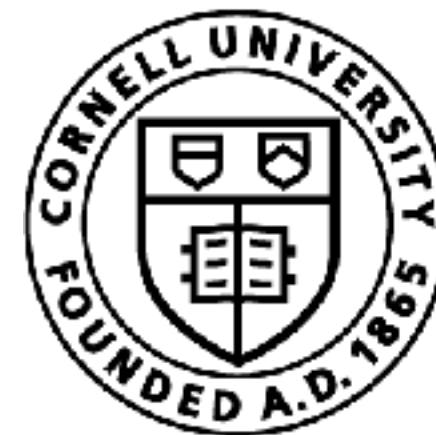
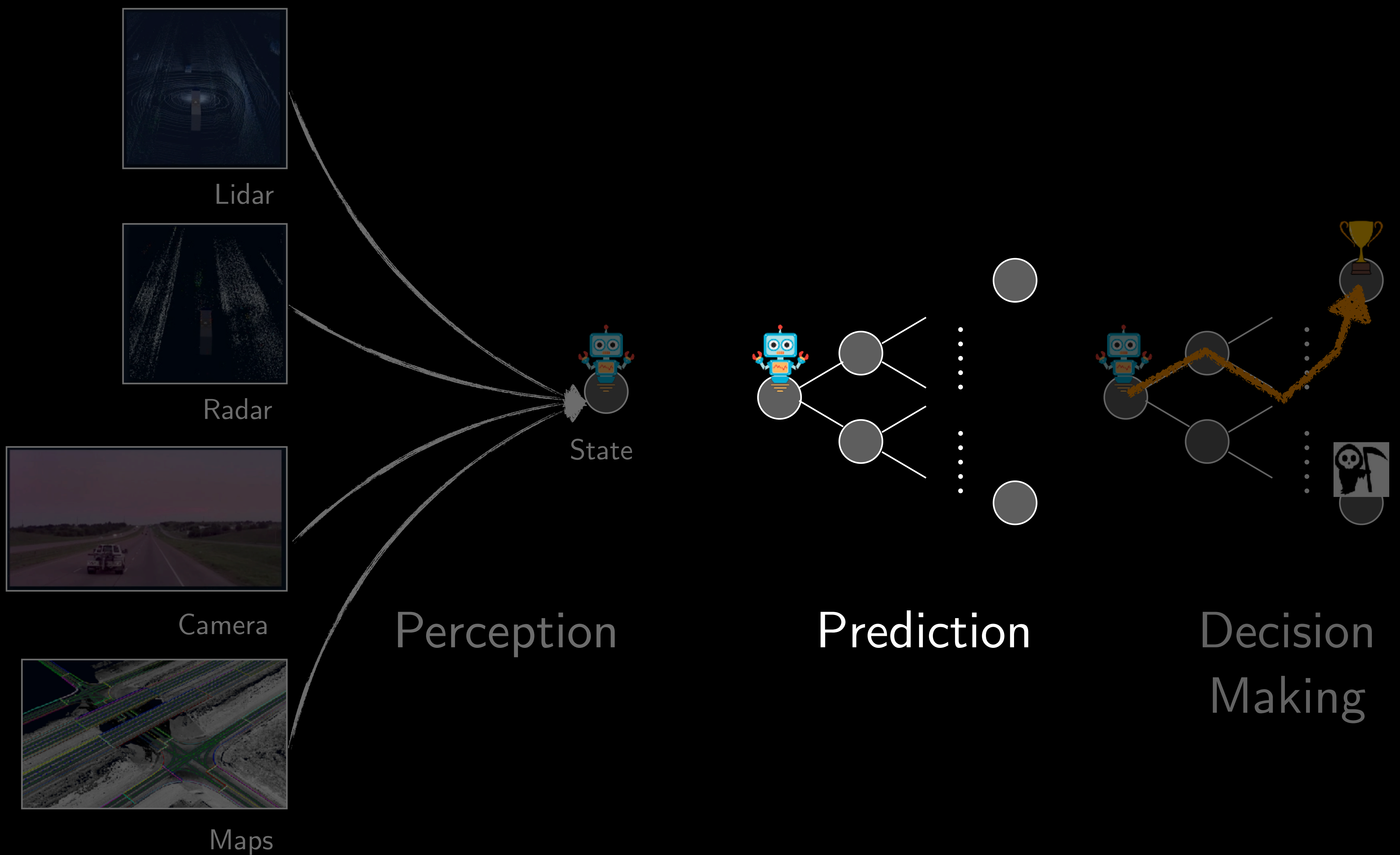


Generative World Models: The Dreamer Models

Sanjiban Choudhury

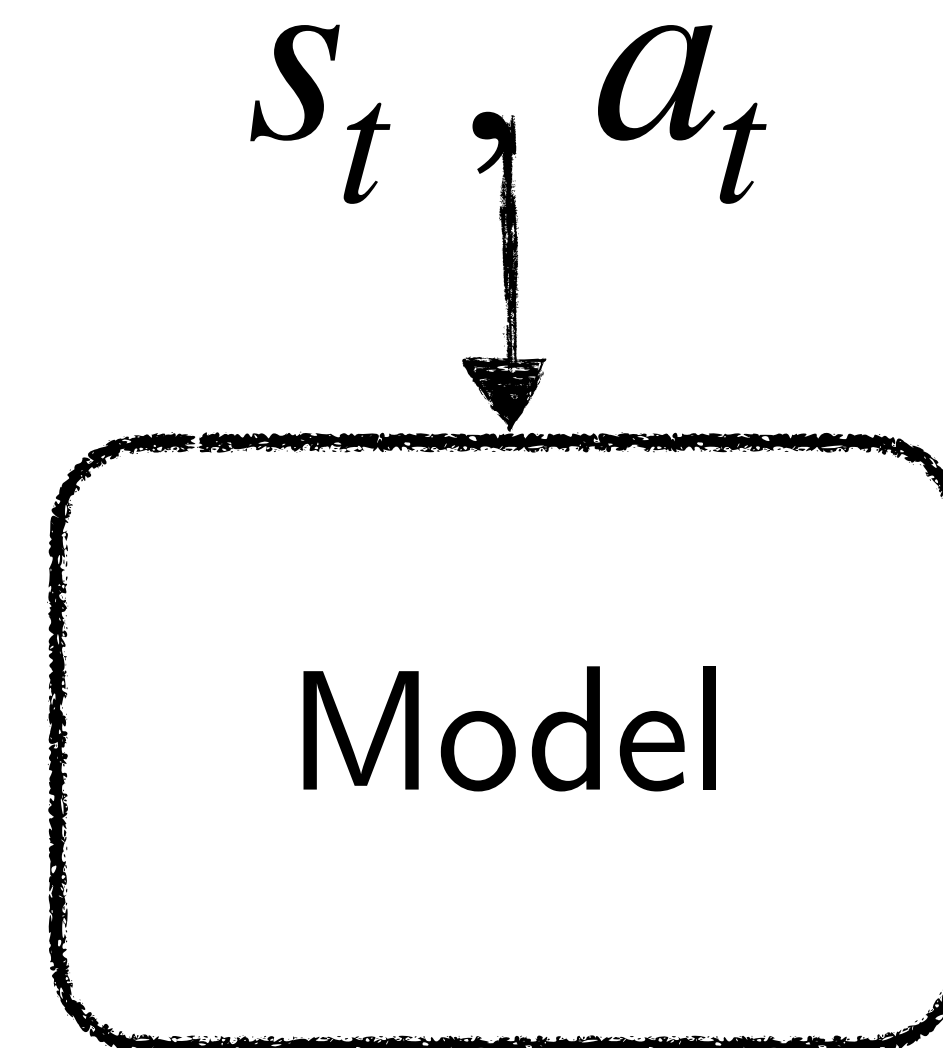


Cornell Bowers CIS
Computer Science

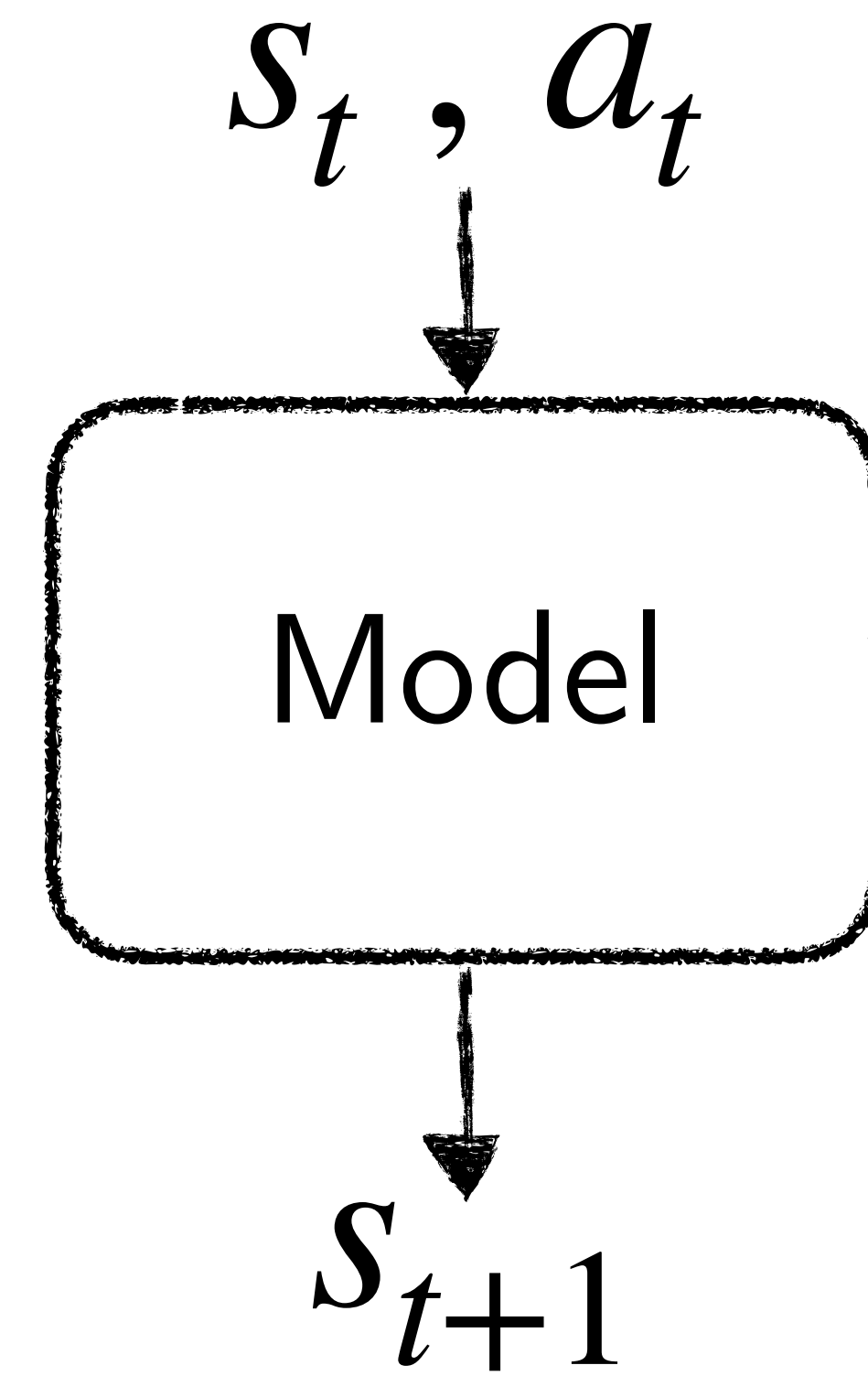


Models.

What is a model?



What is a model?



What is a model?

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

Why Model?

Models are *necessary*

Robots can't just try out random actions in the world!



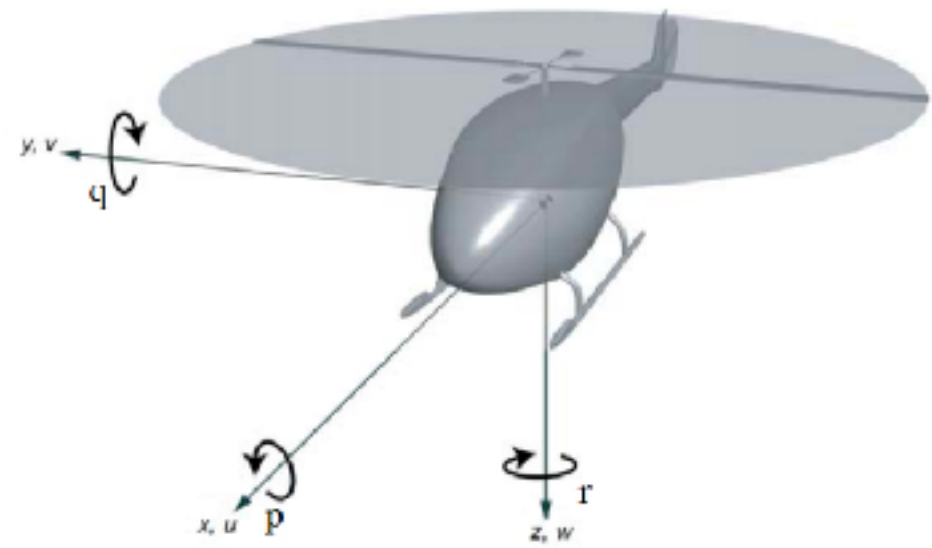
Learning Models.

(Early work in Model Based RL by Pieter Abbeel et al. 2010
https://people.eecs.berkeley.edu/~pabbeel/autonomous_helicopter.html)

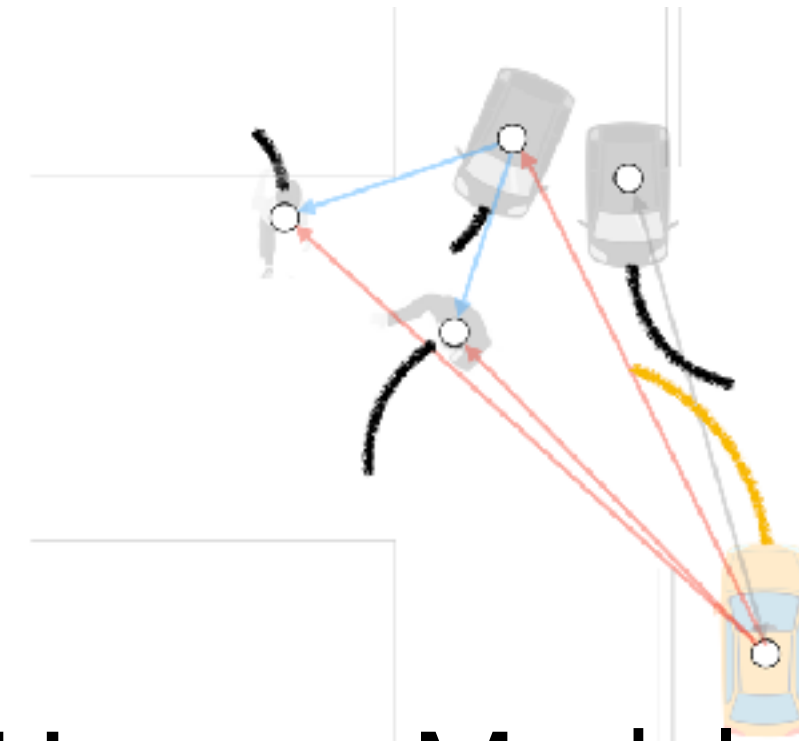


Stanford University Autonomous Helicopter

Models: From Simple to Complex



Physics Models



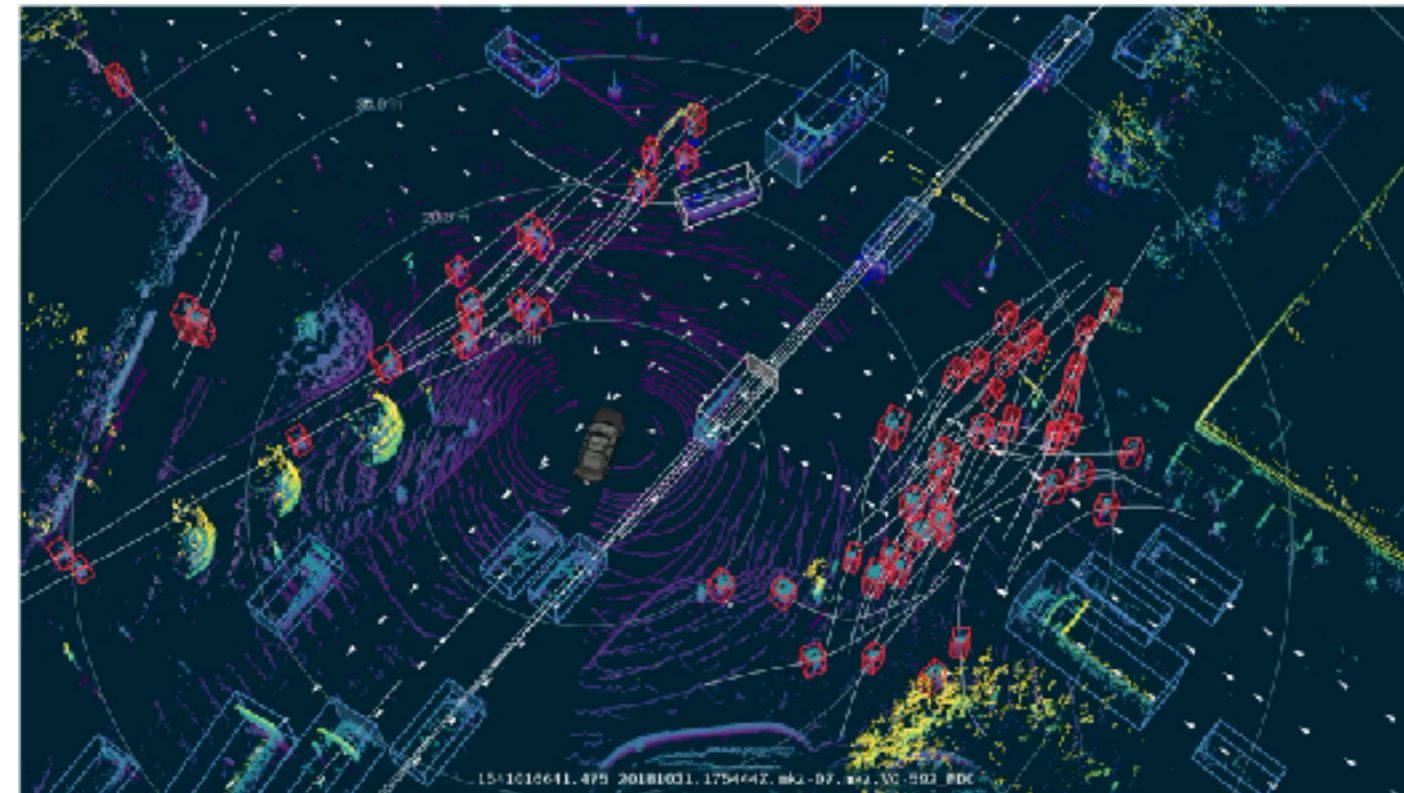
Human Models



Open World Models

Simple

Complex

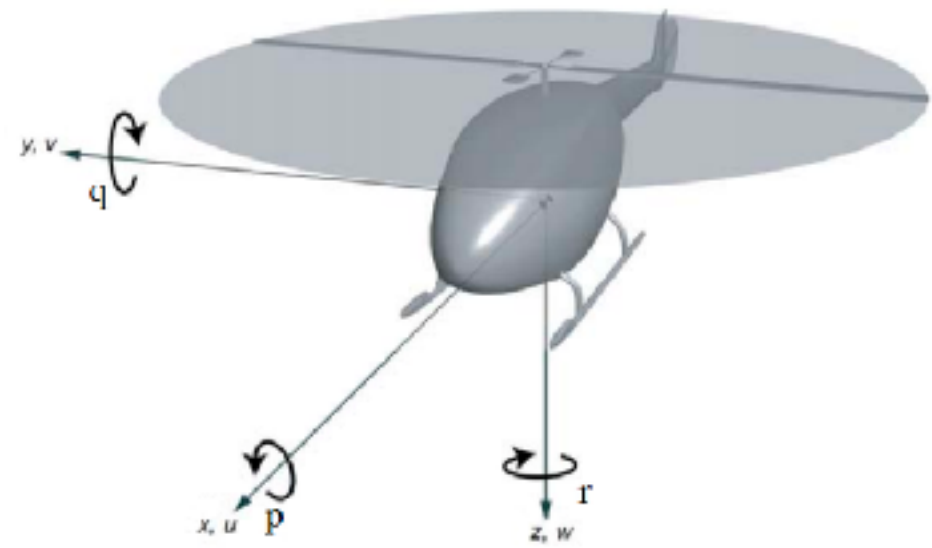


Models: From Simple to Complex

Simple

Complex

Models: From Simple to Complex



Physics Models

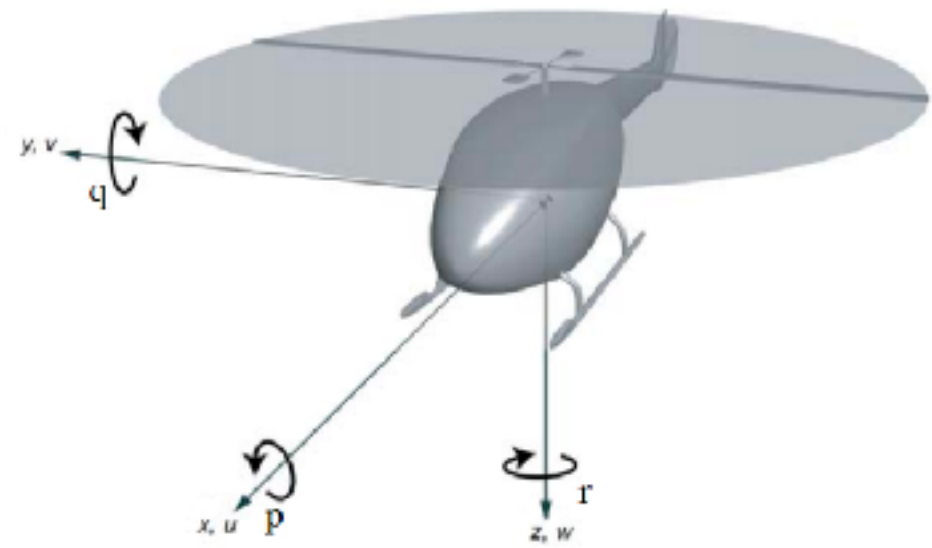
Simple

Know state

Strong prior
on dynamics



Models: From Simple to Complex

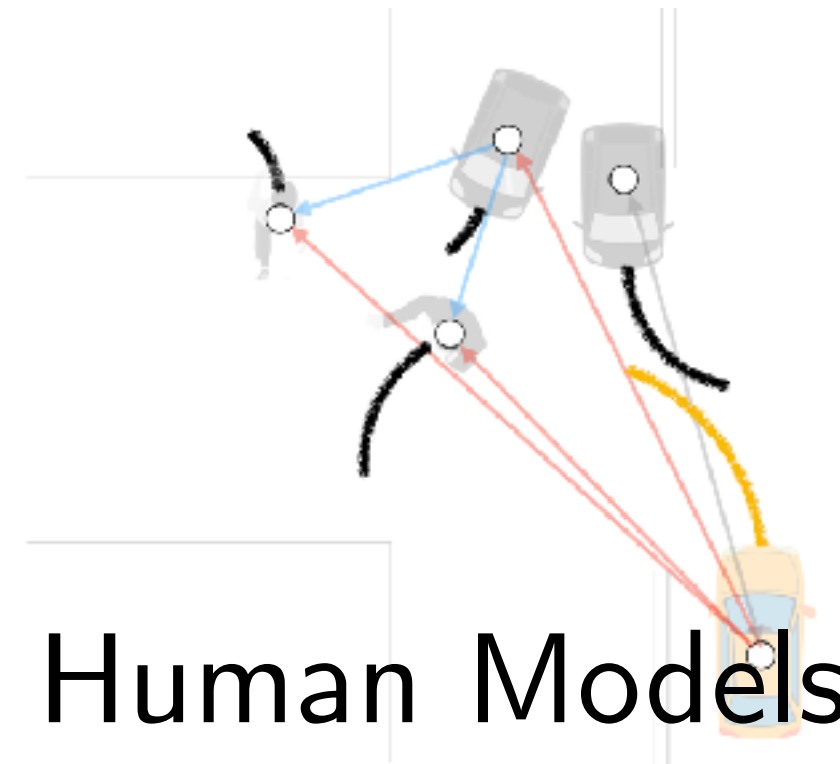


Physics Models

Simple

Know state

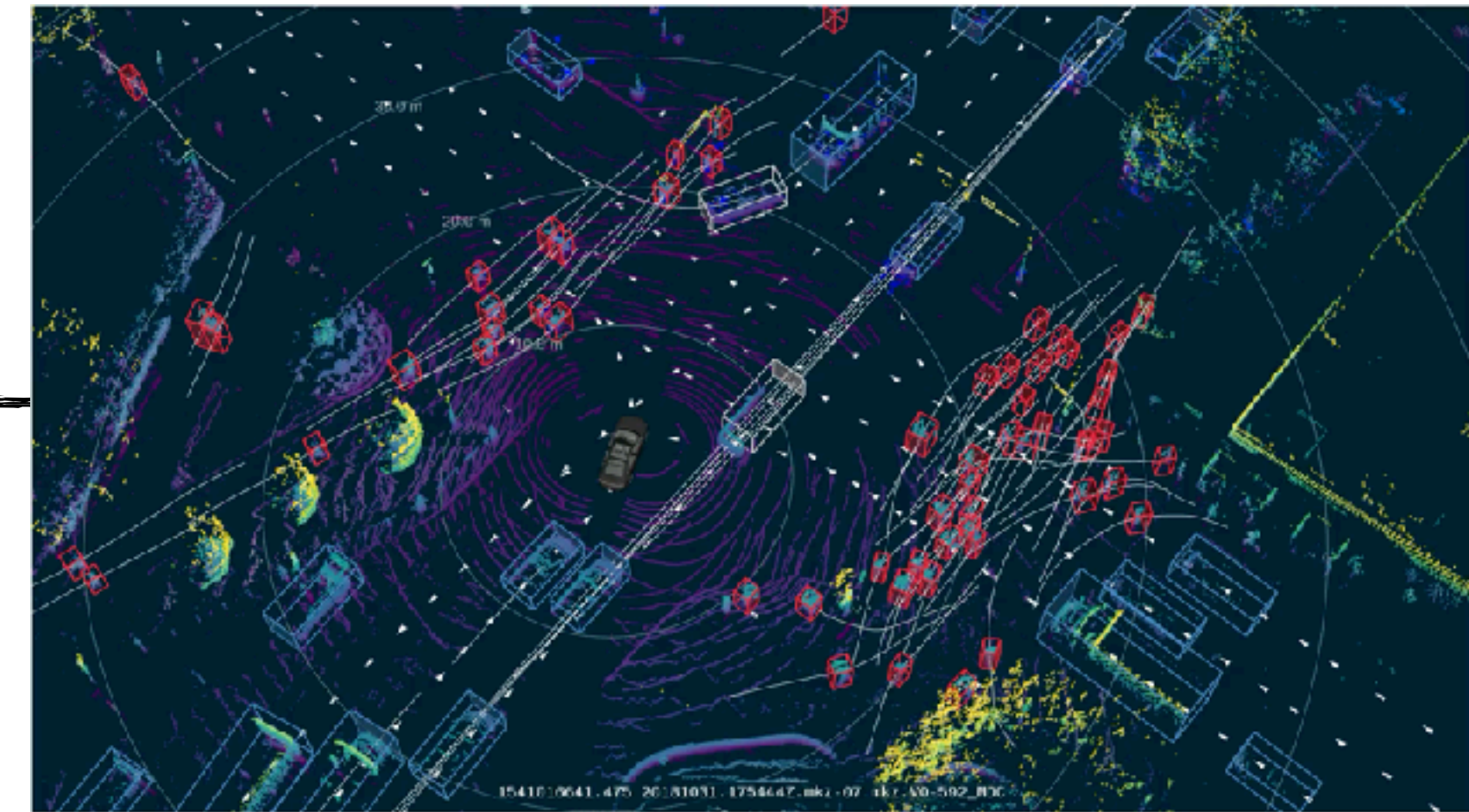
Strong prior
on dynamics



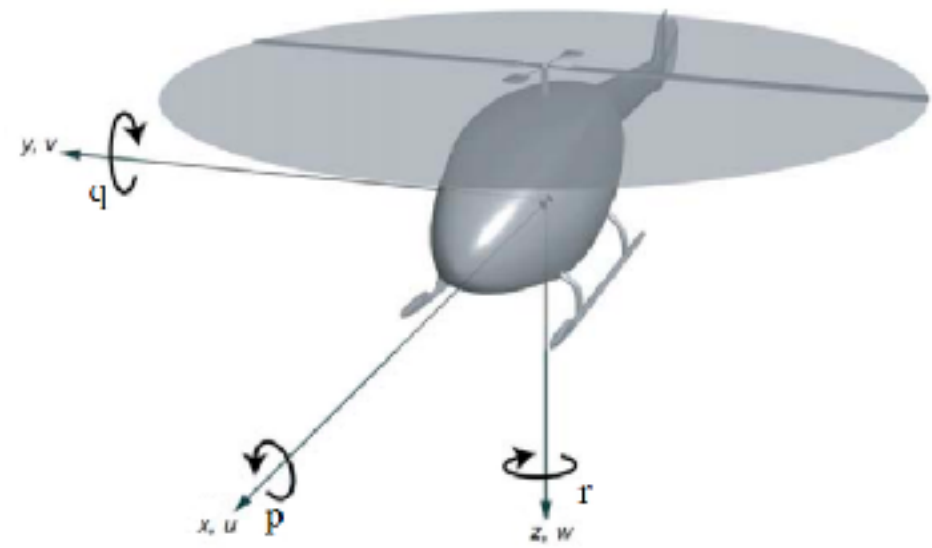
Human Models

Know state

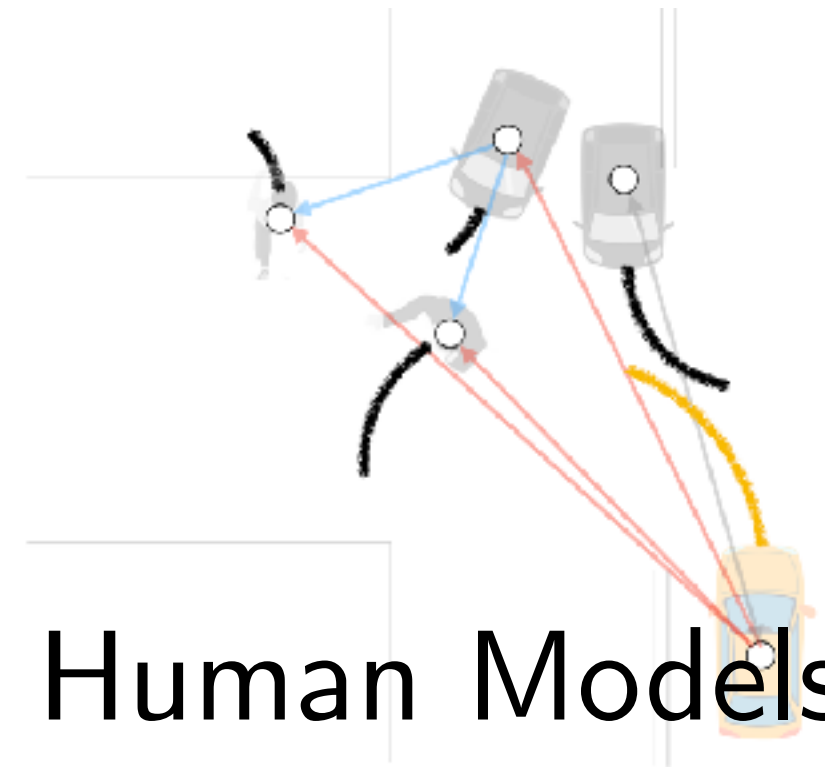
Unknown
dynamics



Models: From Simple to Complex



Physics Models



Human Models



Open World Models

Simple

Complex

Know state

Know state

Unknown state

Strong prior on dynamics

Unknown dynamics

Unknown dynamics

Activity!



Modelling Tamago Sushi



Think-Pair-Share!

Think (30 sec): How would you model making tamago sushi?

Pair: Find a partner

Share (45 sec): Partners exchange ideas



Challenges with learning complex models

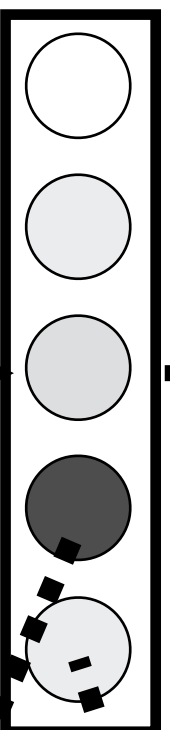
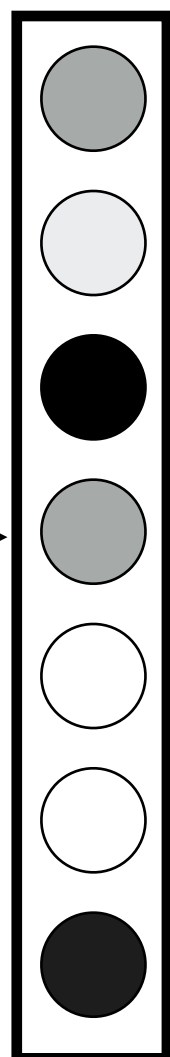
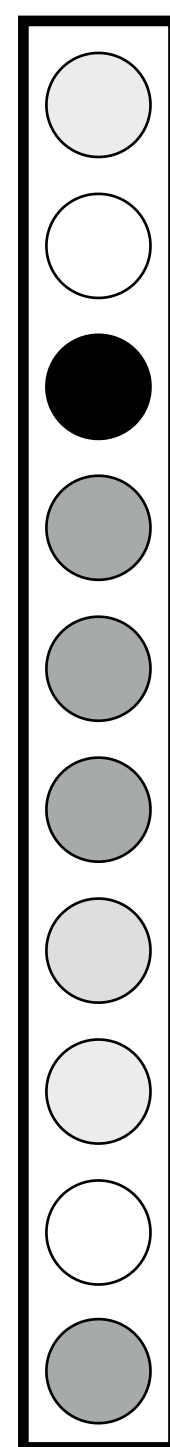
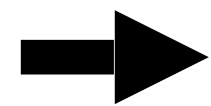
Challenge 1: High-dimensional observations

Challenge 2: Planning with Complex Dynamics

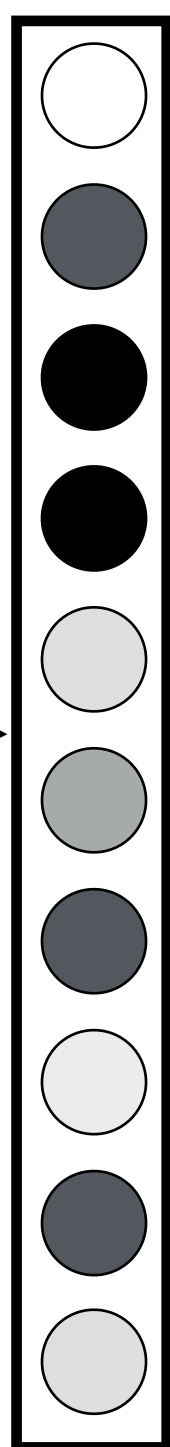
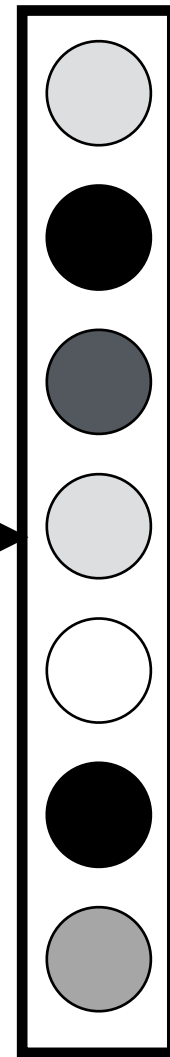
\mathbf{X}



Image



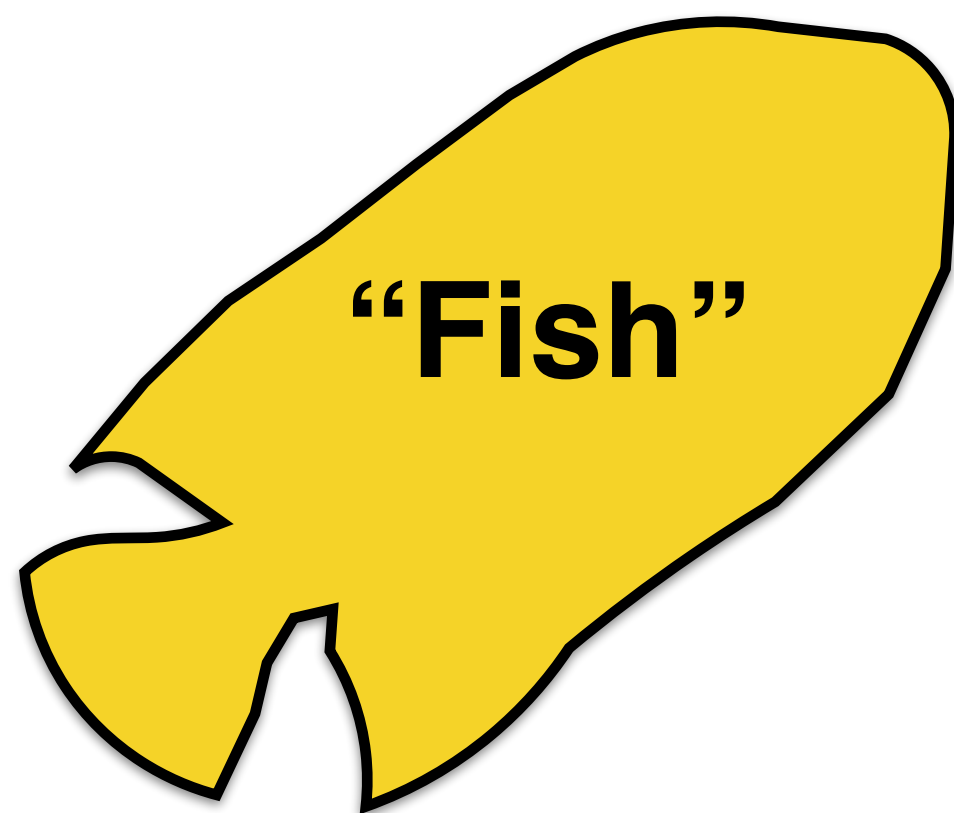
\mathcal{F}



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

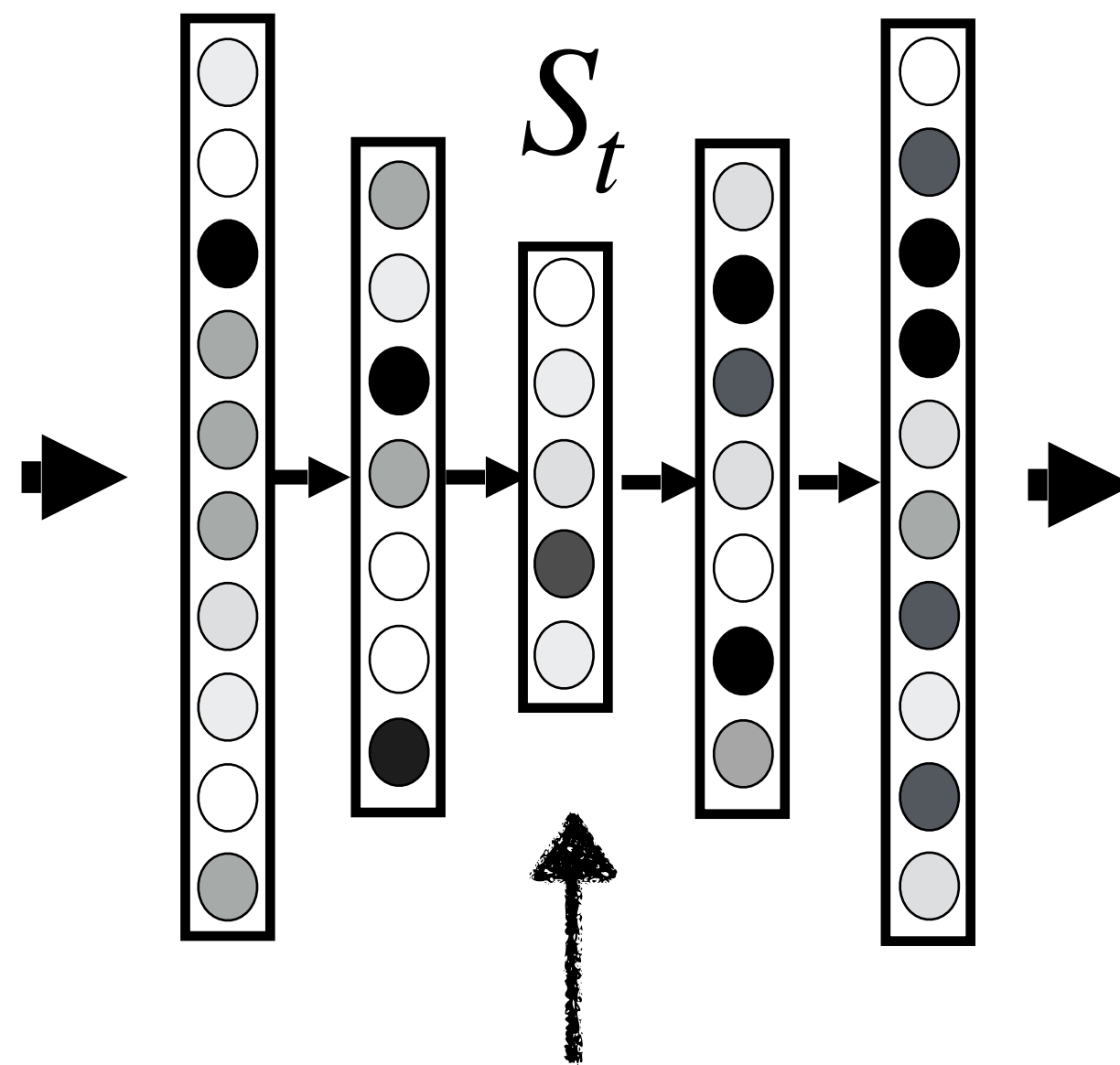


Reconstructed image

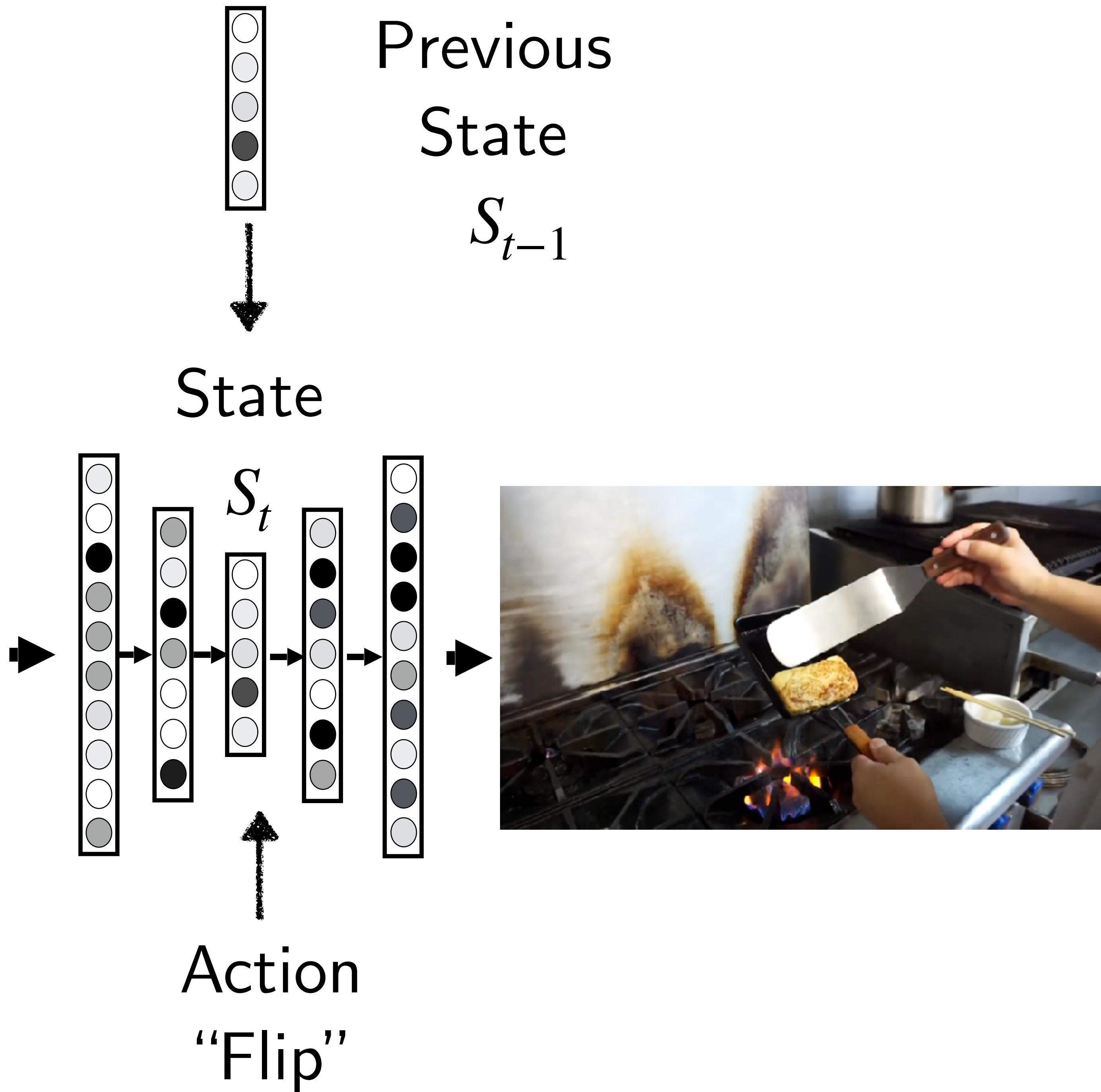




State

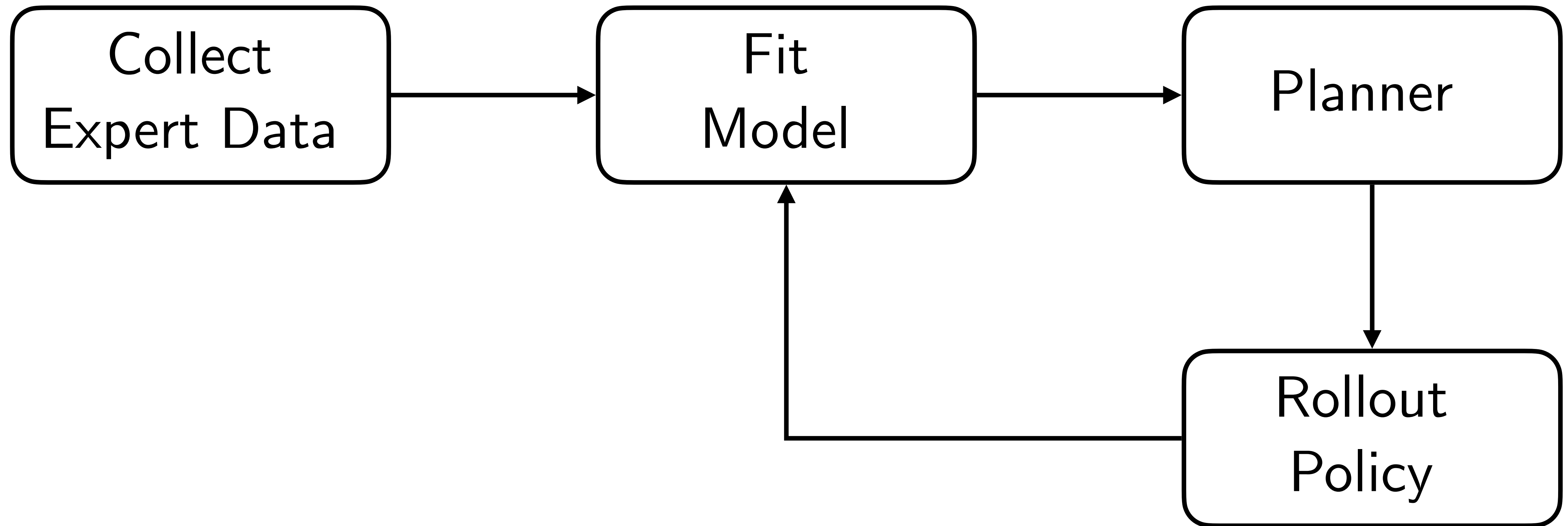


Action
"Flip"



Recall from previous lecture!

(Ross & Bagnell, 2012)



What if we don't have
expert data?





The
DREAMER
Algorithm



DREAM TO CONTROL: LEARNING BEHAVIORS BY LATENT IMAGINATION

Danijar Hafner *

University of Toronto

Google Brain

Timothy Lillicrap

DeepMind

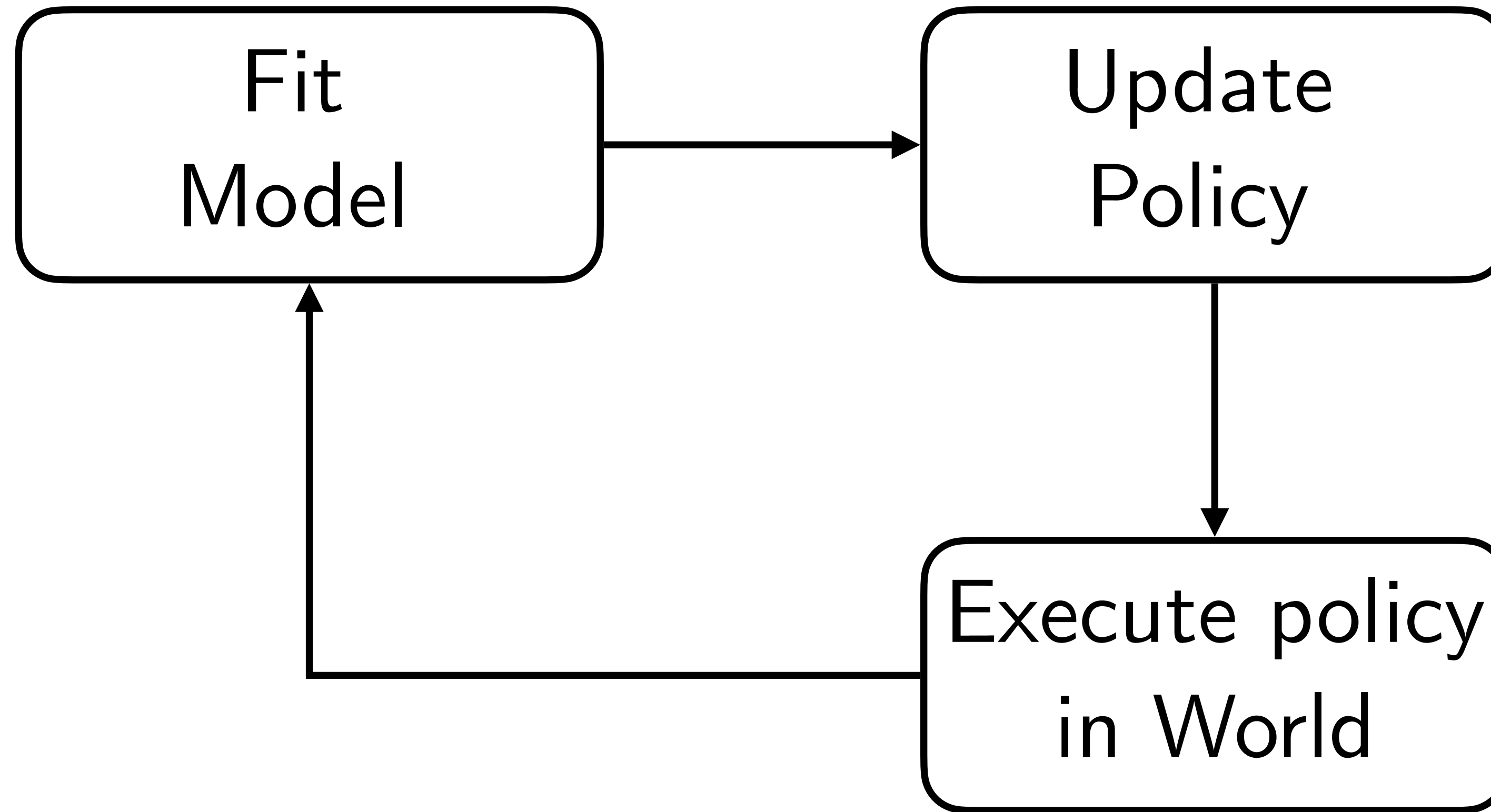
Jimmy Ba

University of Toronto

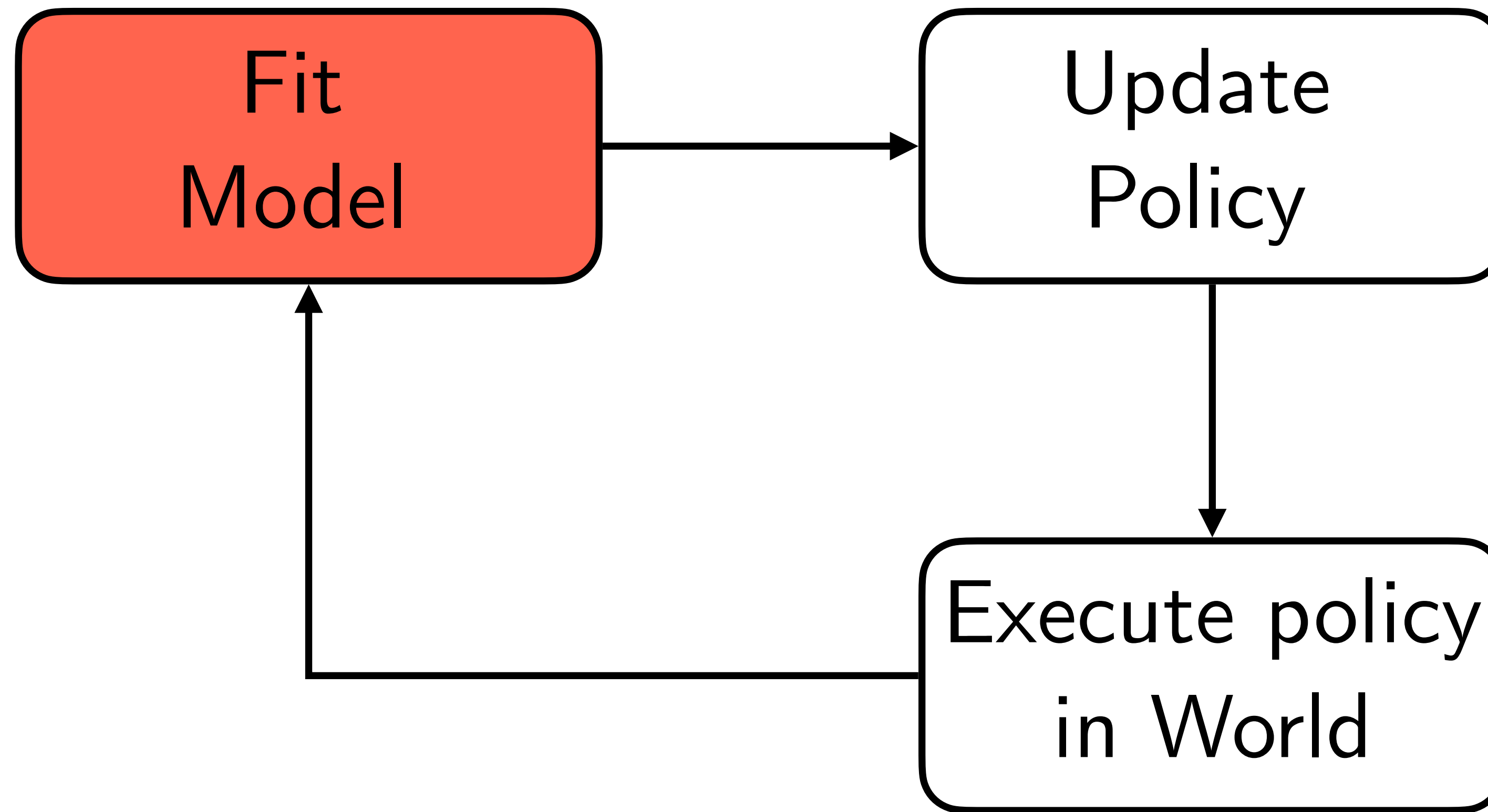
Mohammad Norouzi

Google Brain

DREAMER



DREAMER



Goal: Fit a Model

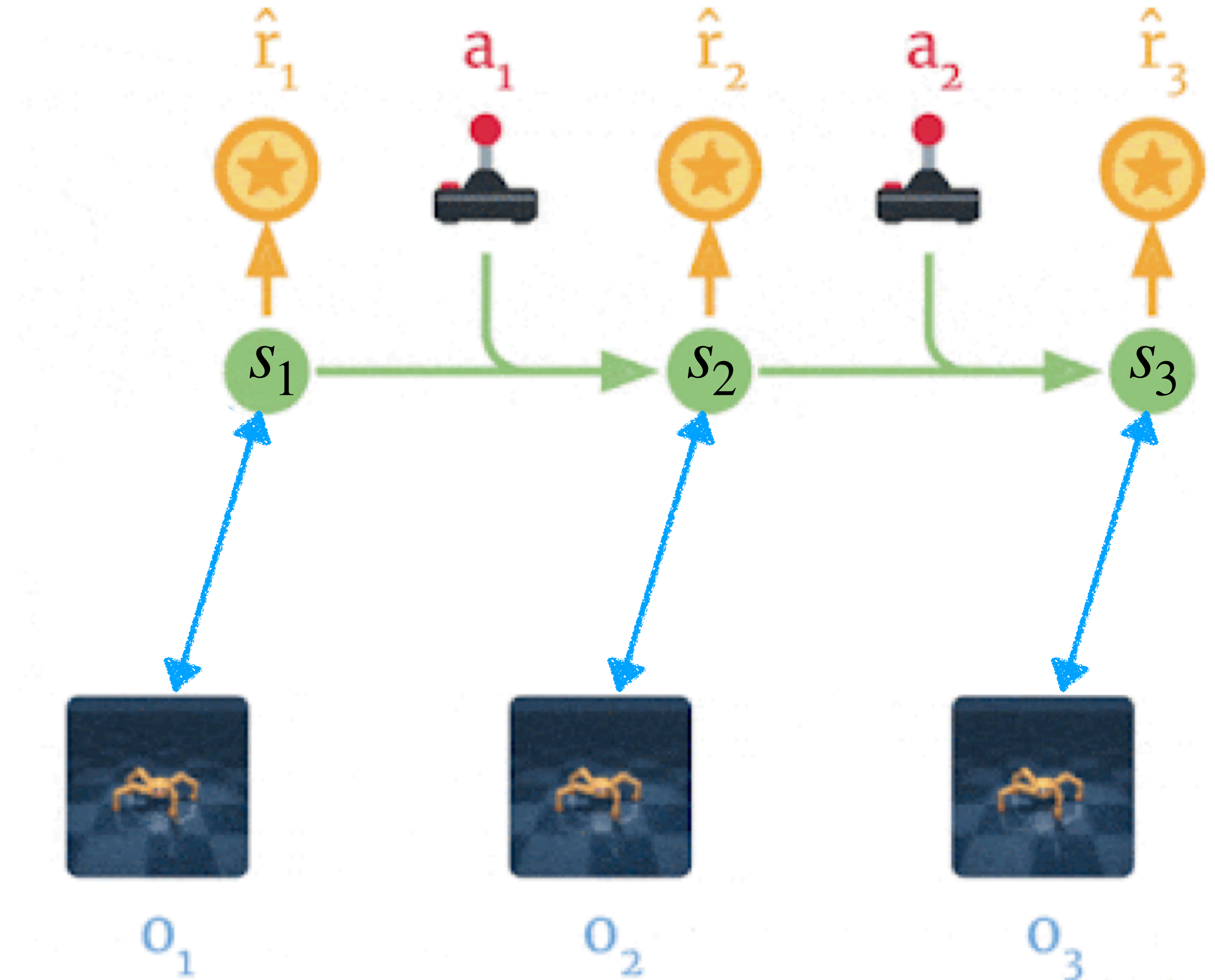
Given:
Observations, rewards,
actions



Goal: Fit a Model

Given:
Observations, rewards,
actions

Predict:
States,
Dynamics Function,
Reward Function

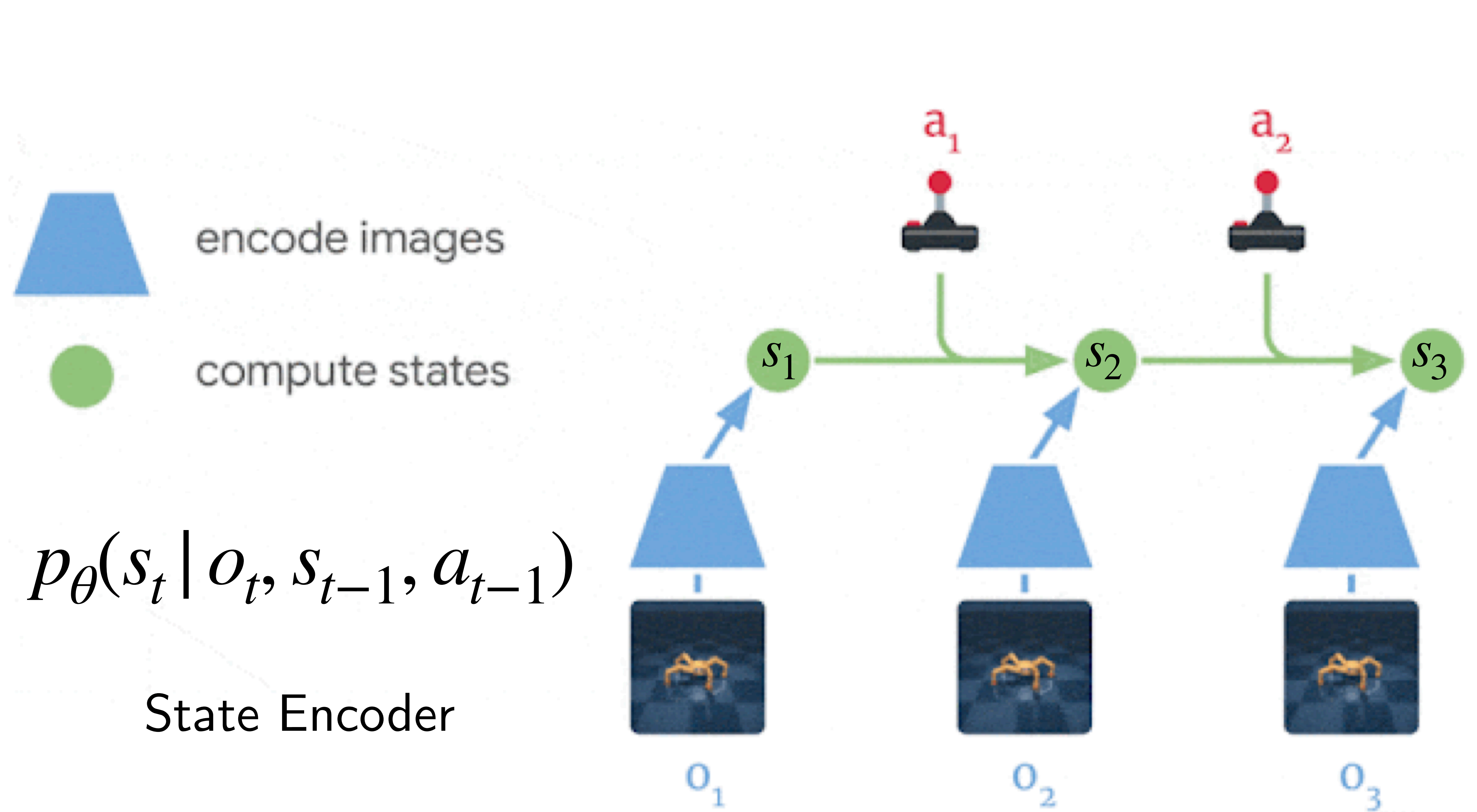


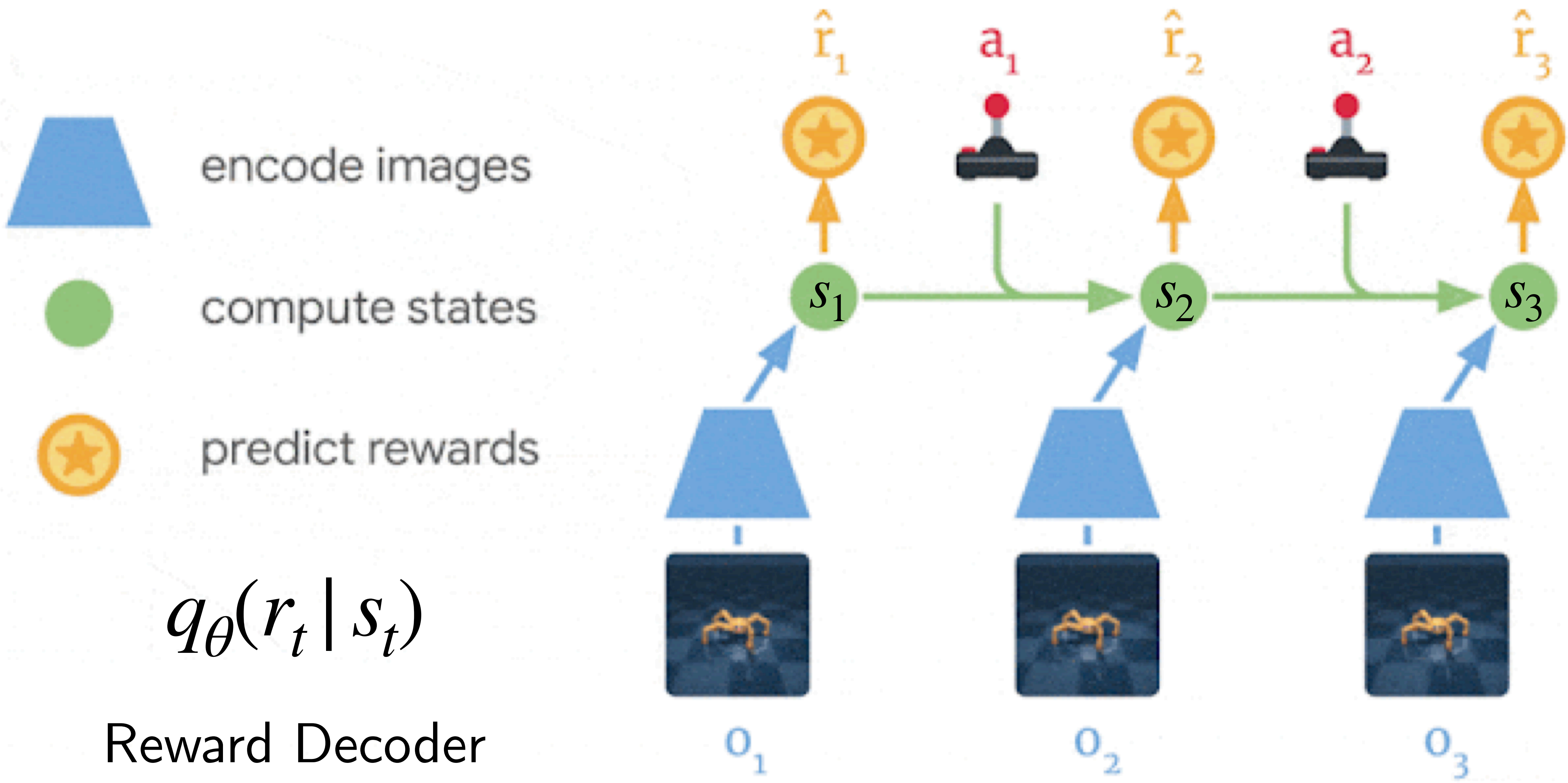
Actions



Observations

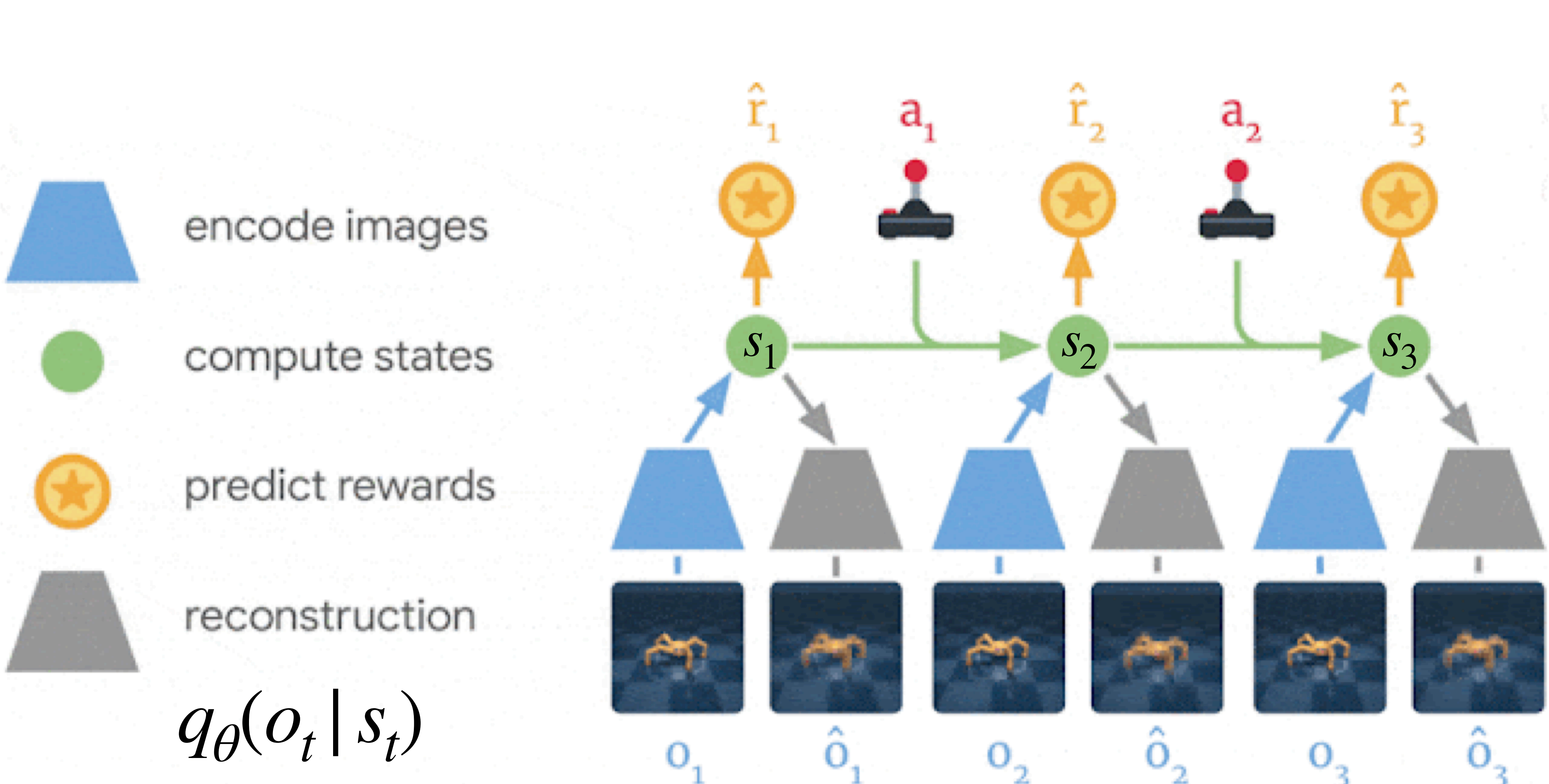






$$q_{\theta}(r_t | s_t)$$

Reward Decoder

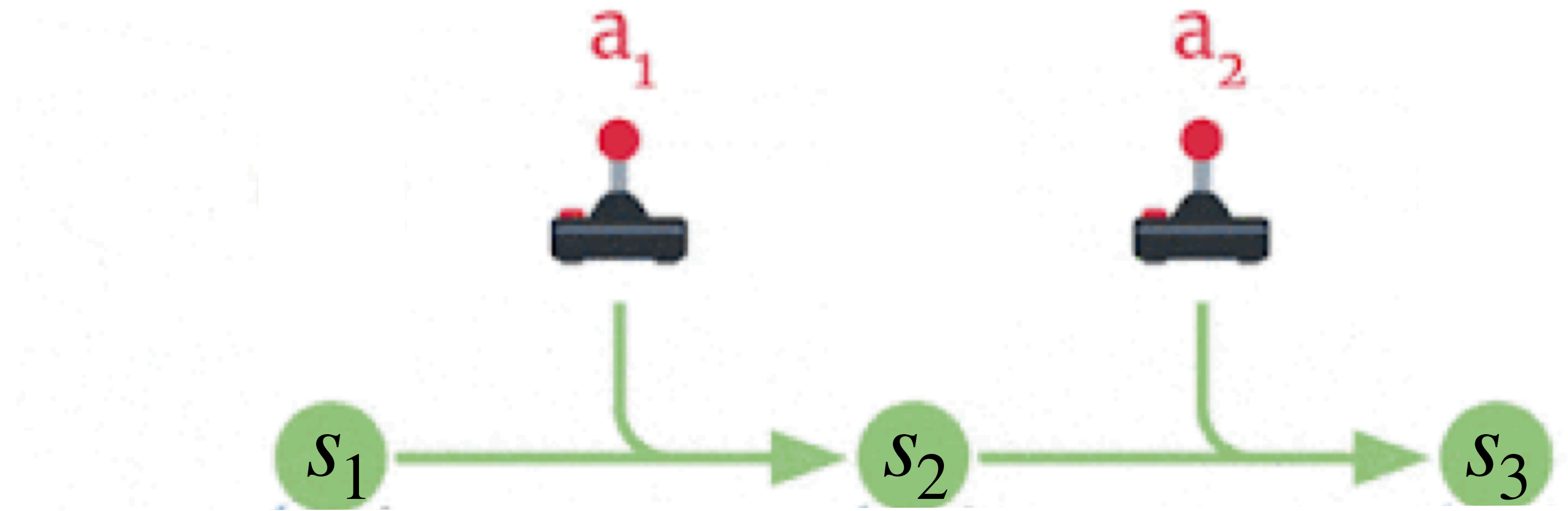


$$q_{\theta}(o_t | s_t)$$

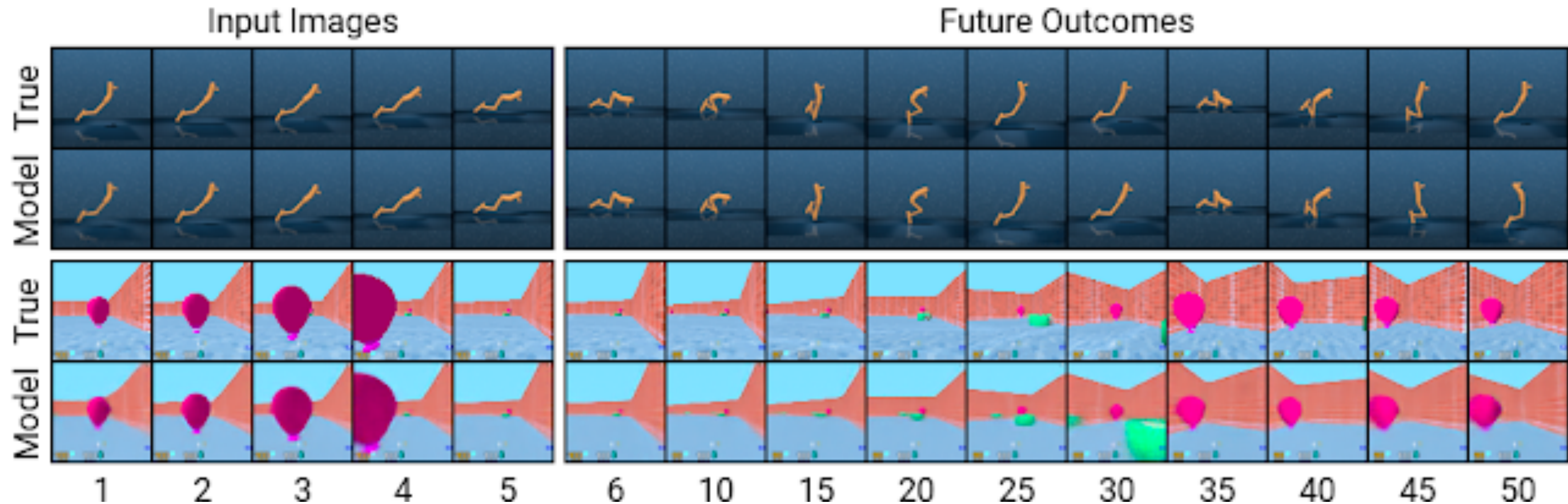
Observation Decoder

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

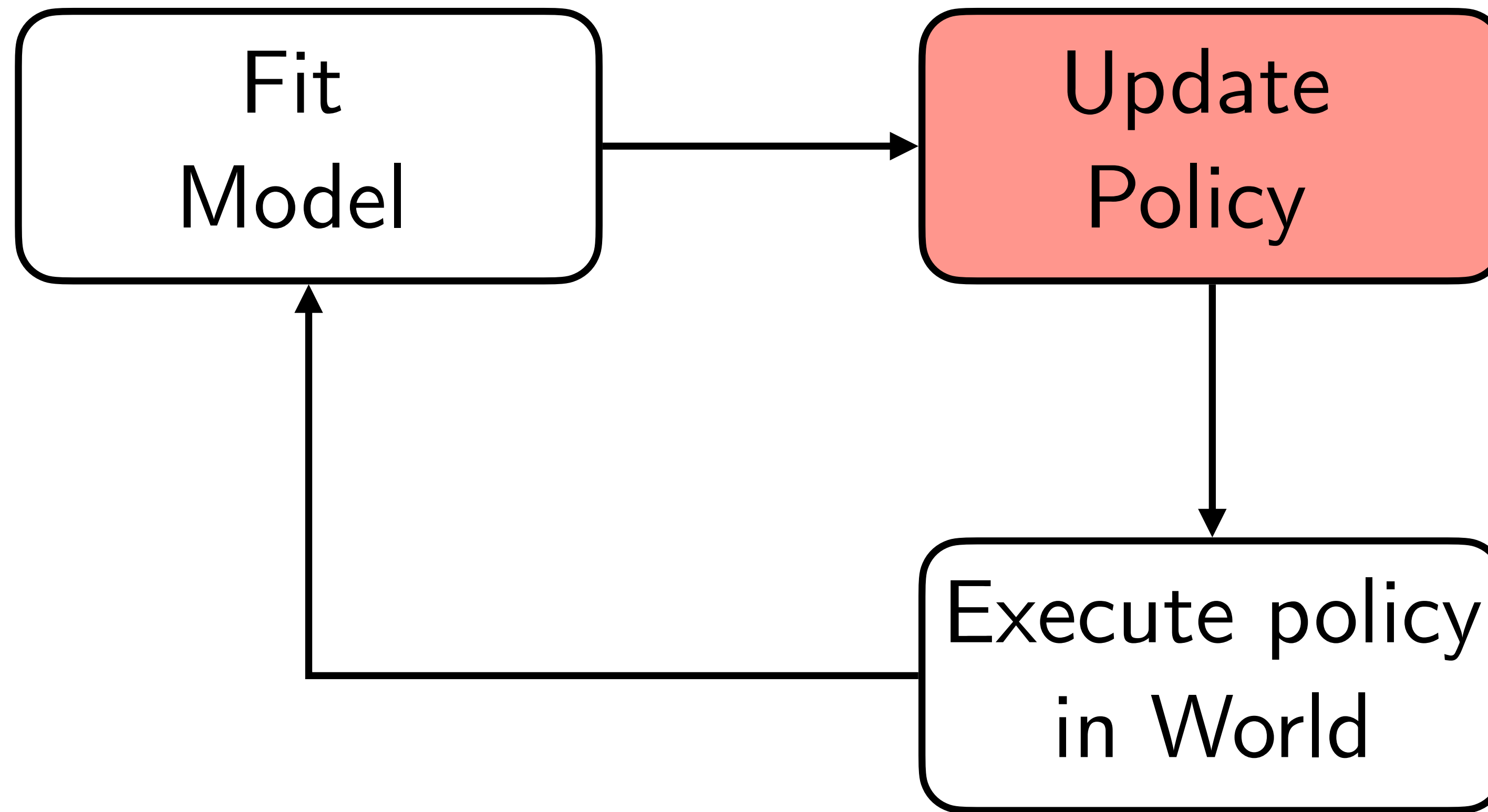
Dynamics
Function



Results: Learning World Model



DREAMER



Goal: Learn a Policy using Actor-Critic

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

Actor

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

Critic

From rollouts in the model

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



O_1



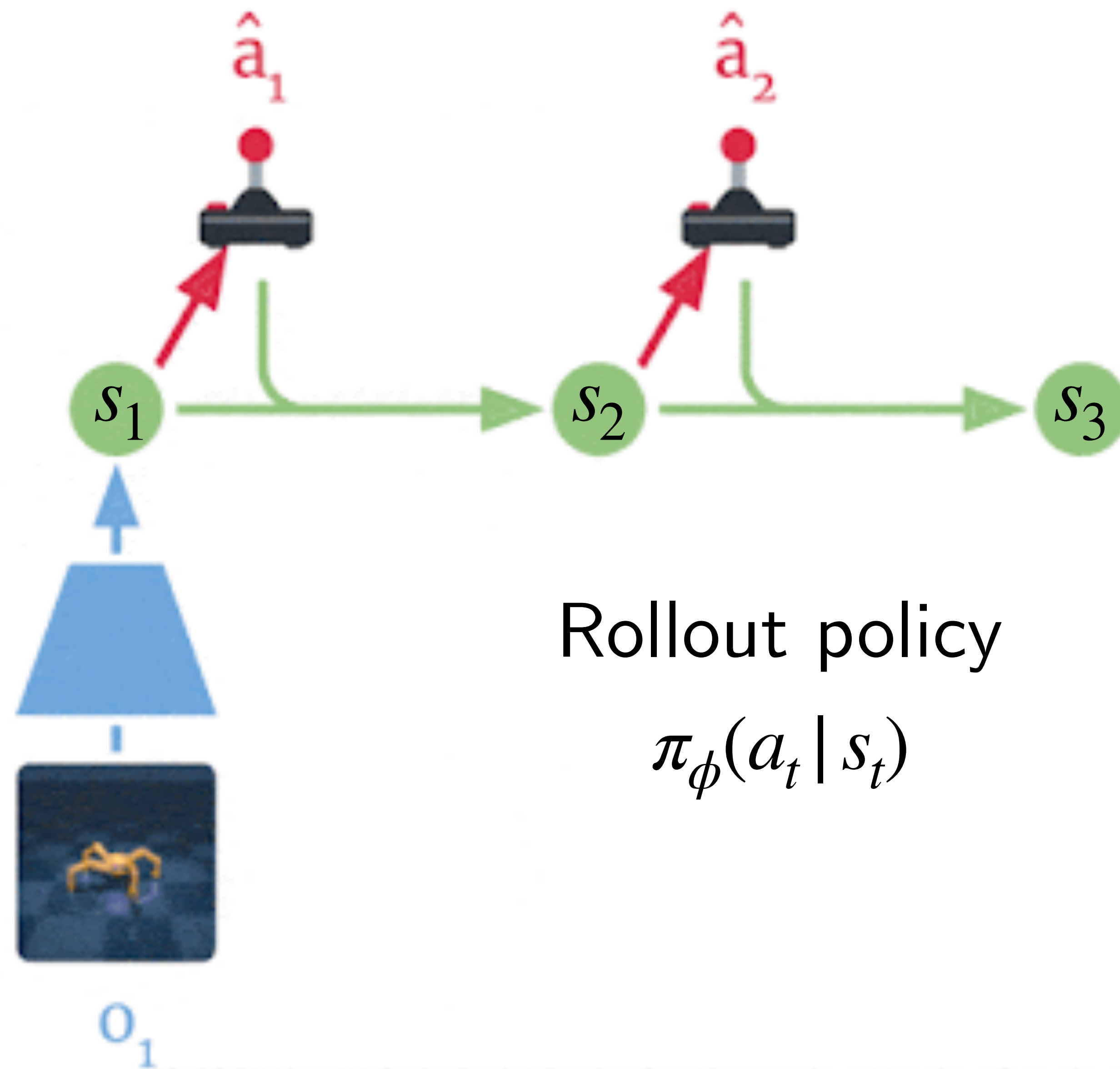
encode images





encode images

imagine ahead





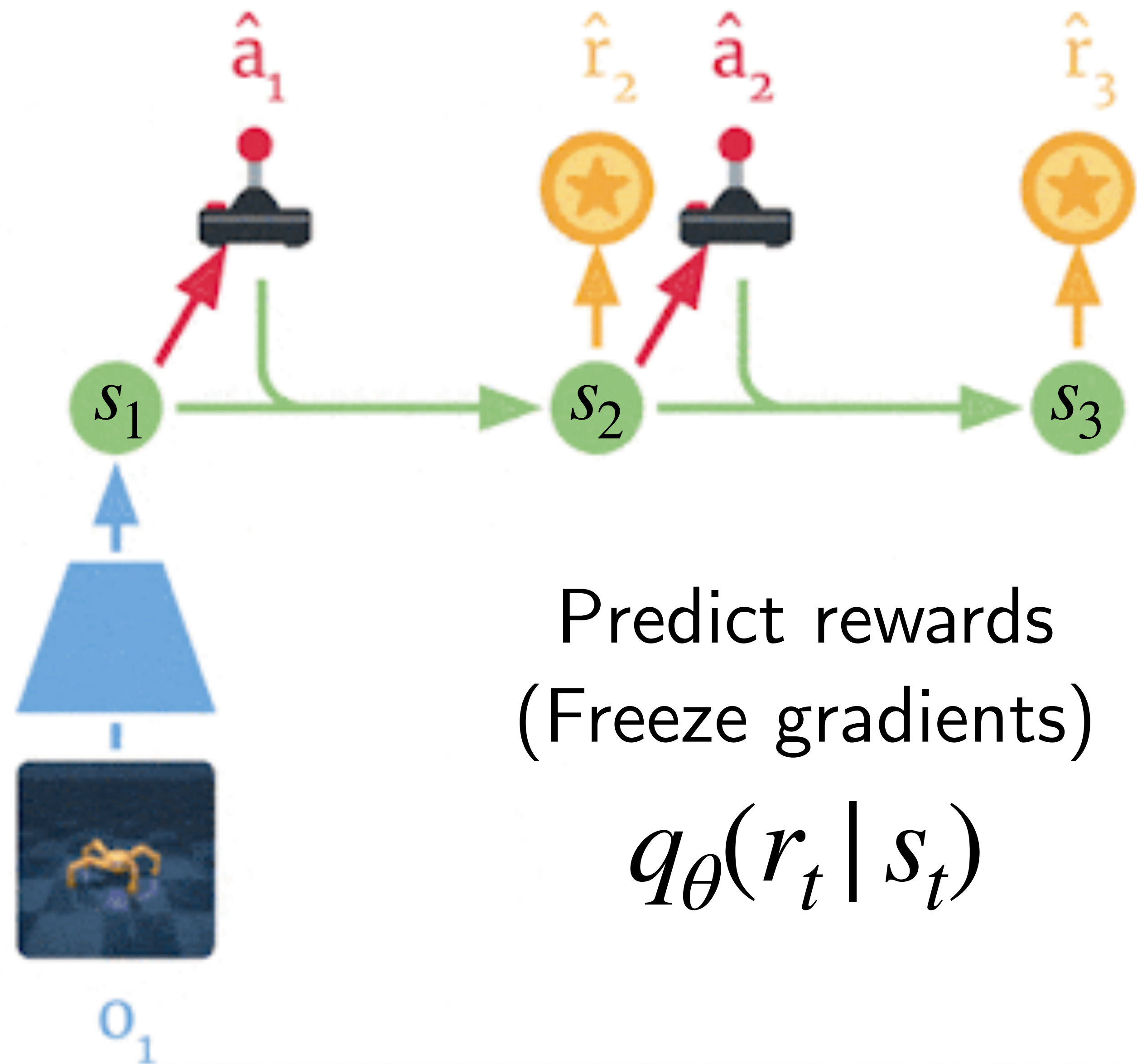
encode images



imagine ahead



predict rewards





encode images



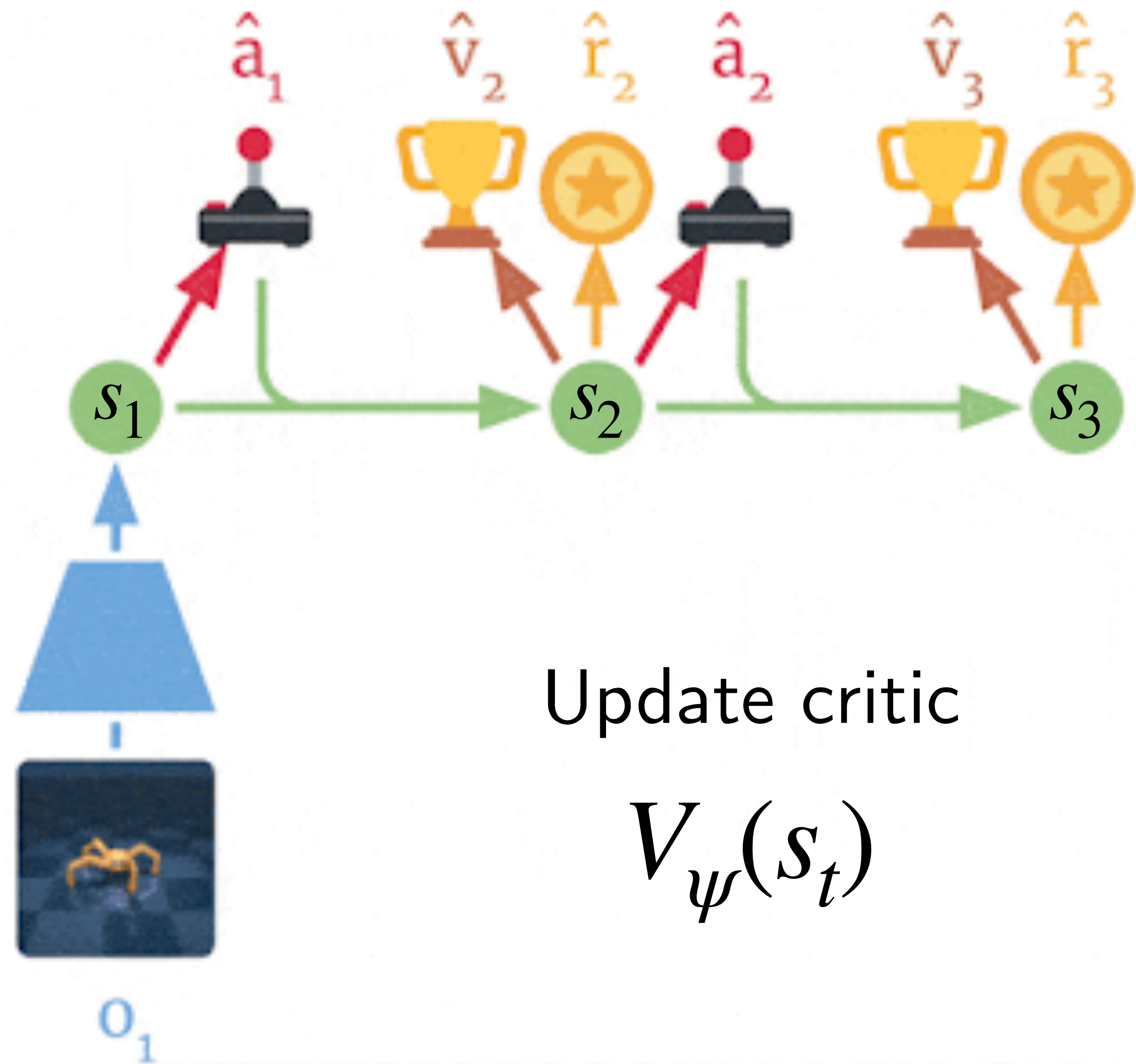
imagine ahead



predict rewards



predict values





encode images



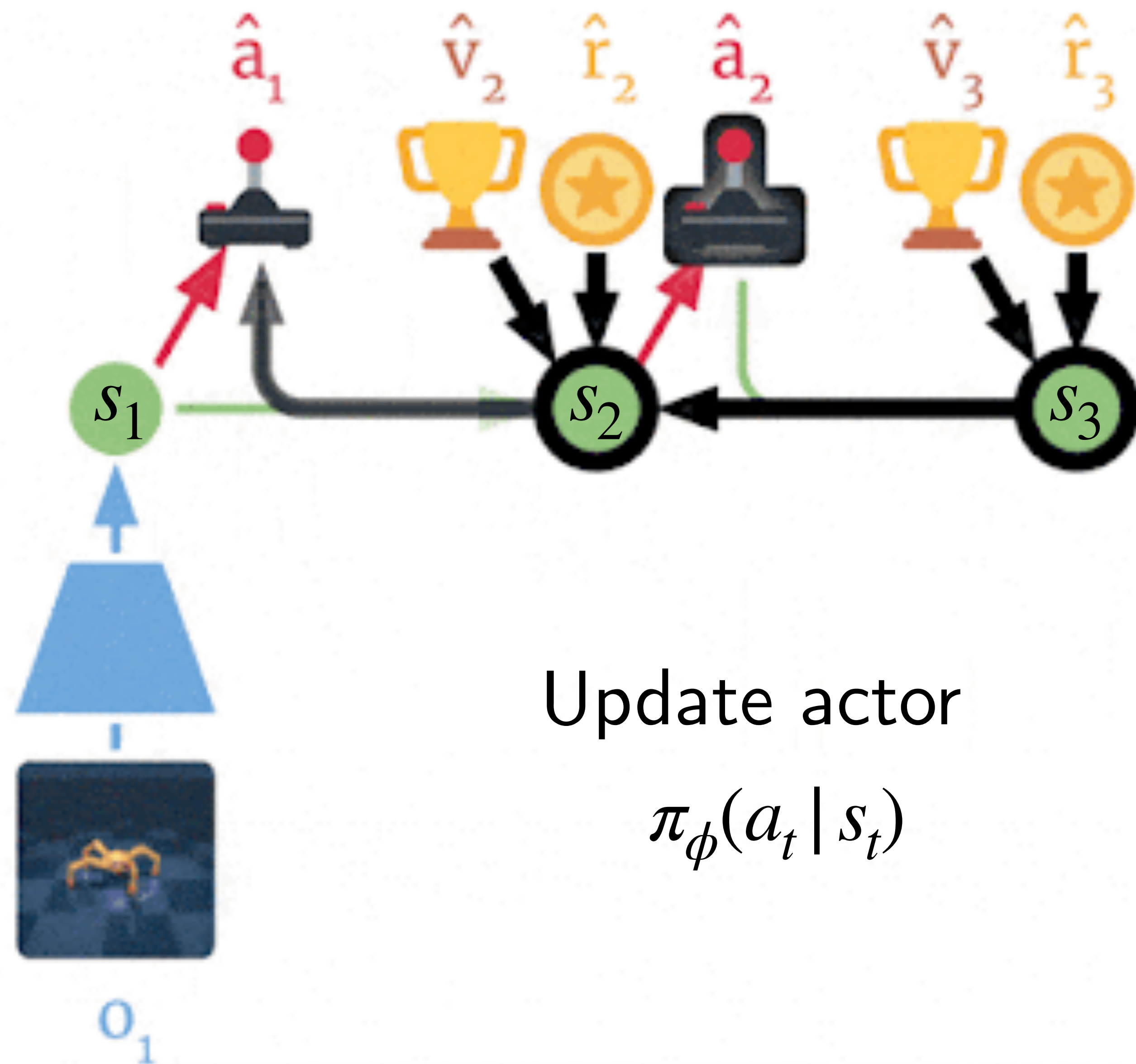
imagine ahead



predict rewards



predict values

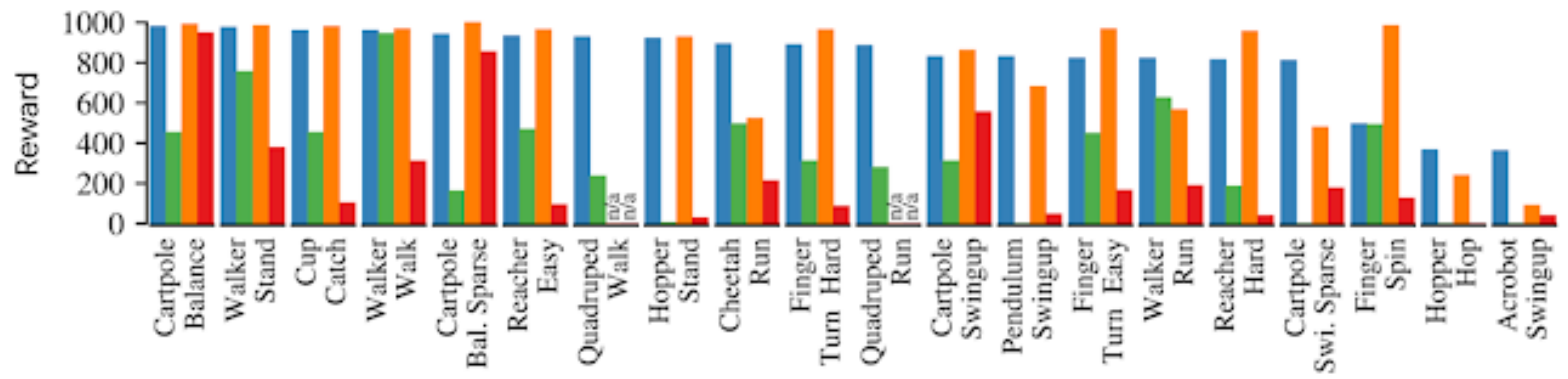


DREAMER: Results



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

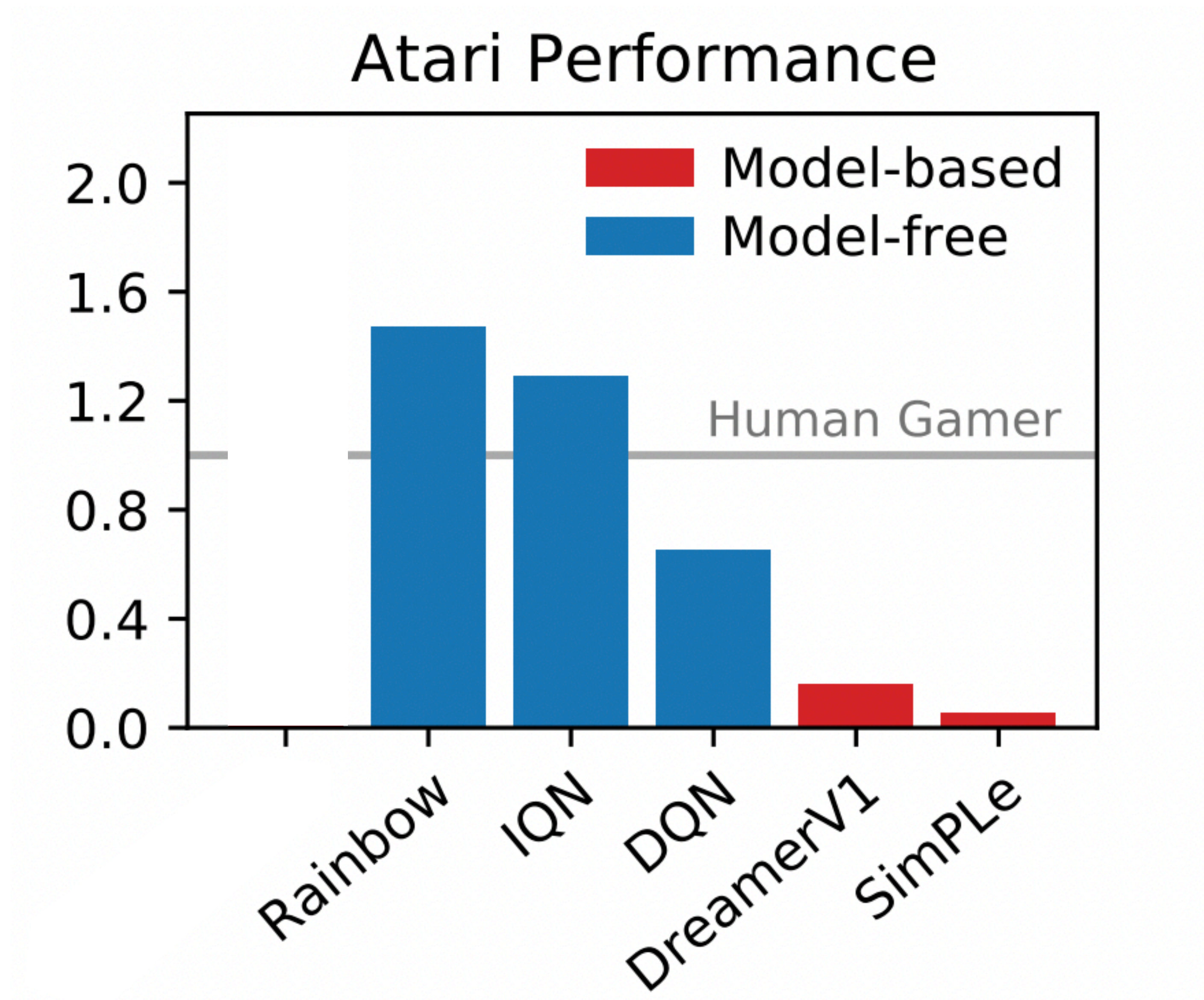
Model-free { D4PG (786) A3C (243)
 23 days of interaction



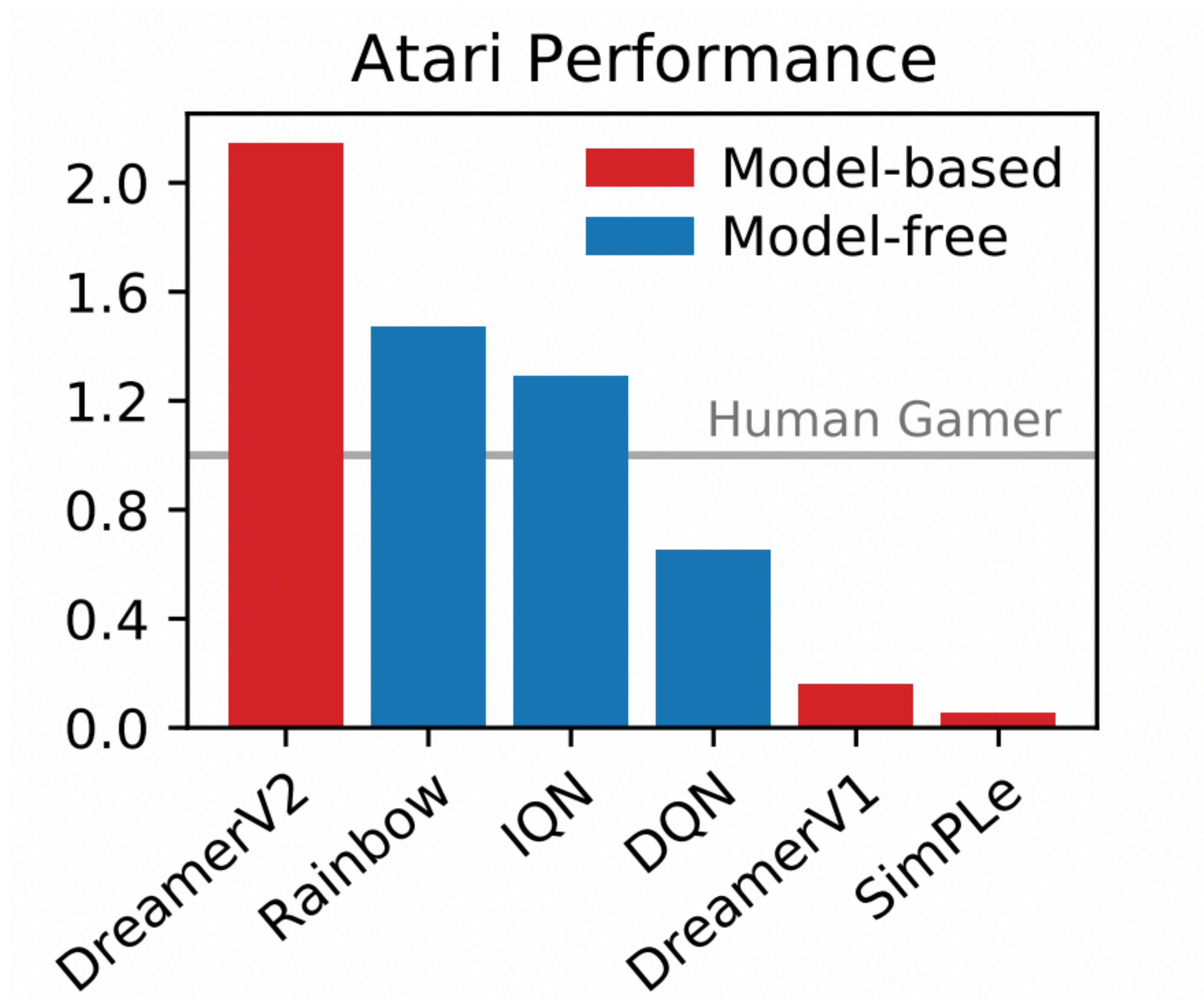
Are we done?



Atari was hard for Model Based RL

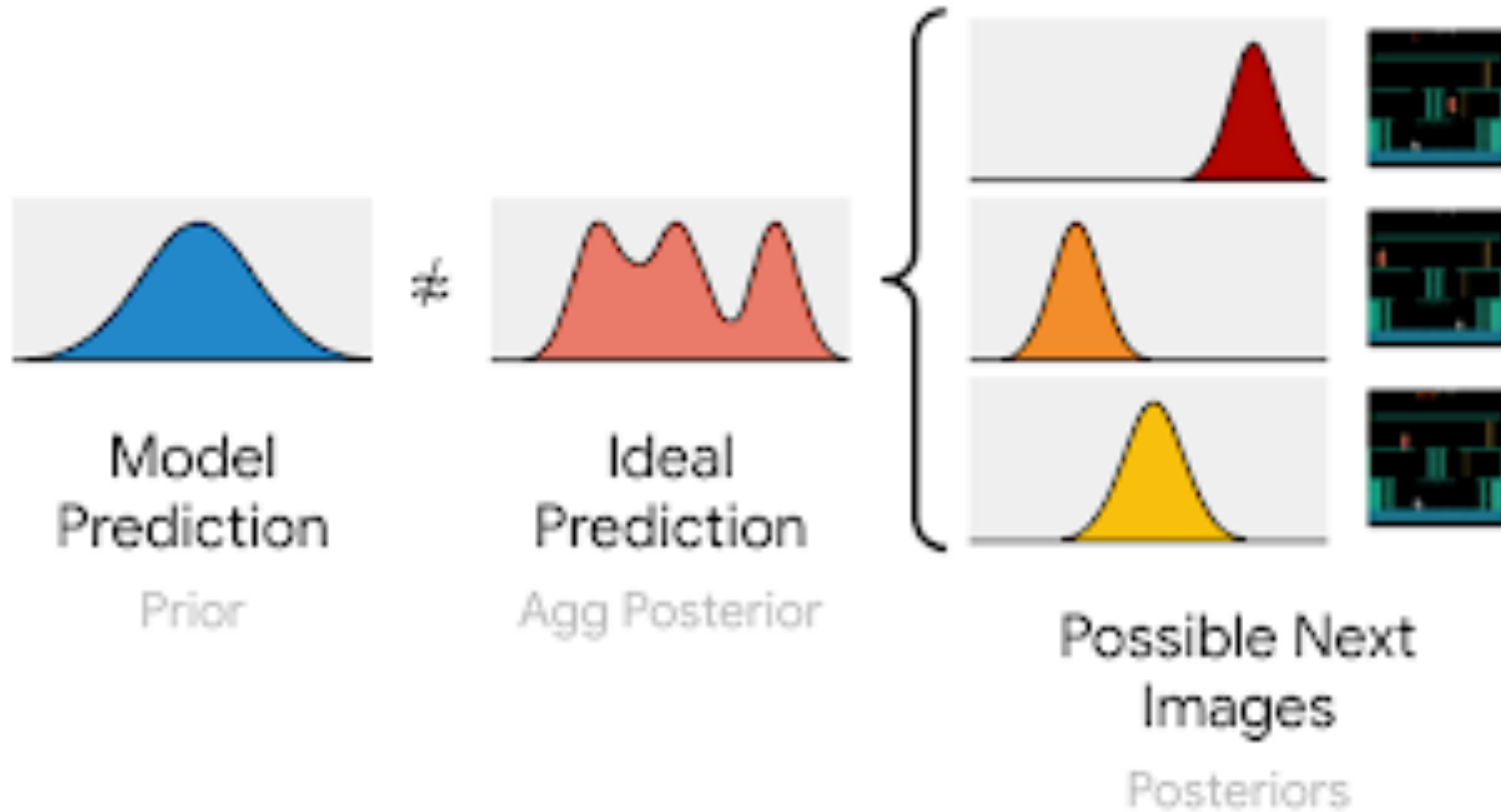


DreamerV2 beats all model free!

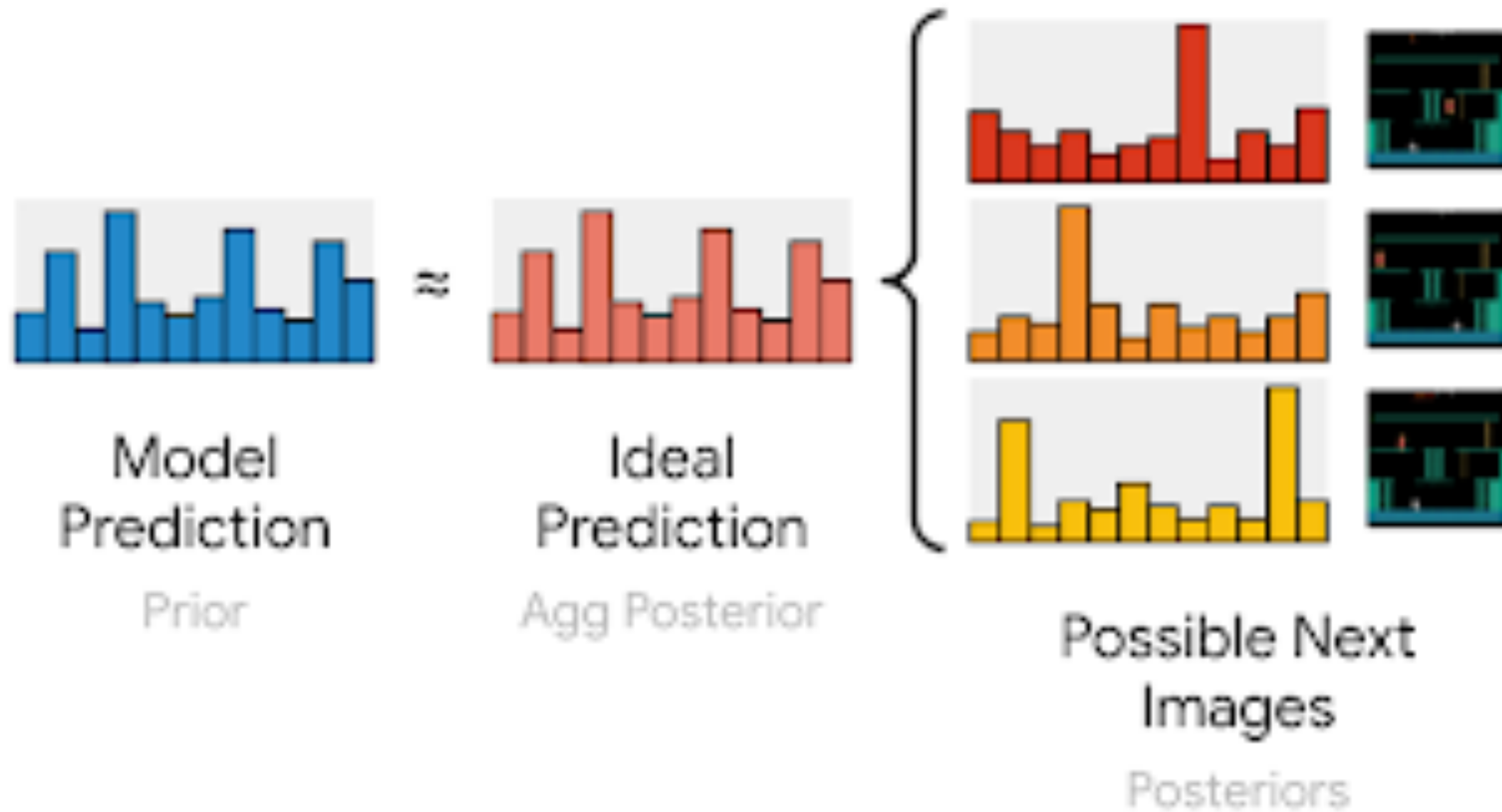


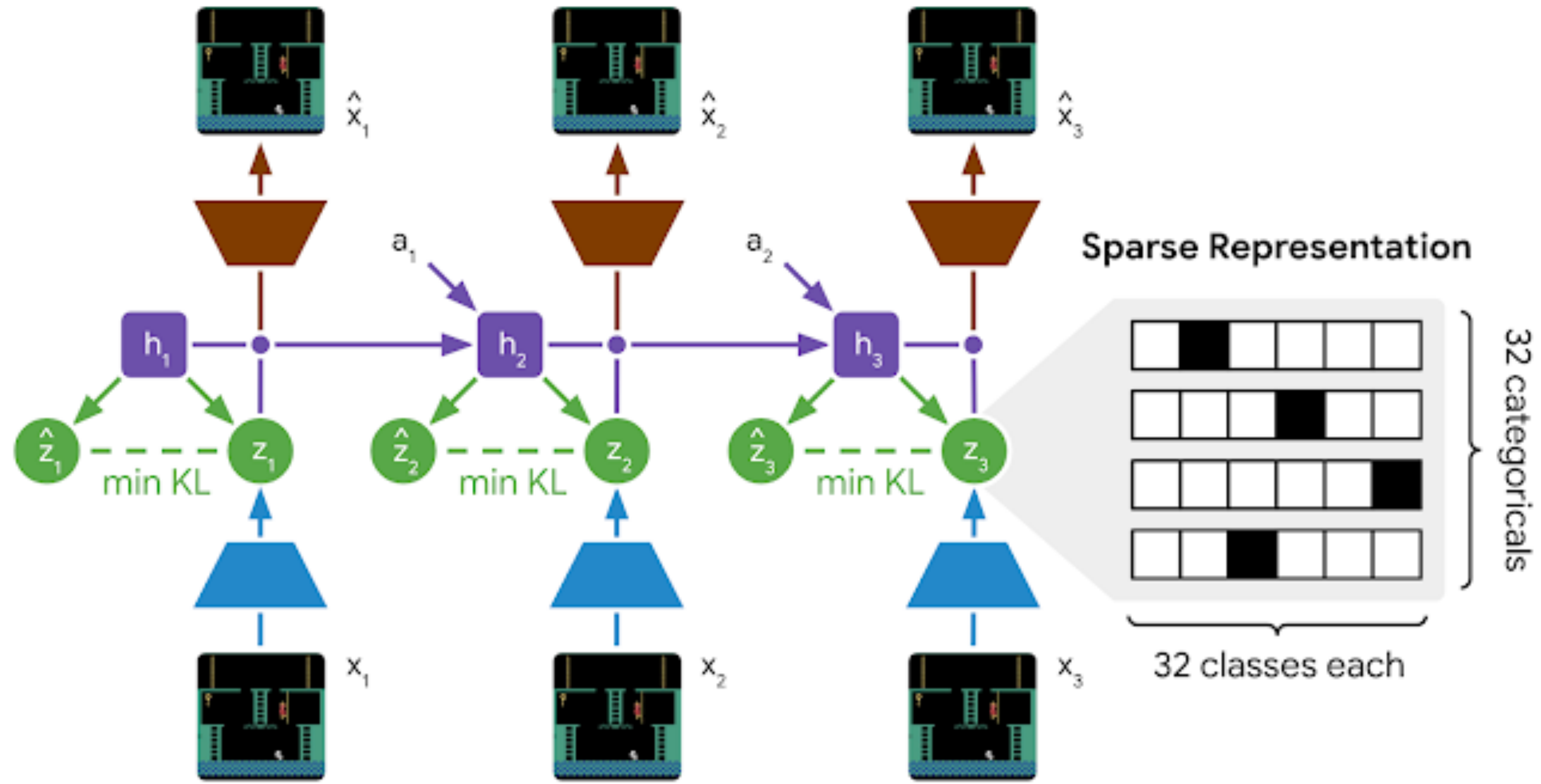
Problem: Dreamer V1
predicts a single mode of
dynamics

Dreamer V1 predicts single mode dynamics



Idea: Predict multiple discrete modes!

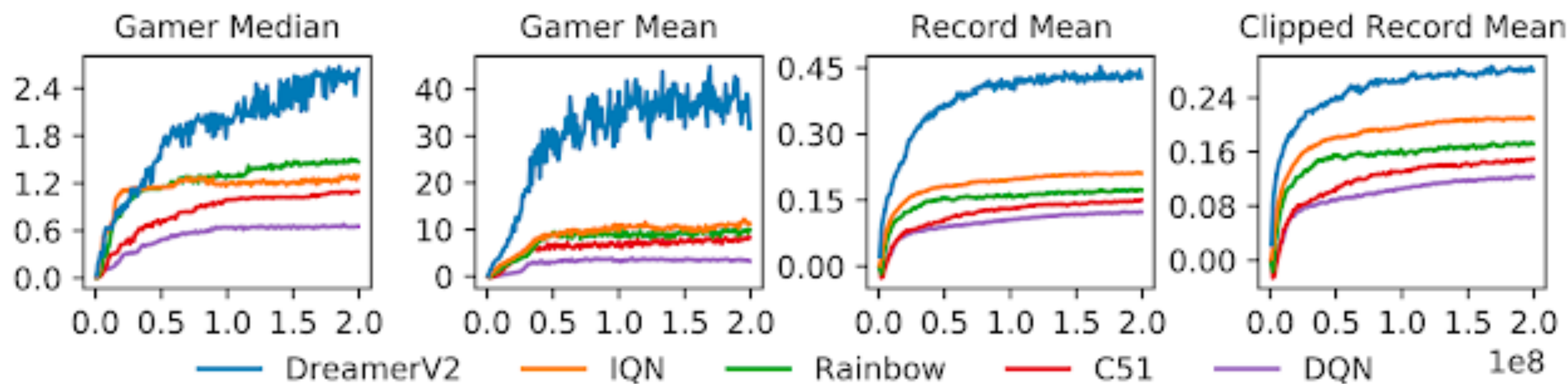




True



Model



Are we done?



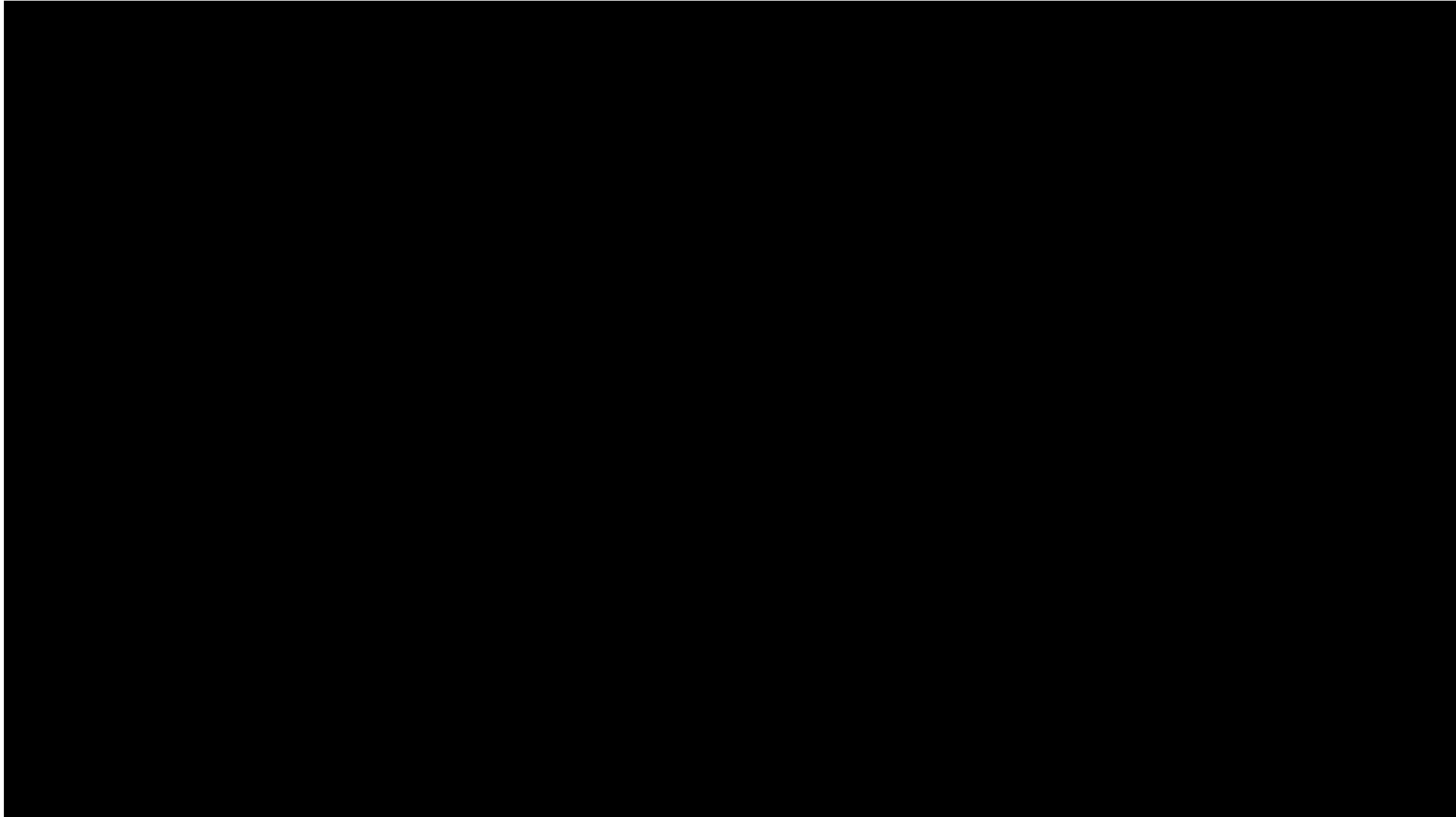
Mastering Diverse Domains through World Models

Danijar Hafner,^{1,2} Jurgis Pasukonis,¹ Jimmy Ba,² Timothy Lillicrap¹

¹DeepMind ²University of Toronto

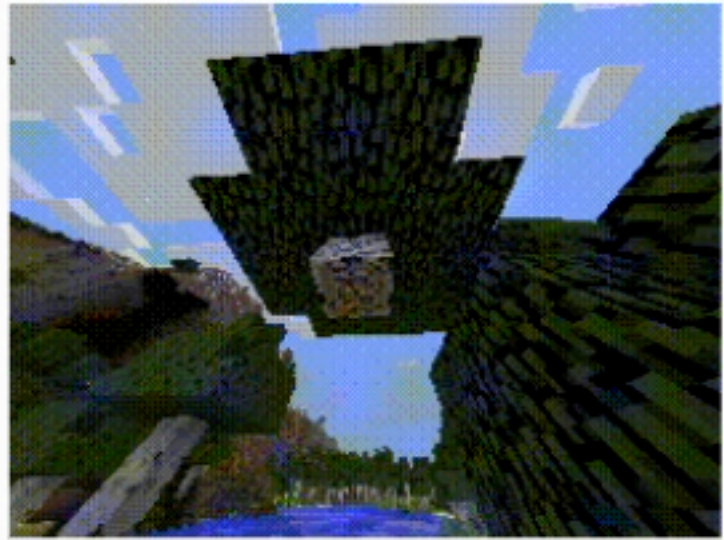


MineRL Diamond Challenge



MineRL Diamond Challenge

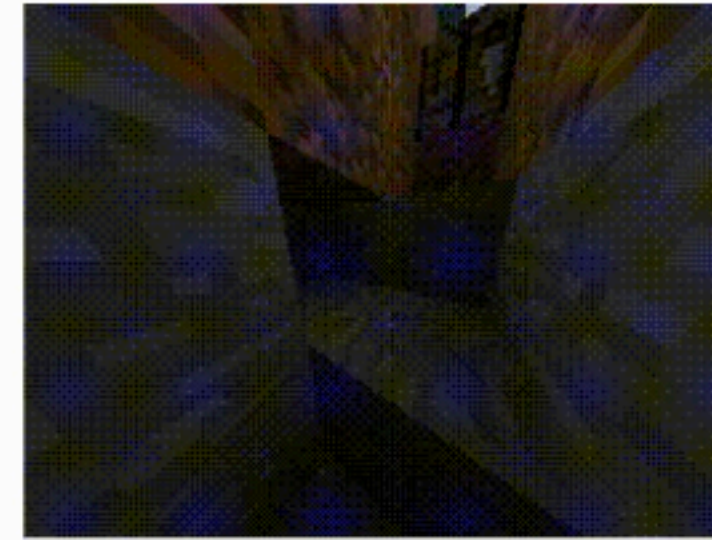
Gather Wood



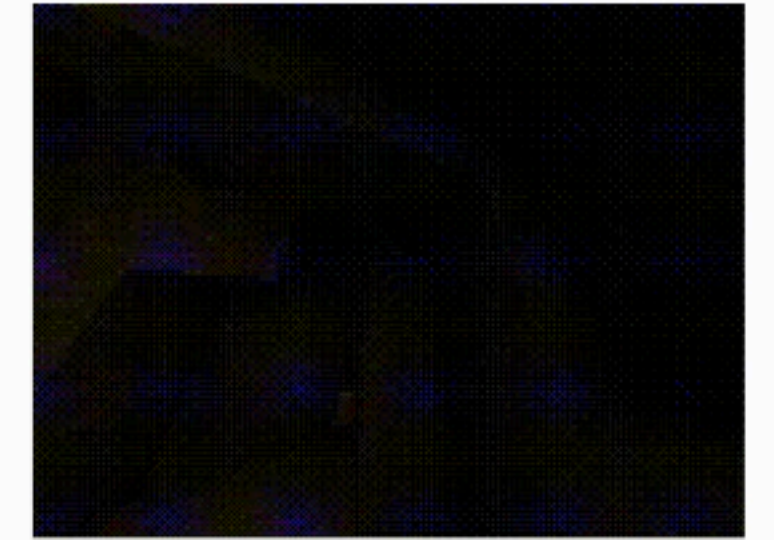
Create Wood Pickaxe



Mine Stone and Create Stone Pickaxe



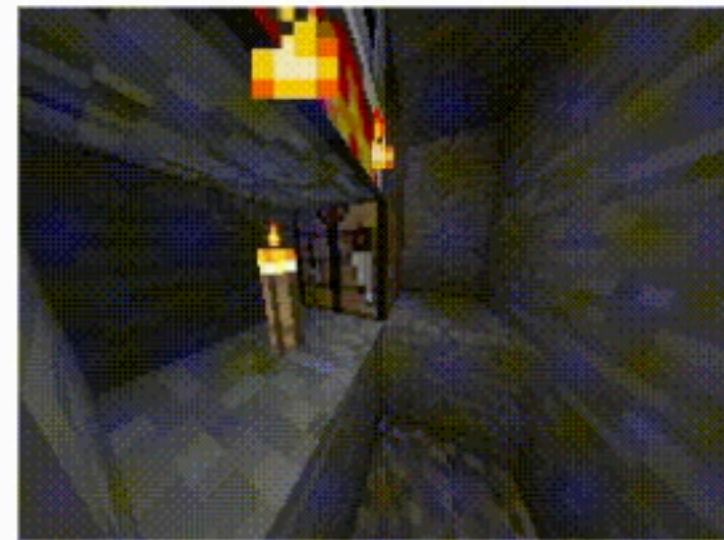
Mine Iron Ore



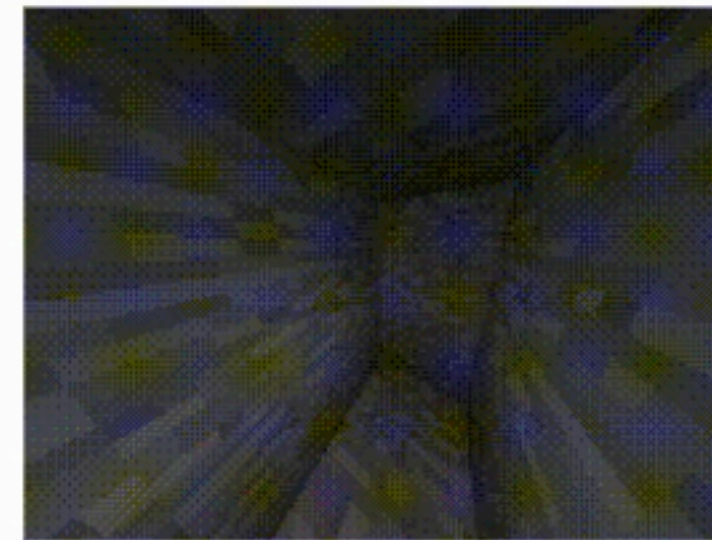
Create Furnace



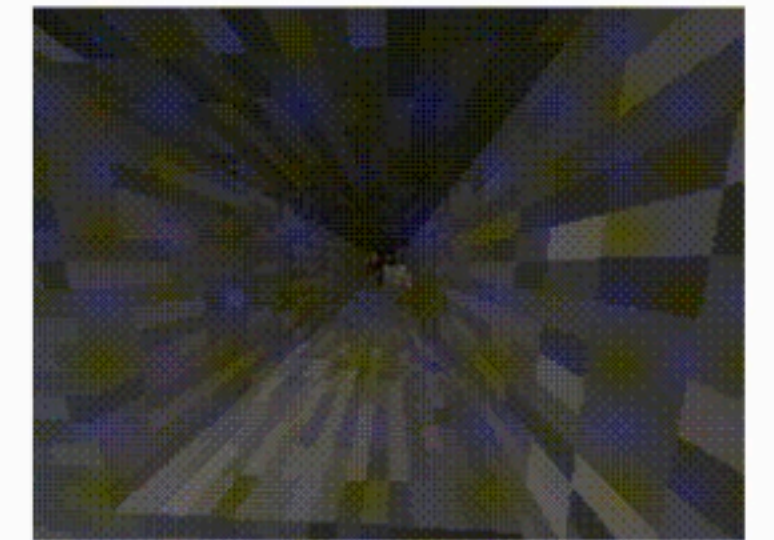
Smelt Iron and Create Iron Pickaxe



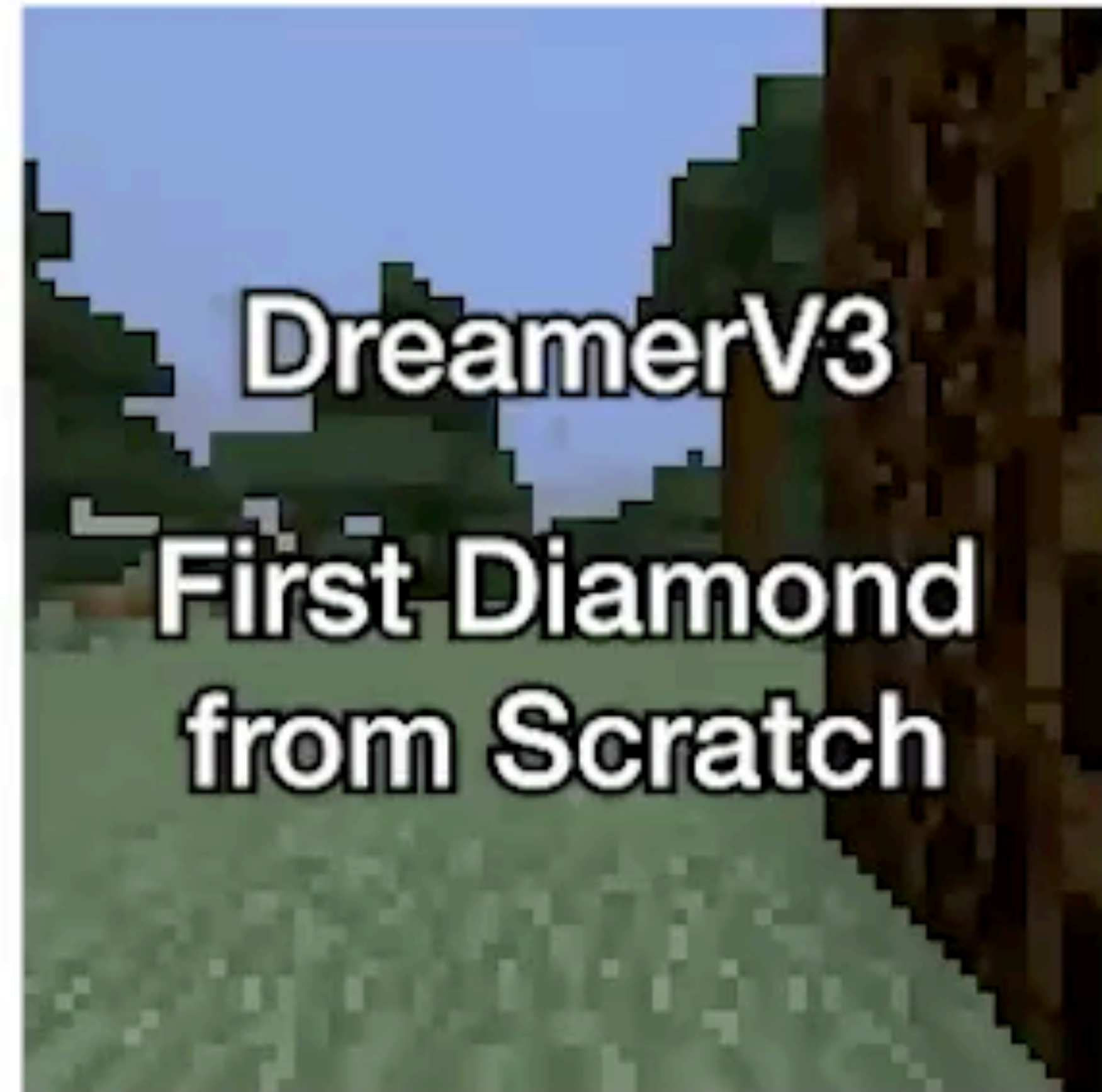
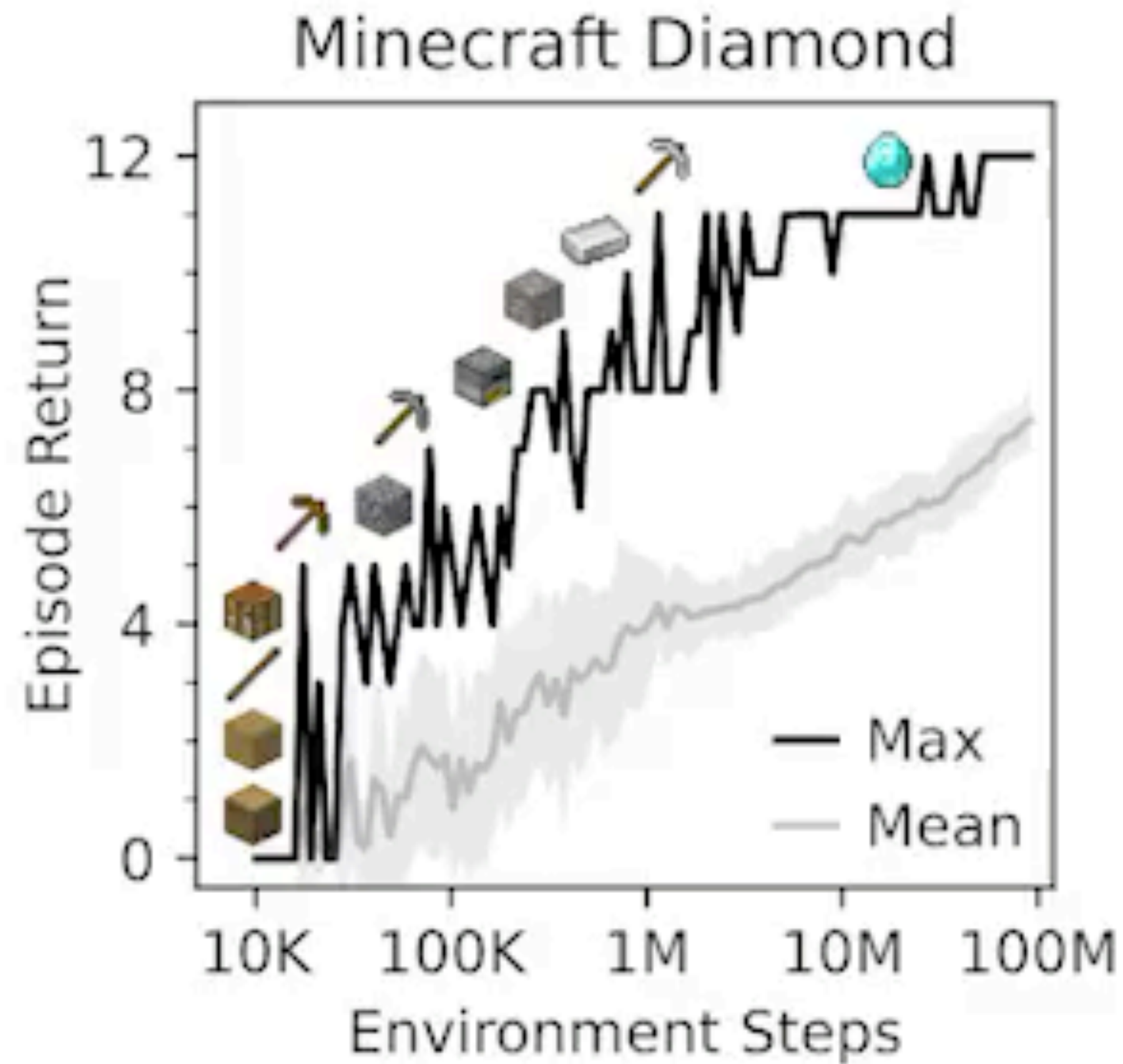
Search



Mine Diamond

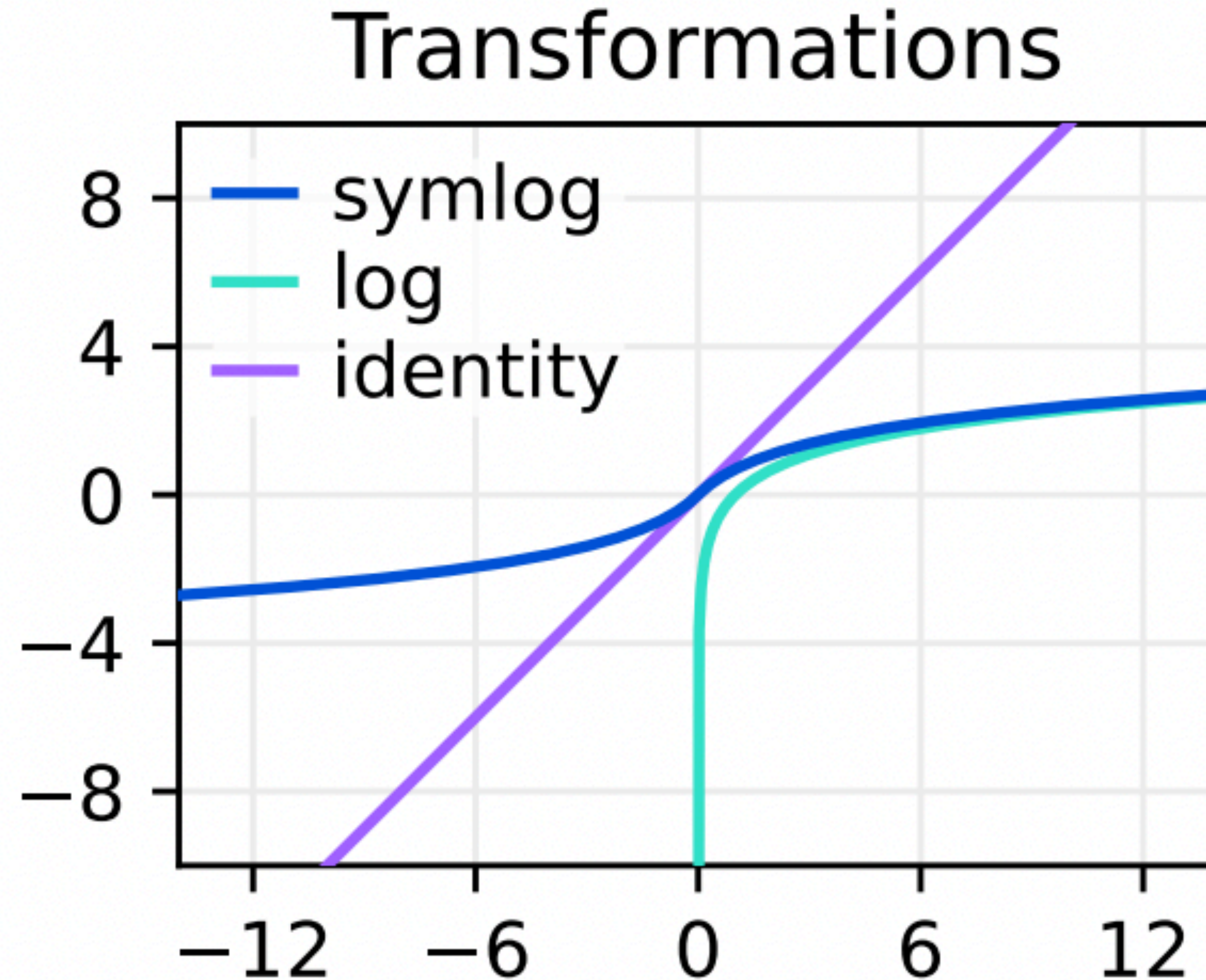


DreamerV3 solved this task!

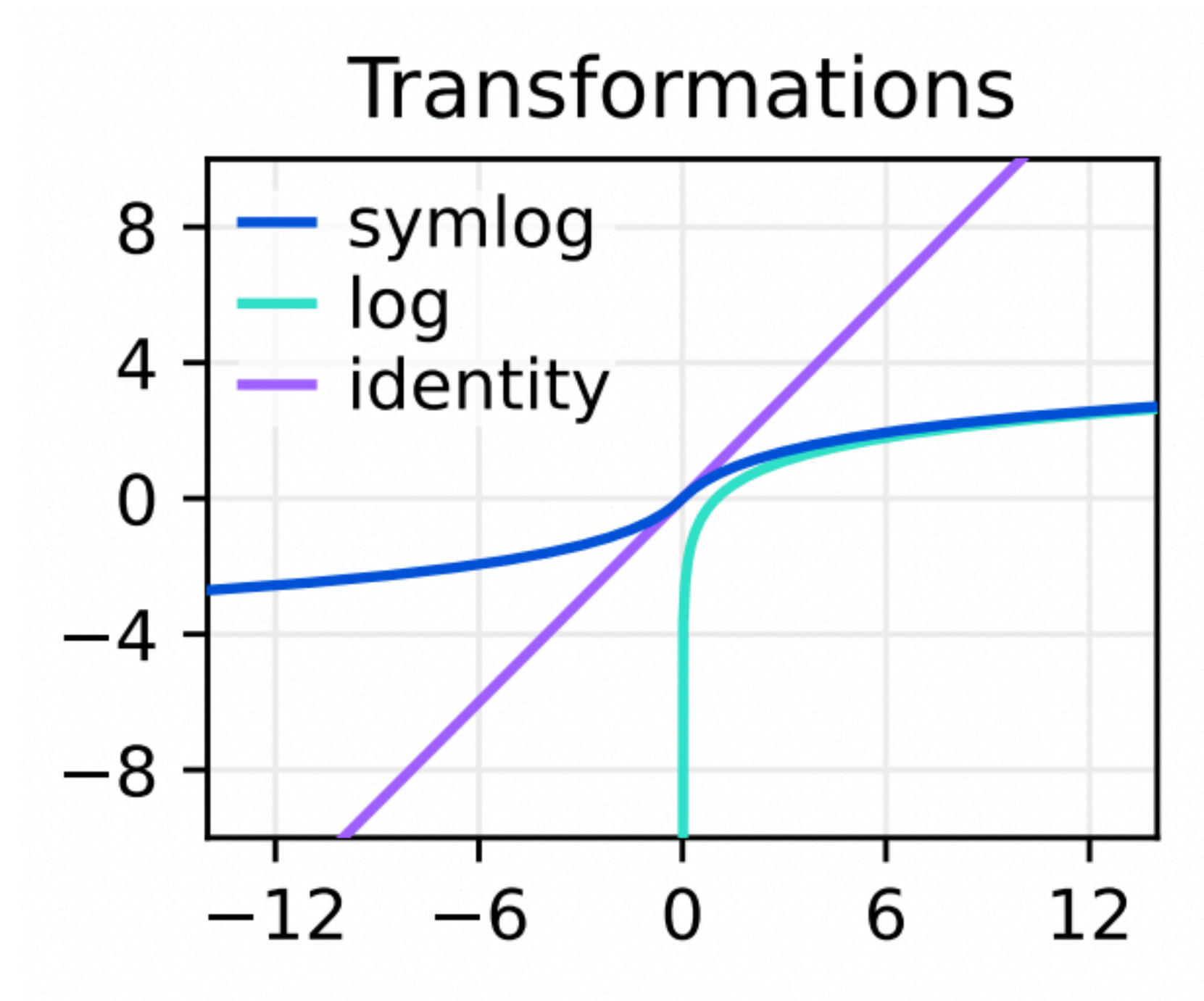


Problem: Scale of rewards,
values vary wildly across
domains

Solution: “Squash” predictions with symlog function



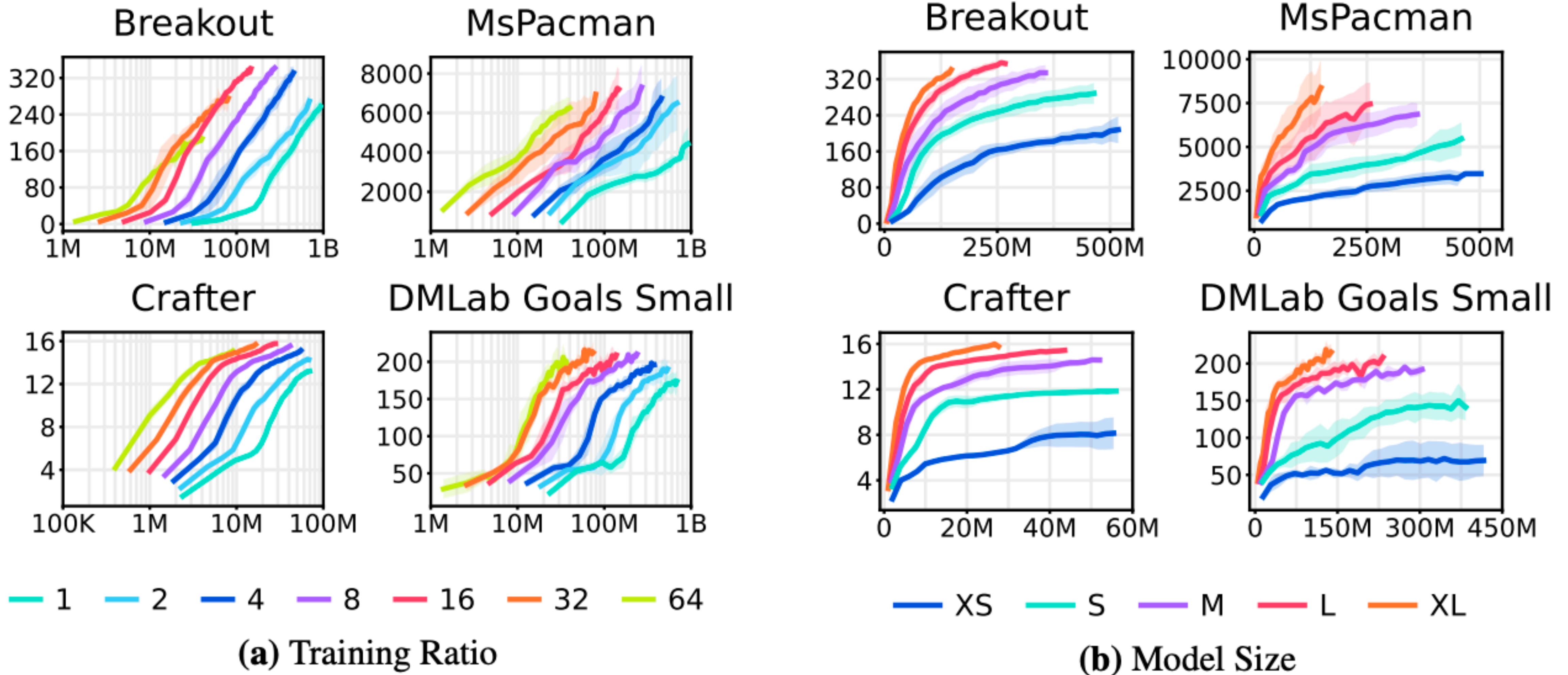
Solution: “Squash” predictions with symlog function



$$\mathcal{L}(\theta) \doteq \frac{1}{2} (f(x, \theta) - \text{symlog}(y))^2 \quad \hat{y} \doteq \text{symexp}(f(x, \theta))$$

$$\text{symlog}(x) \doteq \text{sign}(x) \ln(|x| + 1) \quad \text{symexp}(x) \doteq \text{sign}(x) (\exp(|x|) - 1)$$

DreamerV3 scales really well!



tl;dr

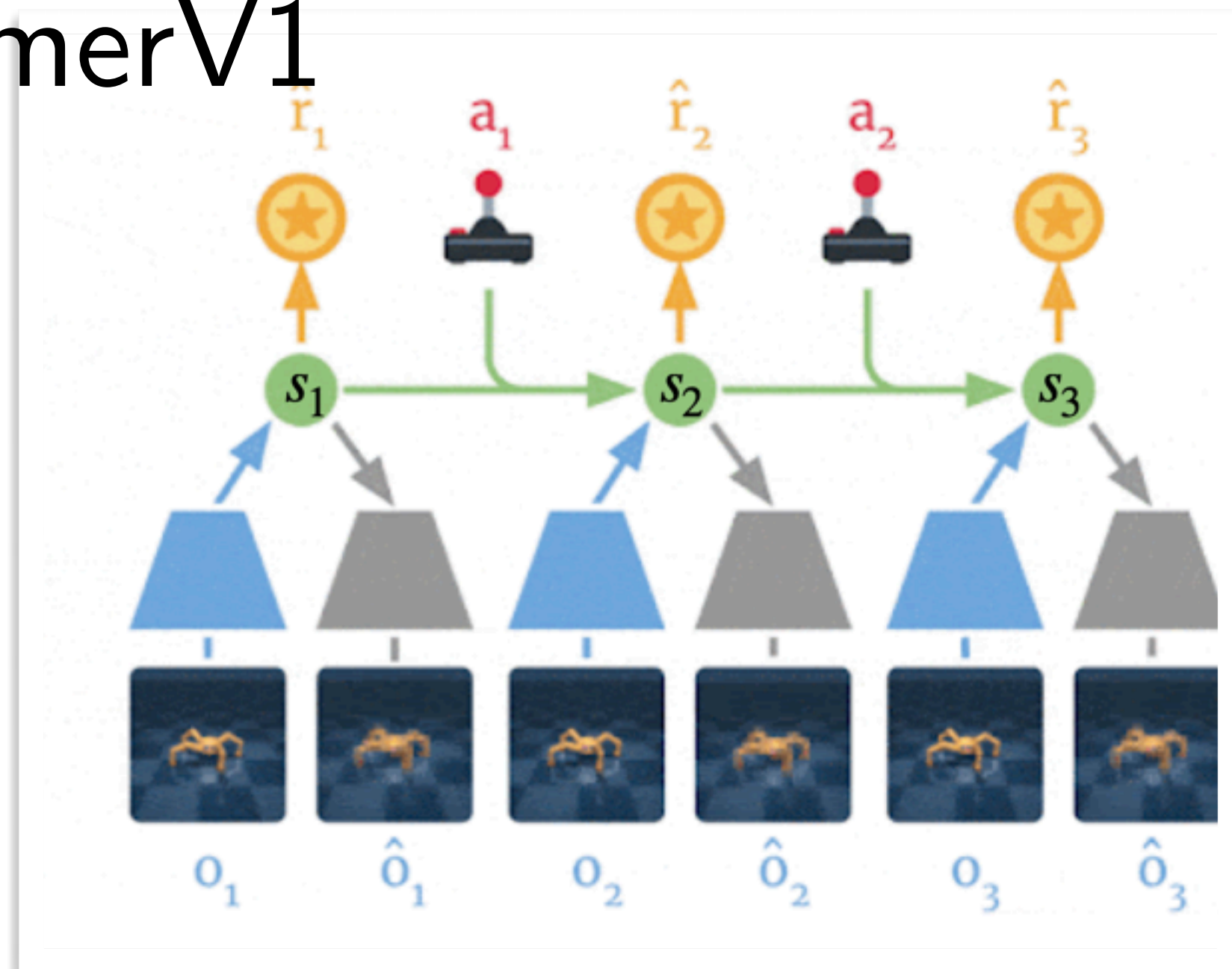
Challenges with learning complex models

Challenge 1: Partial Observability

Challenge 2: Planning with Complex Dynamics

16

DreamerV1



Extensions (V2, V3)

