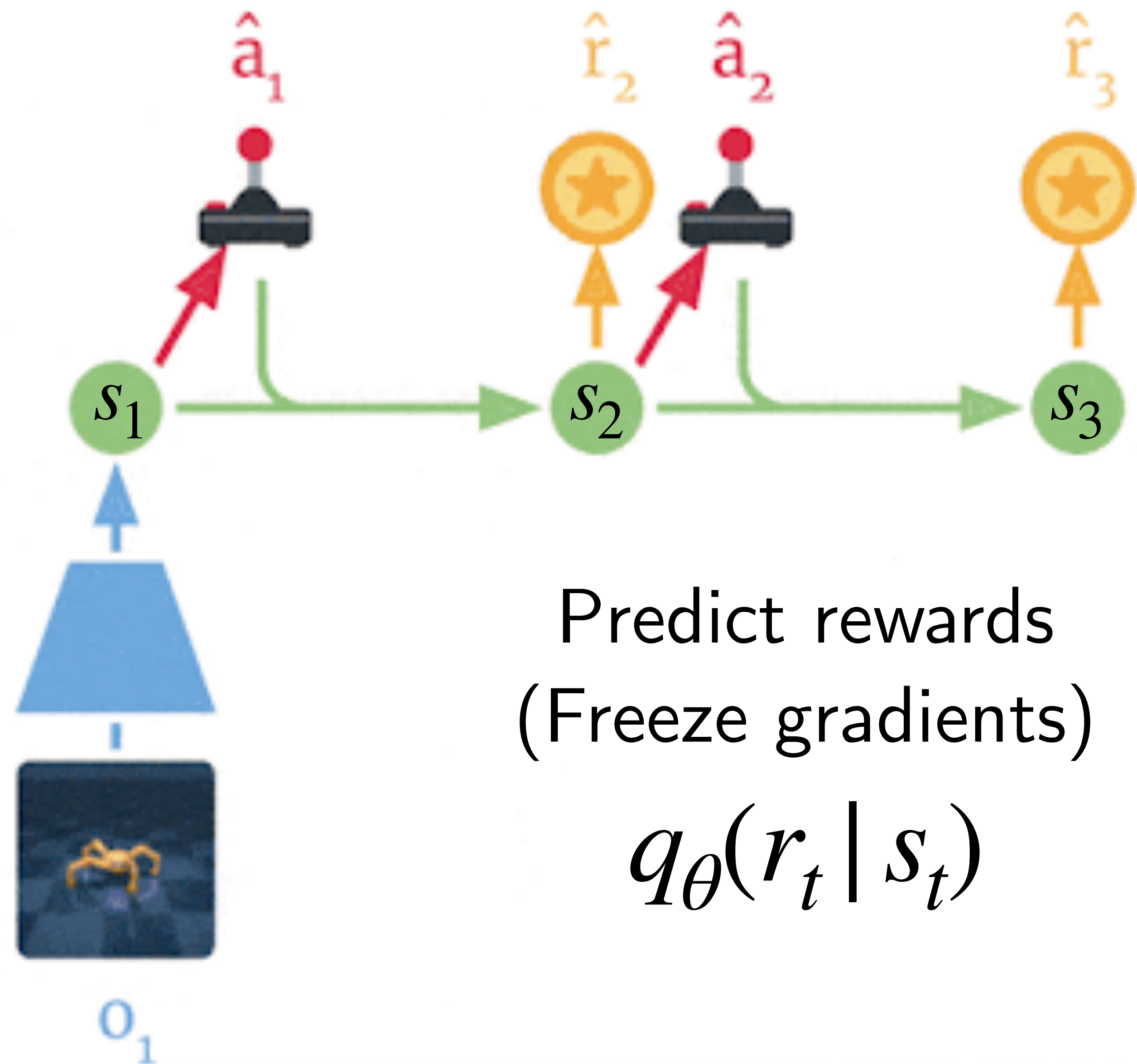
encode images



imagine ahead



predict rewards





encode images



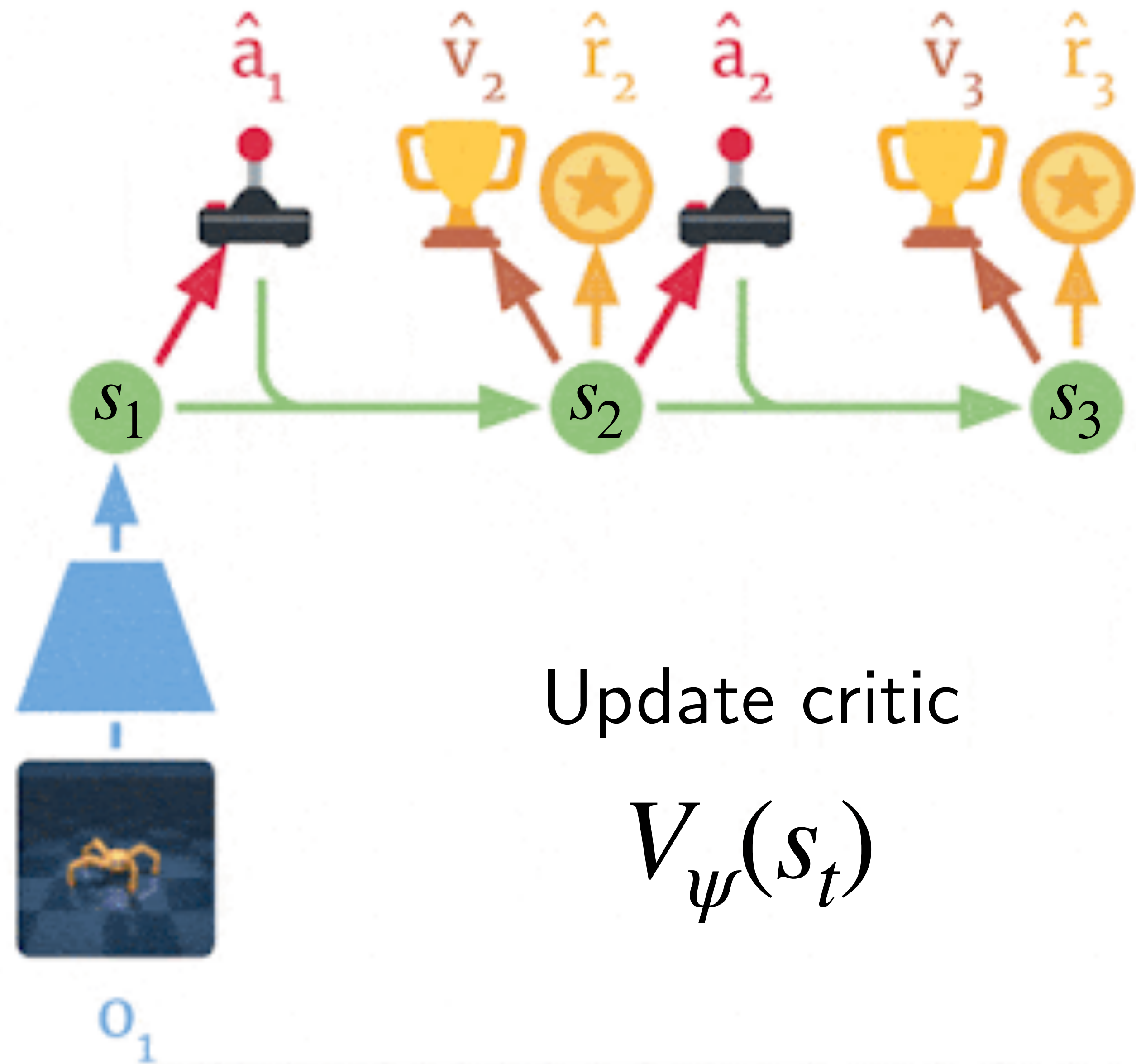
imagine ahead



predict rewards



predict values





encode images



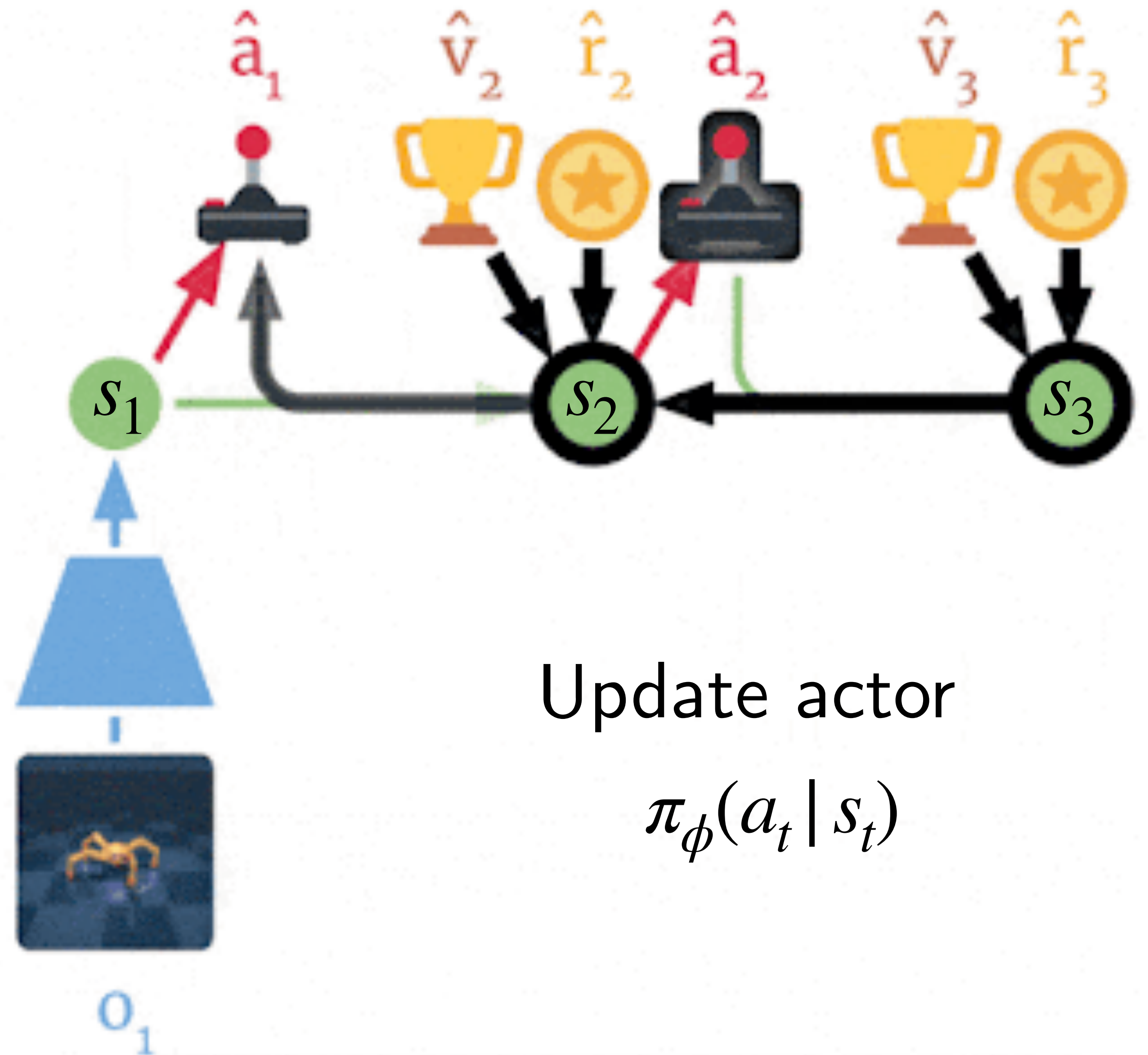
imagine ahead



predict rewards



predict values

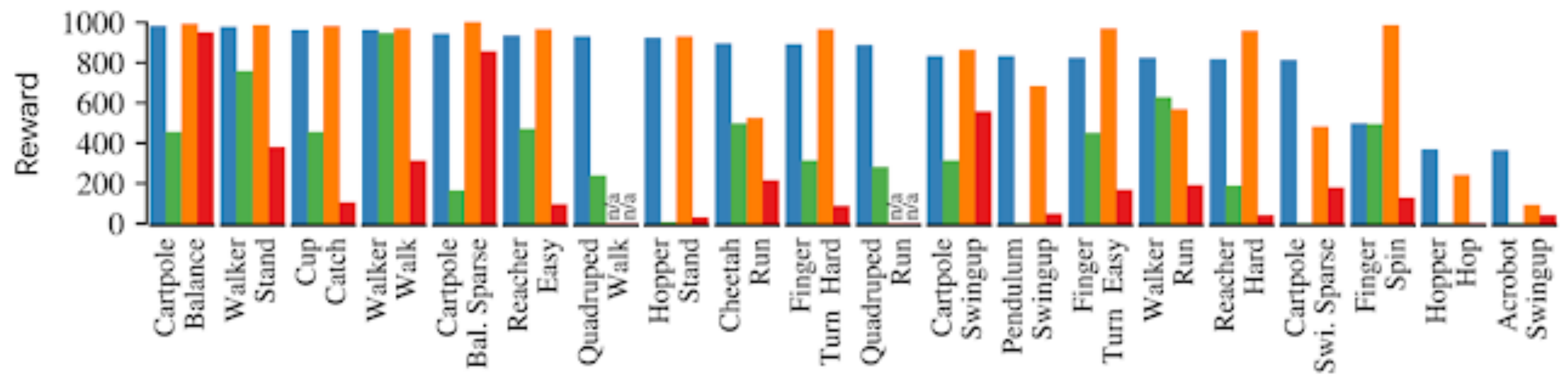


DREAMER: Results



Model-based { Dreamer (823) PlaNet (332)
 28 hours of interaction

Model-free { D4PG (786) A3C (243)
 23 days of interaction



DREAMER is a template
for Model-based RL

But there are many challenges as we
scale to harder real-world applications

DREAMER V2, V3, etc

Today's class

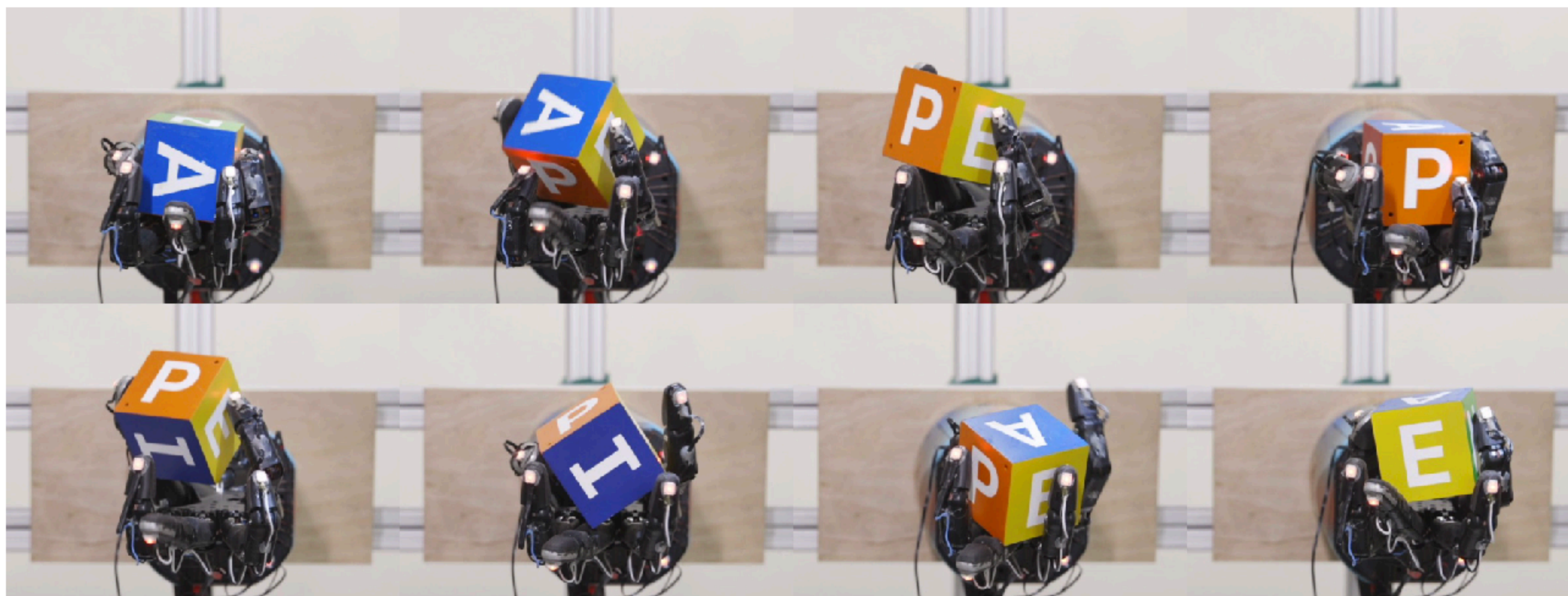
- ☑ Practical MBRL
(Only observations, complex dynamics)
- ☑ Learning Models: The DREAMER algorithm
- ☐ Leveraging Physics: Sim2Real
Case study: OpenAI Dactyl Hand

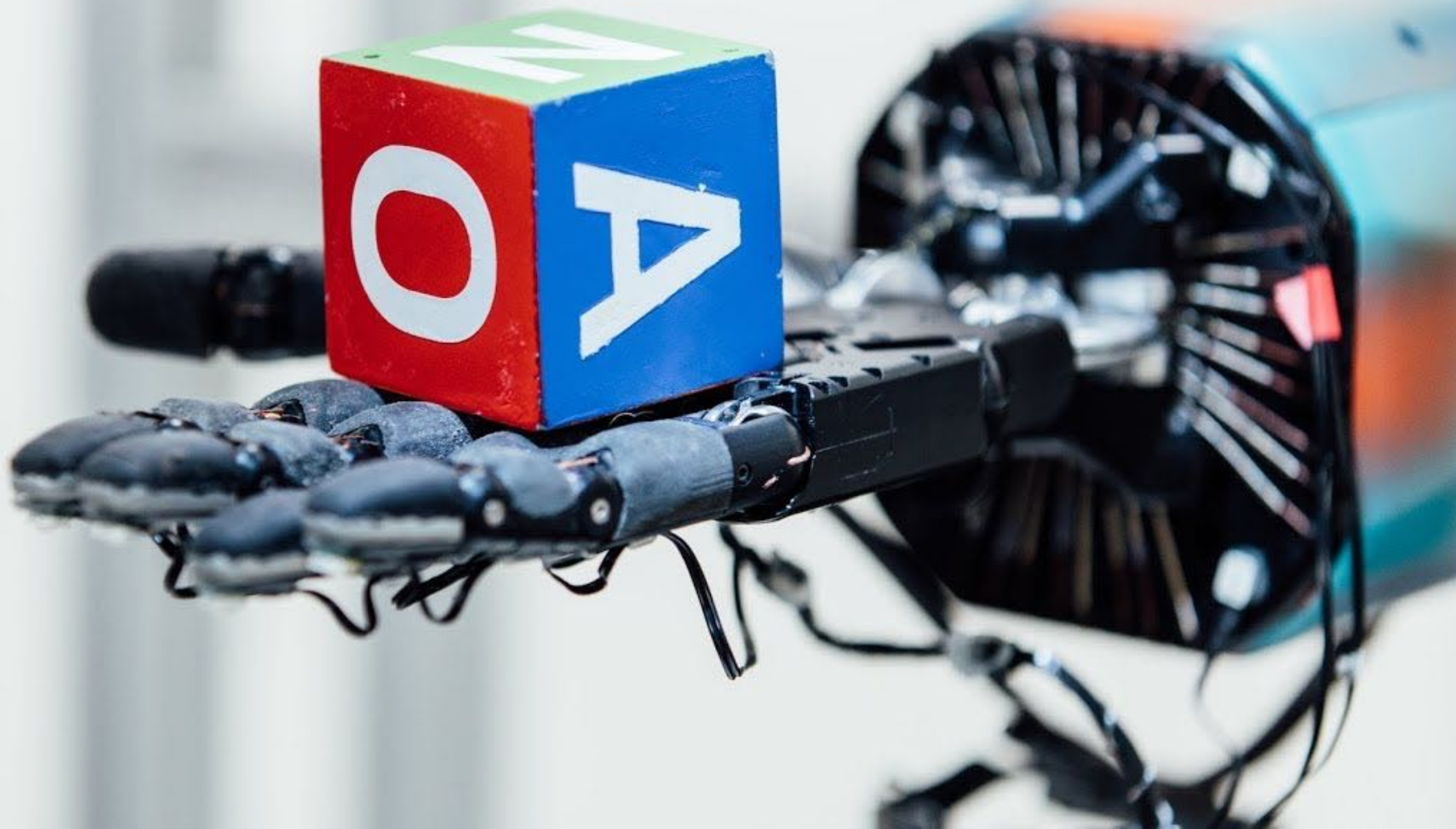
Learning Dexterity

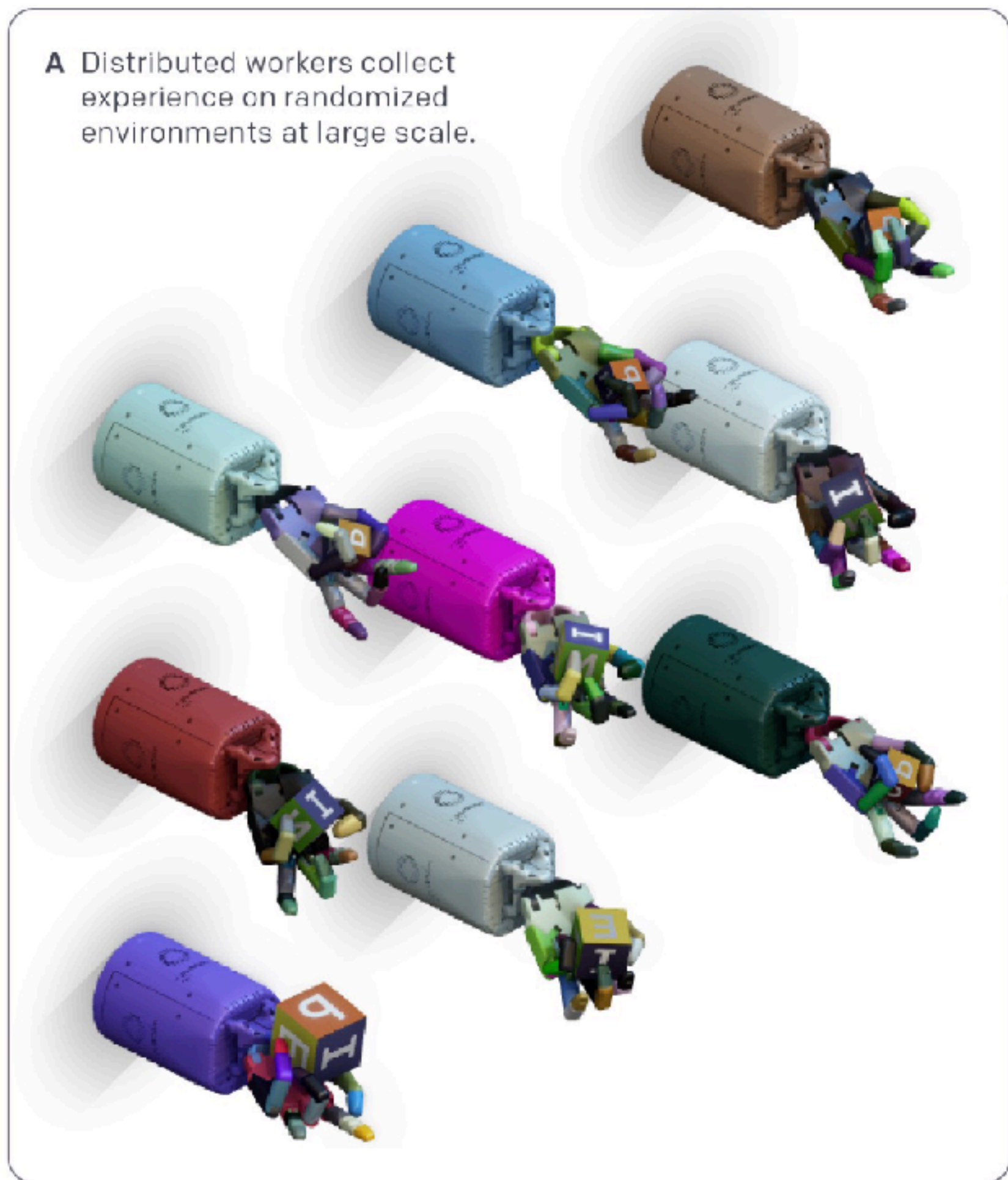
(Open AI)

Learning Dexterous In-Hand Manipulation

OpenAI*, Marcin Andrychowicz, Bowen Baker, Maciek Chociej,
Rafał Józefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron,
Matthias Plappert, Glenn Powell, Alex Ray, Jonas Schneider, Szymon Sidor,
Josh Tobin, Peter Welinder, Lilian Weng, Wojciech Zaremba







Sim

Train a policy in simulation
(RL)



Real

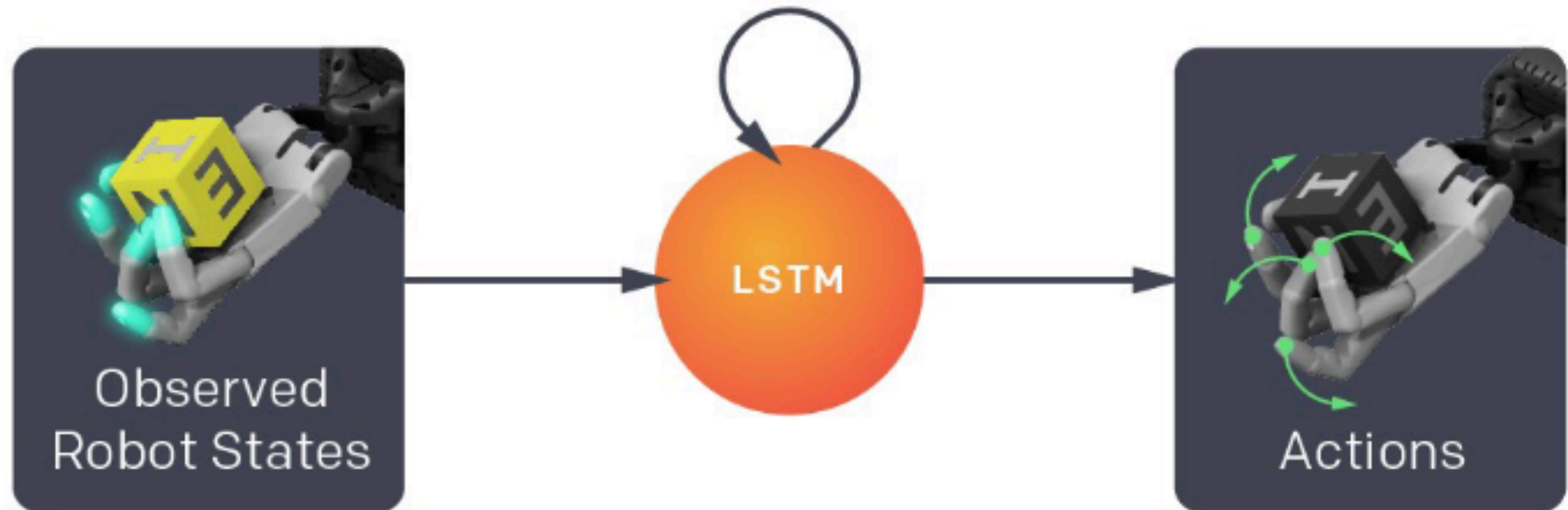
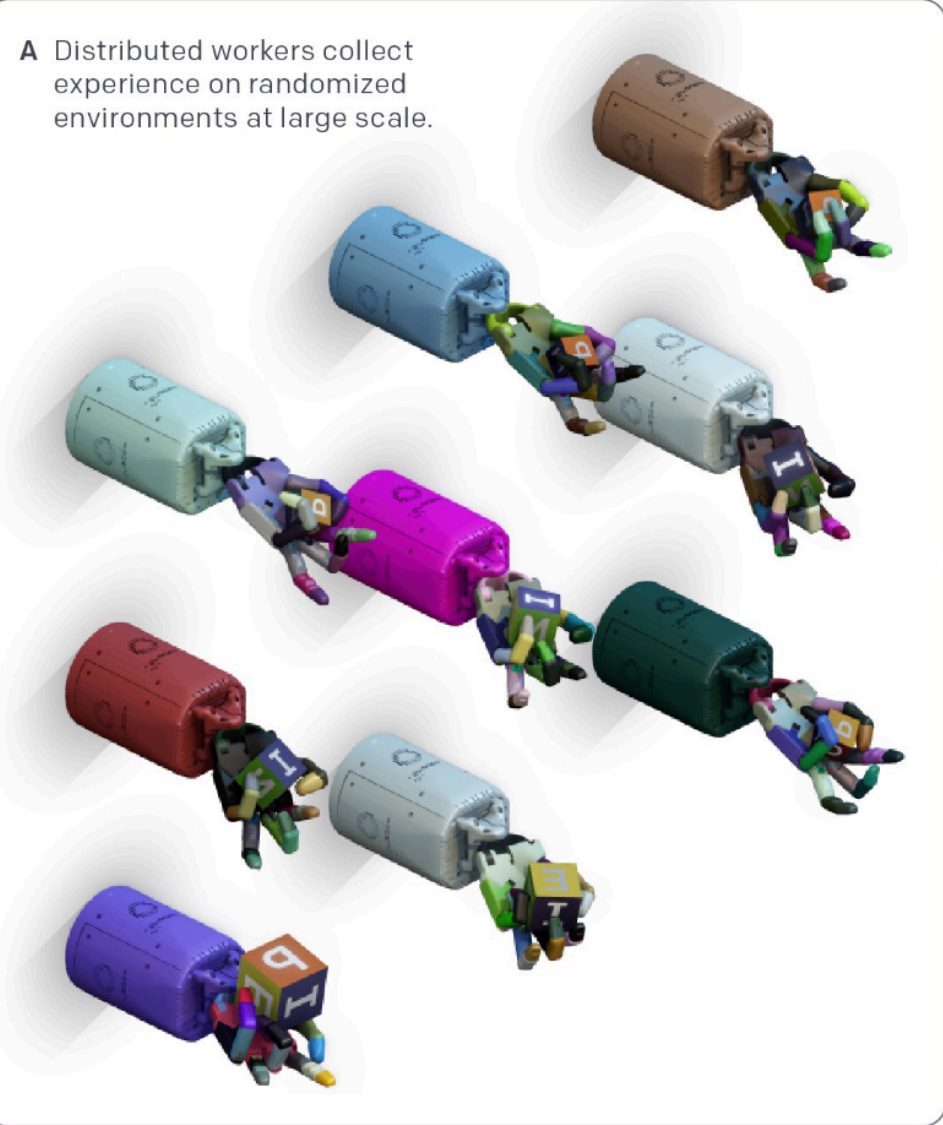
Test in real world

A Distributed workers collect experience on randomized environments at large scale.



B We train a control policy using reinforcement learning. It chooses the next action based on fingertip positions and the object pose.

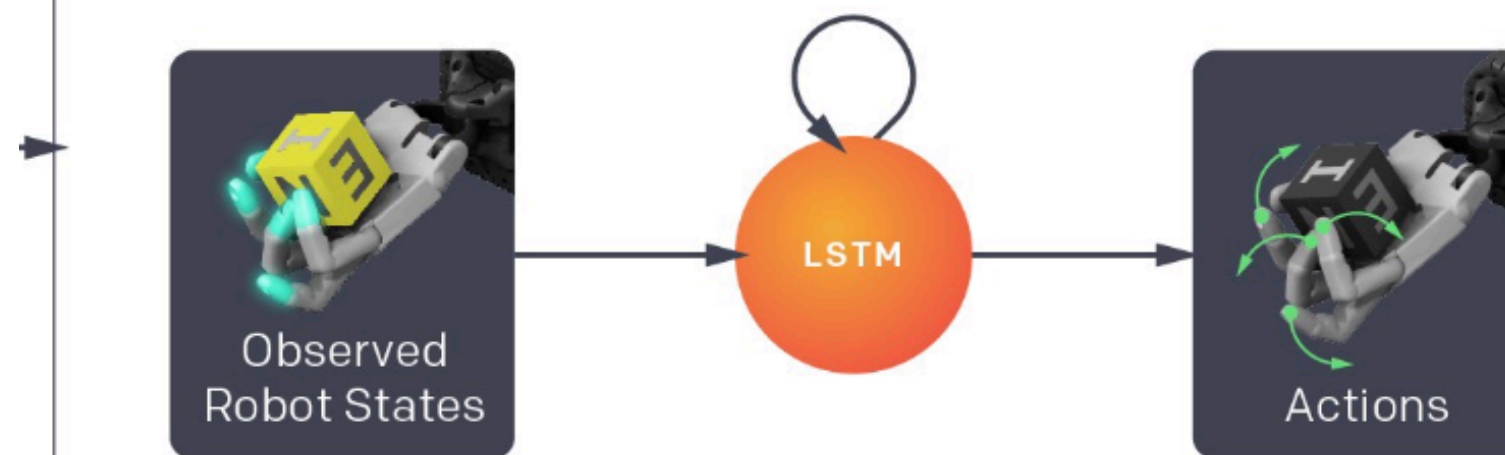
A Distributed workers collect experience on randomized environments at large scale.



A Distributed workers collect experience on randomized environments at large scale.

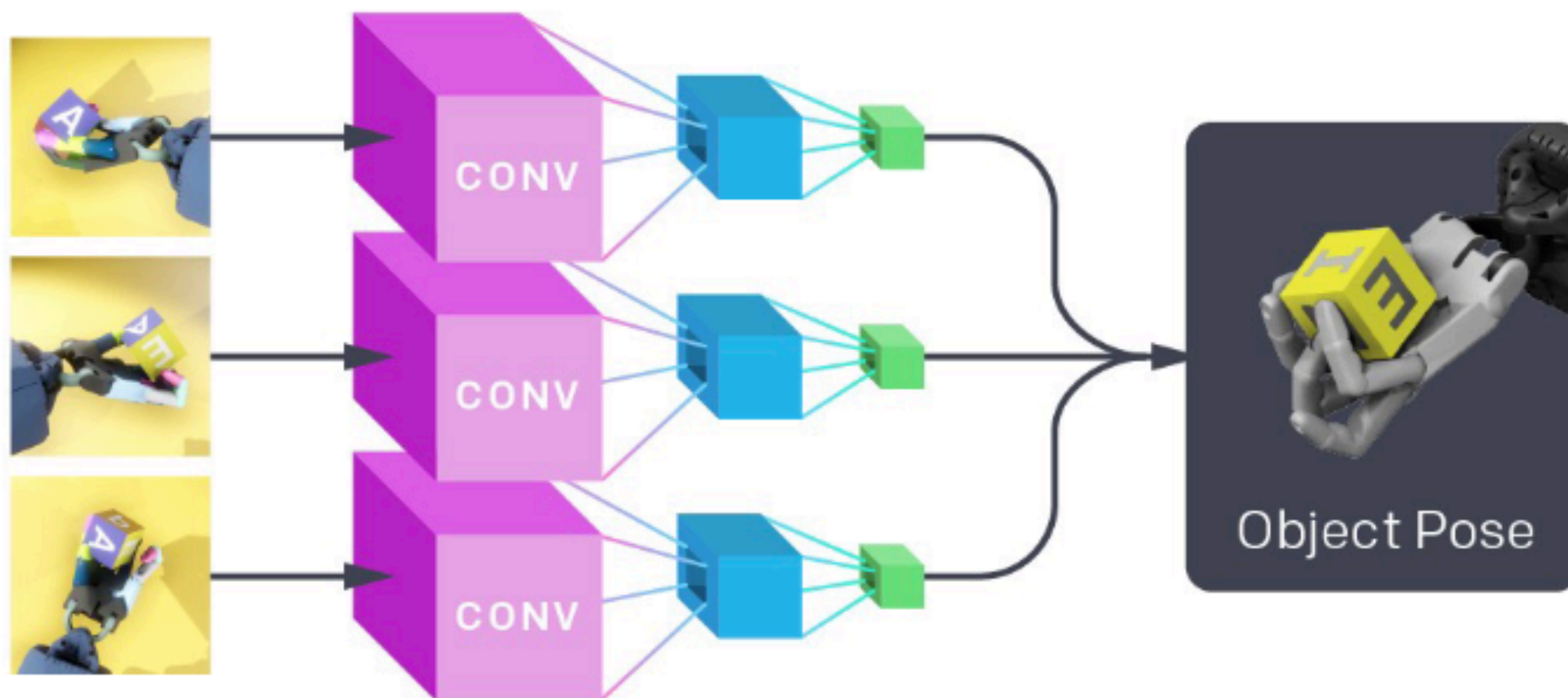


B We train a control policy using reinforcement learning. It chooses the next action based on fingertip positions and the object pose.

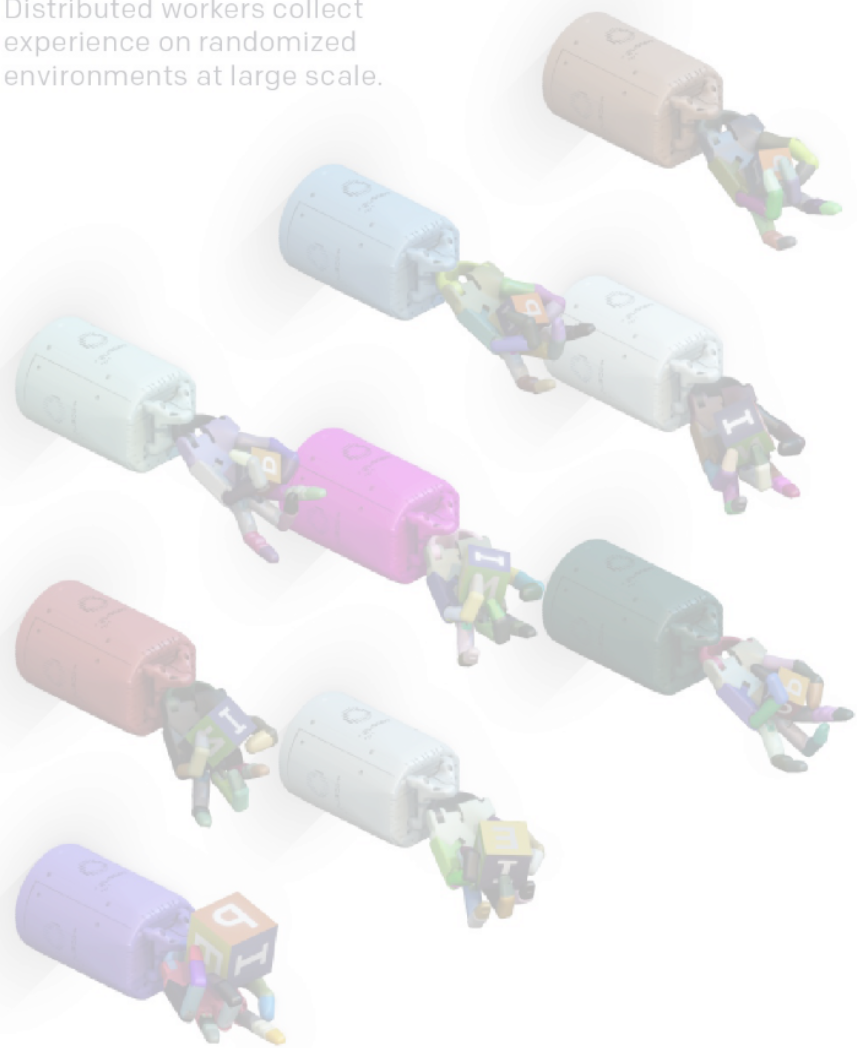


C We train a convolutional neural network to predict the object pose given three simulated camera images.

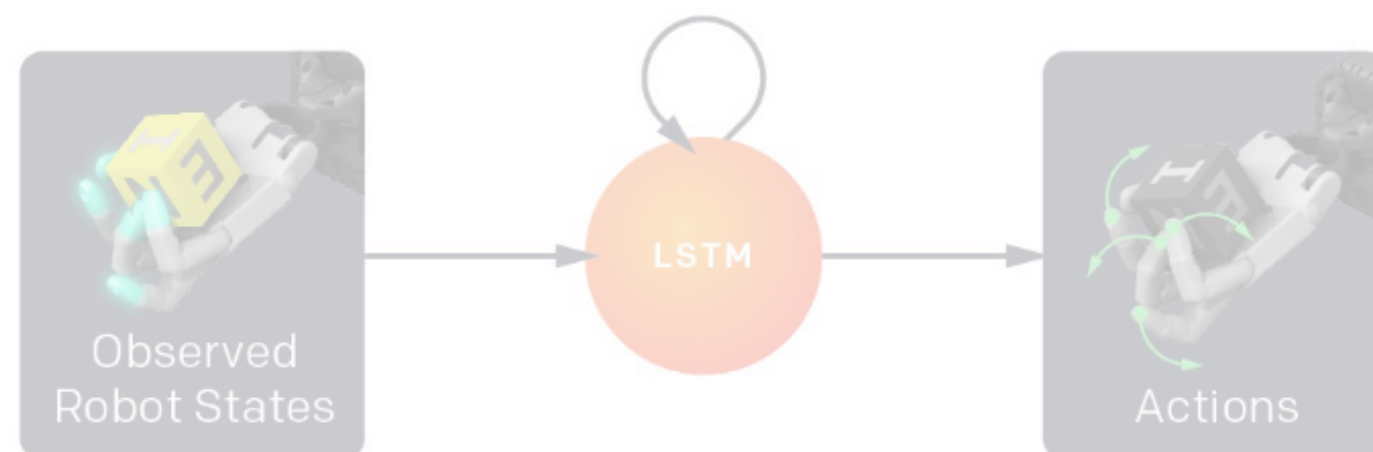
Sim



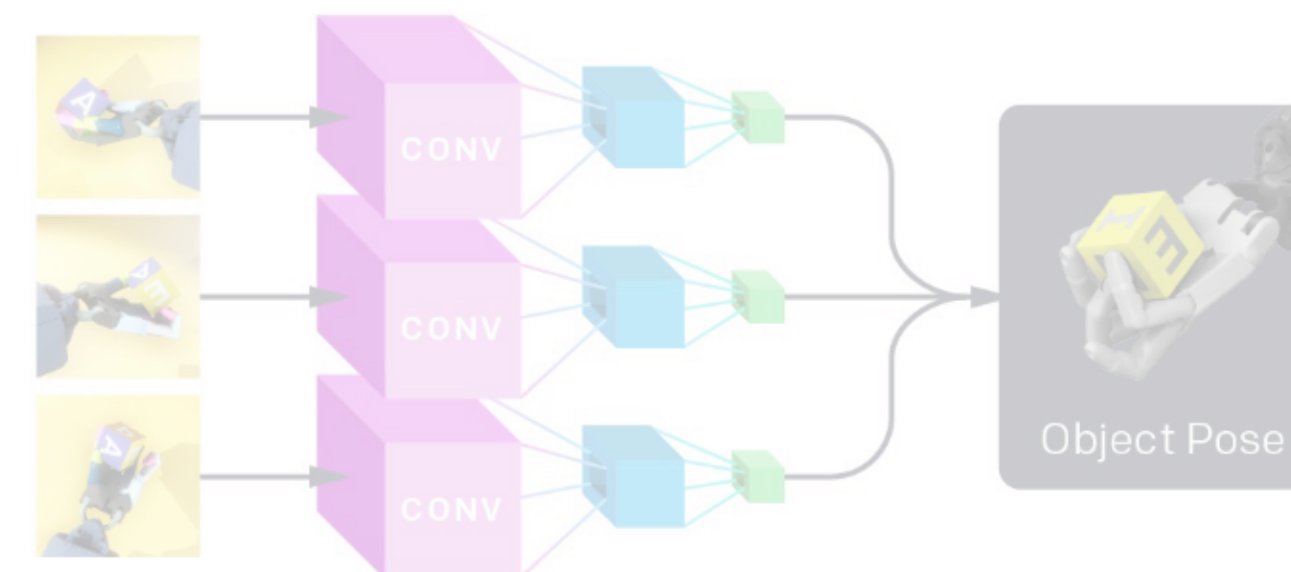
A Distributed workers collect experience on randomized environments at large scale.



B We train a control policy using reinforcement learning. It chooses the next action based on fingertip positions and the object pose.

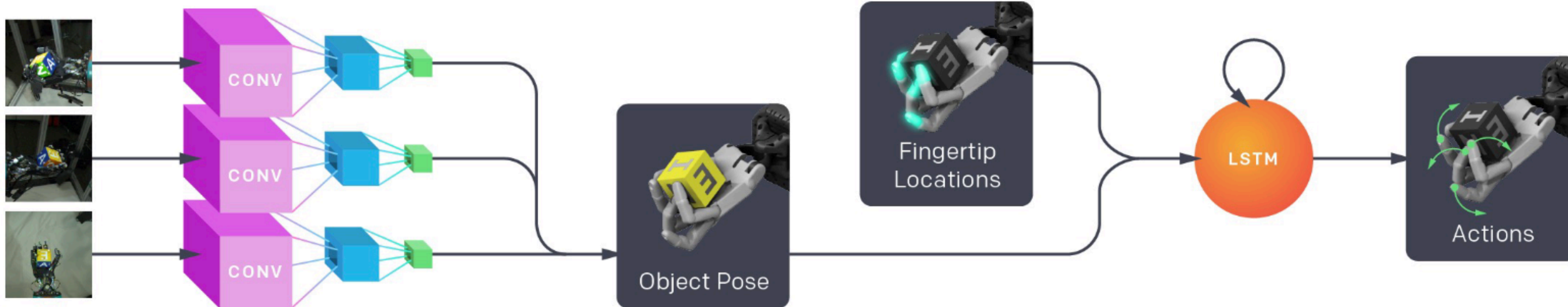


C We train a convolutional neural network to predict the object pose given three simulated camera images.



D We combine the pose estimation network and the control policy to transfer to the real world.

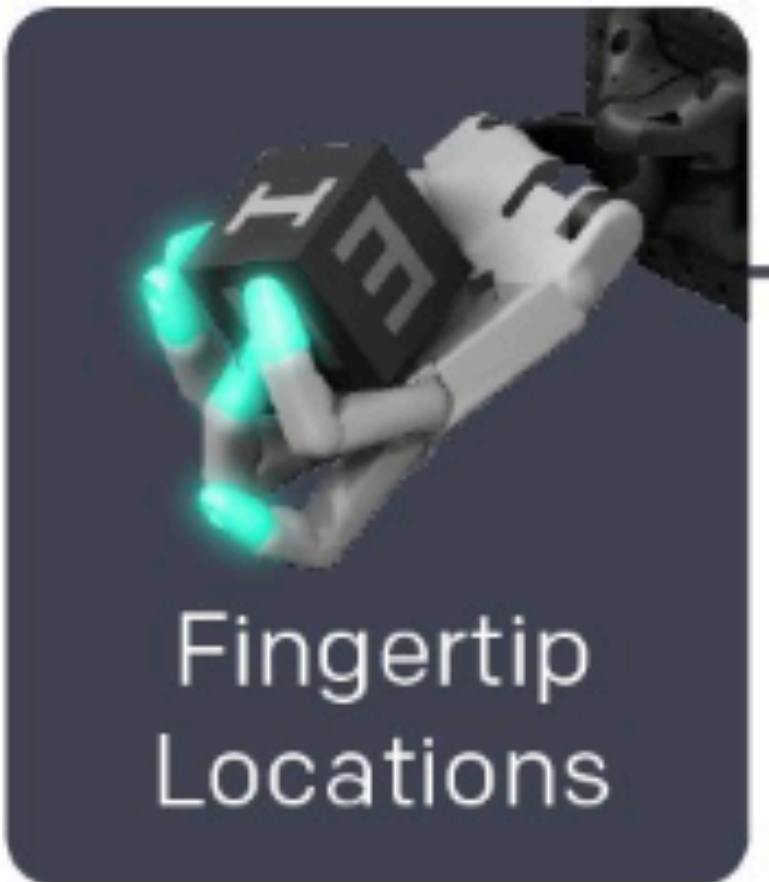
Real



S, *A*, *R*, *T*

S

, A , R , T



Question: Is the current object pose and fingertip location sufficient to capture state?

S

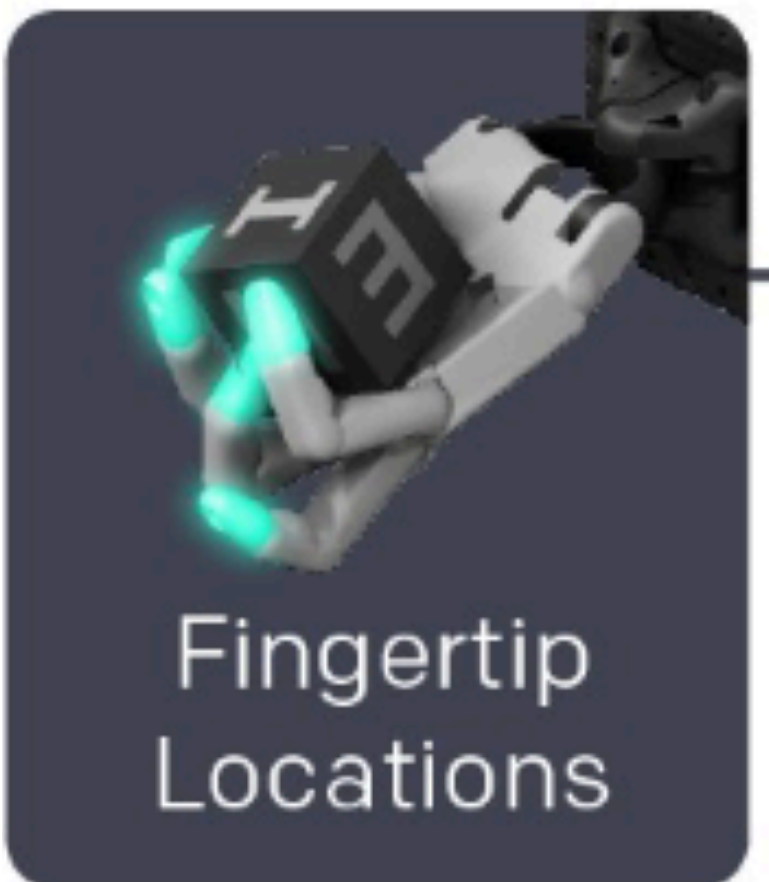
, A , R , T

No!

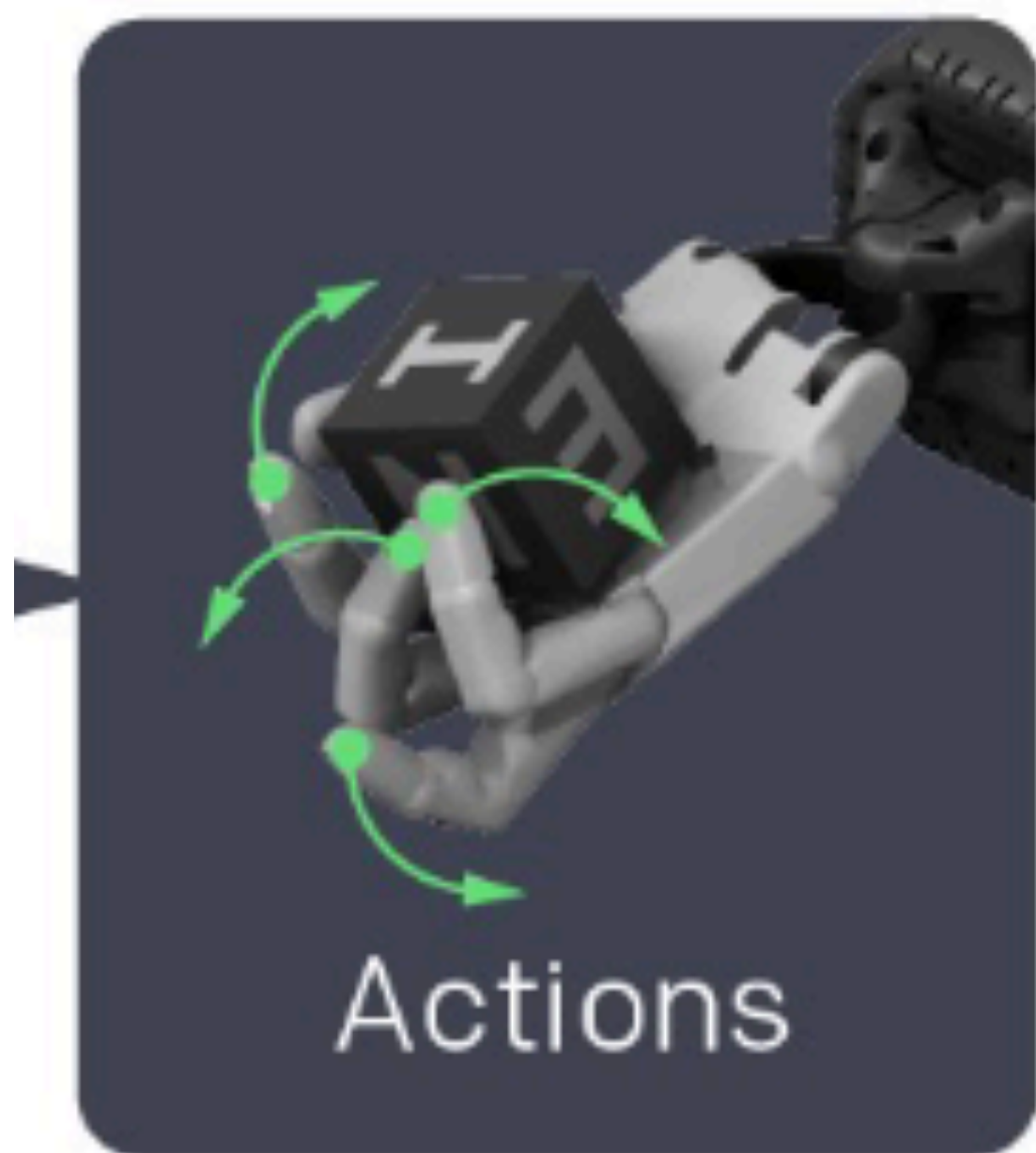
This is merely the current observation of a POMDP

Need to keep a HISTORY

E.g. History of observations can reveal the weight of the object or how fast the index finger can move.



S, A, R, T



S , A , R , \mathcal{T}

The reward given at timestep t is $r_t = d_t - d_{t+1}$, where d_t and d_{t+1} are the rotation angles between the desired and current object orientations before and after the transition, respectively. We give an additional reward of 5 whenever a goal is achieved and a reward of -20 (a penalty) whenever the object is dropped. More information about the simulation environment can be found in Appendix C.1.

Sim2Real as Transferring MDPs

$$\hat{S}, A, R, \hat{\mathcal{T}} \rightarrow S, A, R, \mathcal{T}$$

Sim Real

There will be a mismatch in state representations and transition

Our policy needs to be robust to this mismatch

Key Idea: Add in Randomization in Sim

1. Randomize the observation

Observation noise. To better mimic the kind of noise we expect to experience in reality, we add Gaussian noise to policy observations. In particular, we apply a correlated noise which is sampled once per episode as well as an uncorrelated noise sampled at every timestep.

Key Idea: Add in Randomization in Sim

1. Randomize the observation
2. Randomize the physics

Physics randomizations. Physical parameters like friction are randomized at the beginning of every episode and held fixed. Many parameters are centered on values found during model calibration in an effort to make the simulation distribution match reality more closely. [Table 1](#) lists all physics parameters that are randomized.

Key Idea: Add in Randomization in Sim

1. Randomize the observation
2. Randomize the physics
3. Unmodeled effects

Unmodeled effects. The physical robot experiences many effects that are not modeled by our simulation. To account for imperfect actuation, we use a simple model of motor backlash and introduce action delays and action noise before applying them in simulation. Our motion capture setup sometimes loses track of a marker temporarily, which we model by freezing the position of a simulated marker with low probability for a short period of time in simulation. We also simulate marker occlusion by freezing its simulated position whenever it is close to another marker or the object. To handle additional unmodeled dynamics, we apply small random forces to the object. Details on the concrete implementation are available in Appendix C.2.

Key Idea: Add in Randomization in Sim

Visual appearance randomizations. We randomize the following aspects of the rendered scene: camera positions and intrinsics, lighting conditions, the pose of the hand and object, and the materials and textures for all objects in the scene. [Figure 4](#) depicts some examples of these randomized environments. Details on the randomized properties and their ranges are available in Appendix C.2.



1. Randomize the observation
2. Randomize the physics
3. Unmodeled effects
4. Visual randomization

Today's class

- ☑ Practical MBRL
(Only observations, complex dynamics)
- ☑ Learning Models: The DREAMER algorithm
- ☑ Leveraging Physics: Sim2Real
Case study: OpenAI Dactyl Hand