

Generative World Models: The Dreamer Models

Sanjiban Choudhury



Cornell Bowers CIS
Computer Science

The story so far ...

Robots have to act in the world

Hence, we learned various algorithms for
decision making

But we assumed that we can observe the “state”

The story so far ...

But in the real world, no one tells you the
“state”

All you see are observations

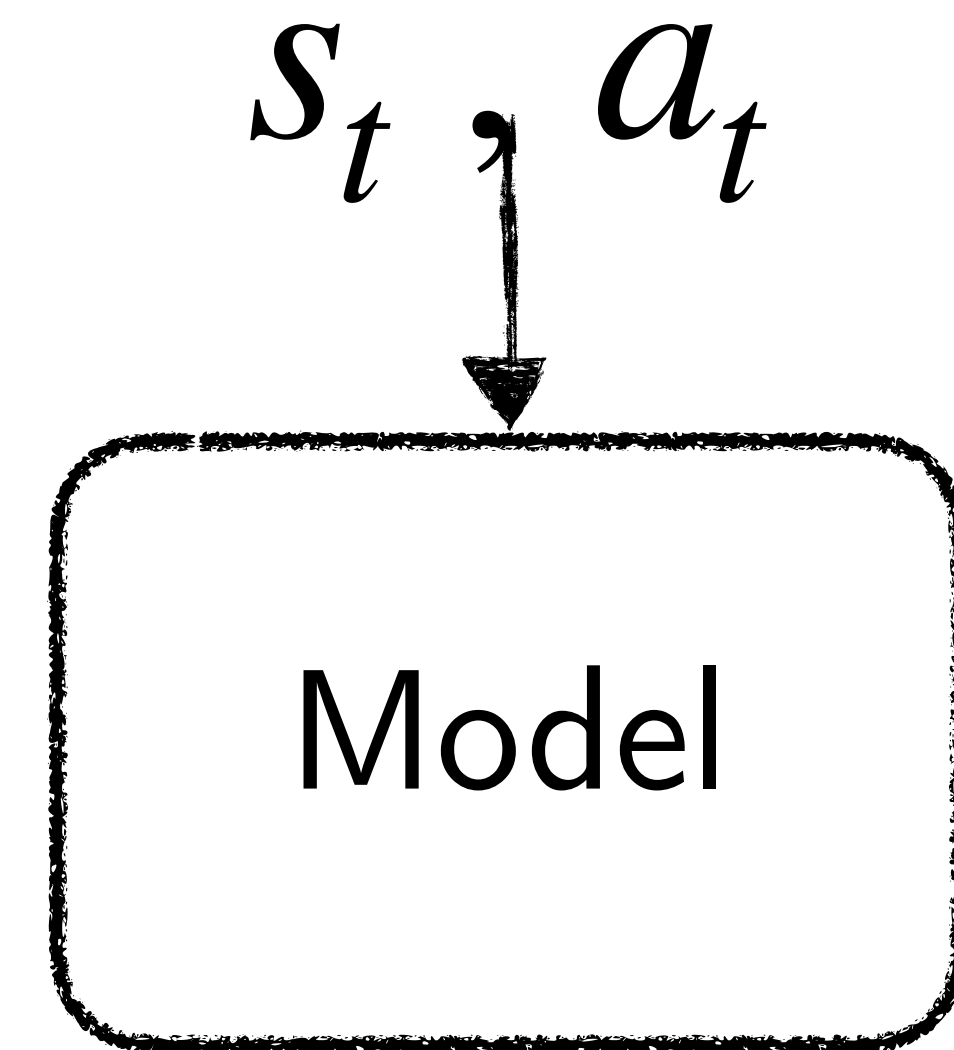
How do we learn from observations?

The story so far ...

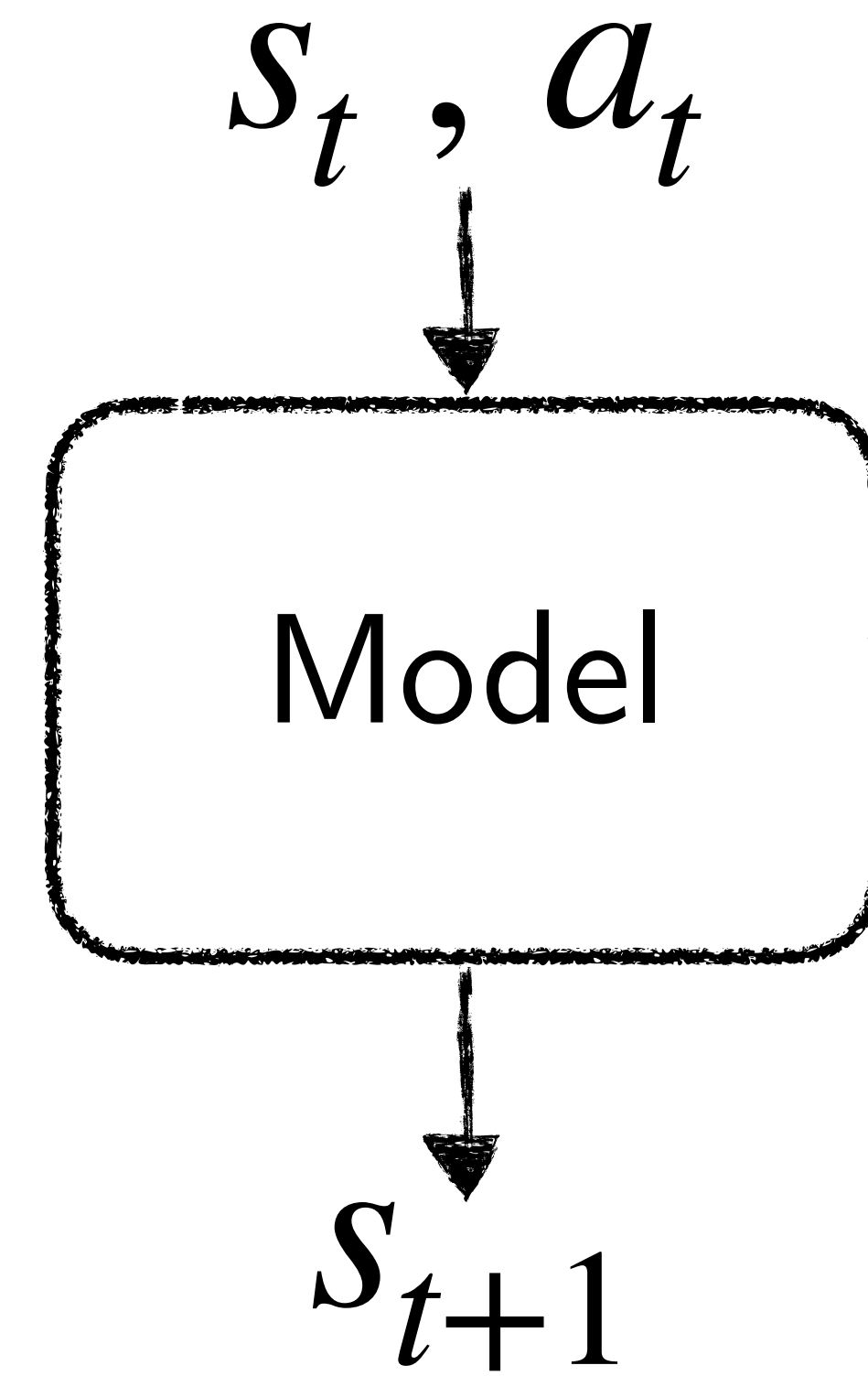
Our focus in this and future lectures
will turn to learning representations

Models.

What is a model?



What is a model?



What is a model?

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

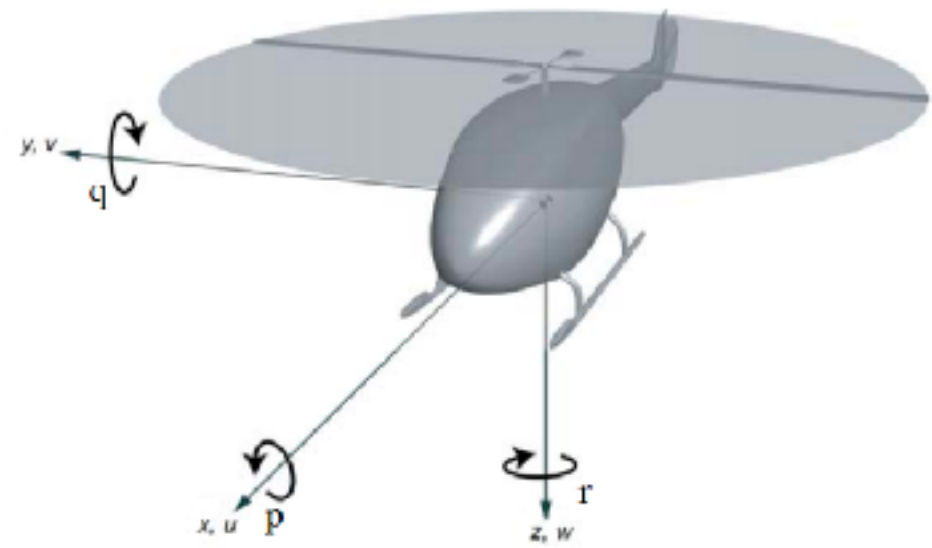
Learning Models

Models: From Simple to Complex

Simple

Complex

Models: From Simple to Complex



Physics Models

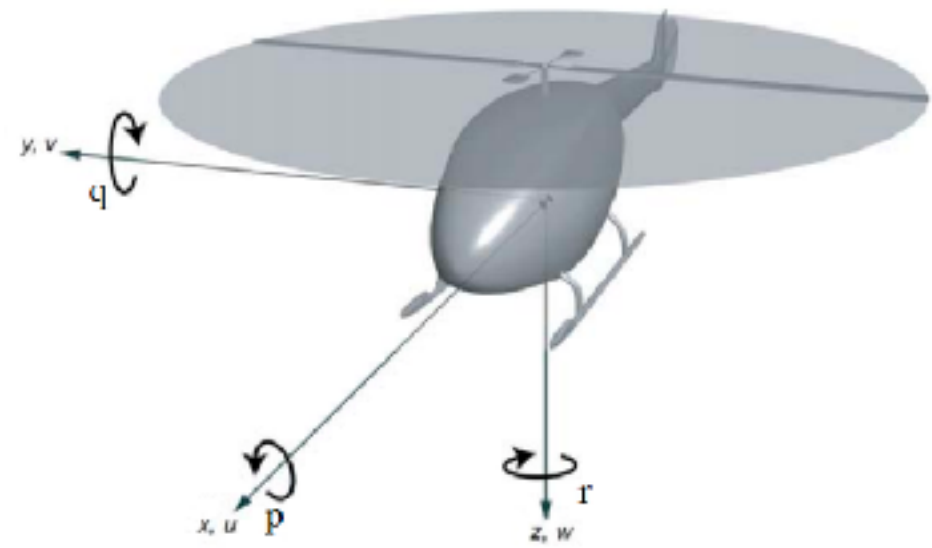
Simple

Known state

Strong prior
on dynamics



Models: From Simple to Complex

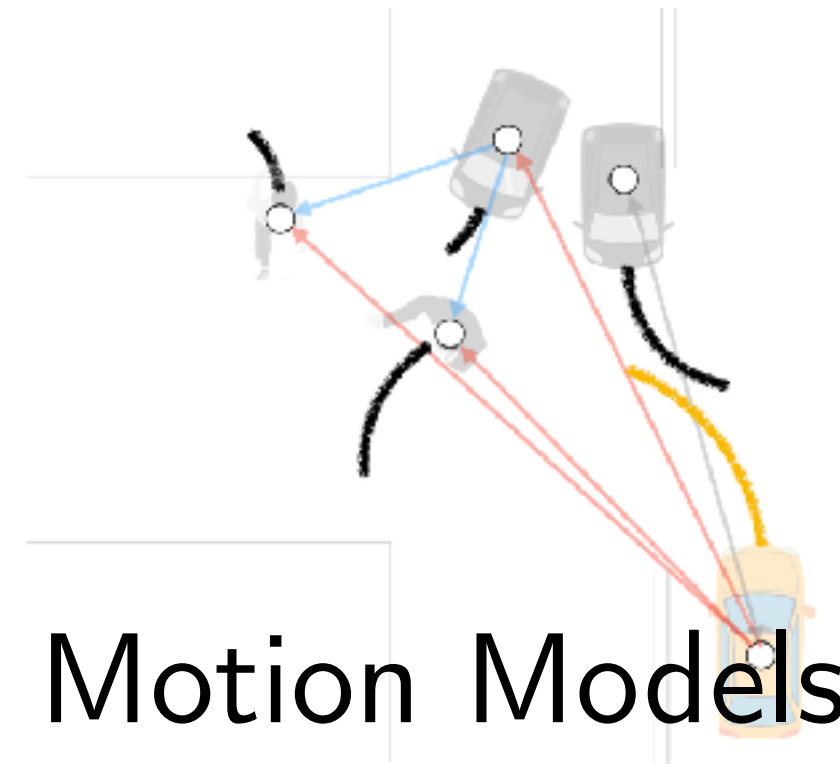


Physics Models

Simple

Known state

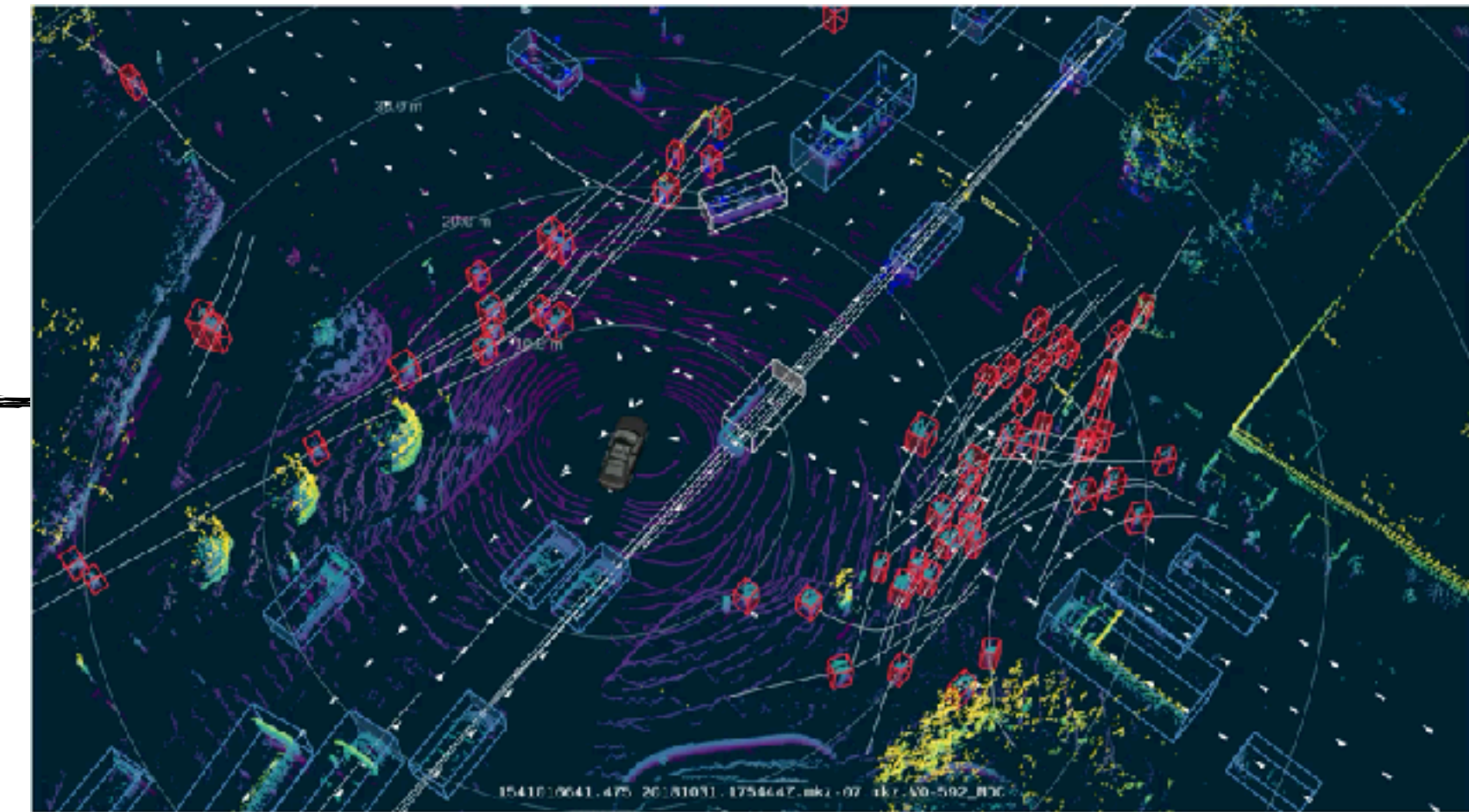
Strong prior
on dynamics



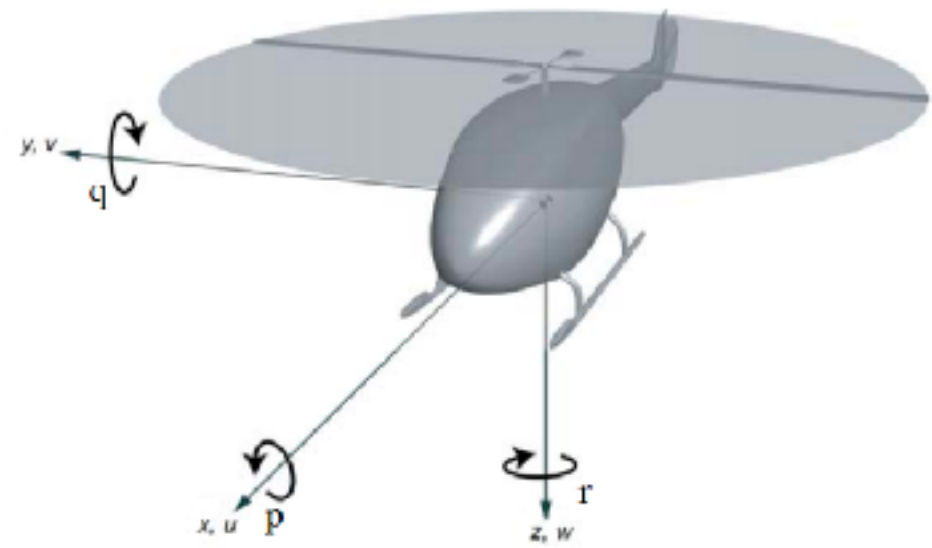
Motion Models

Known state

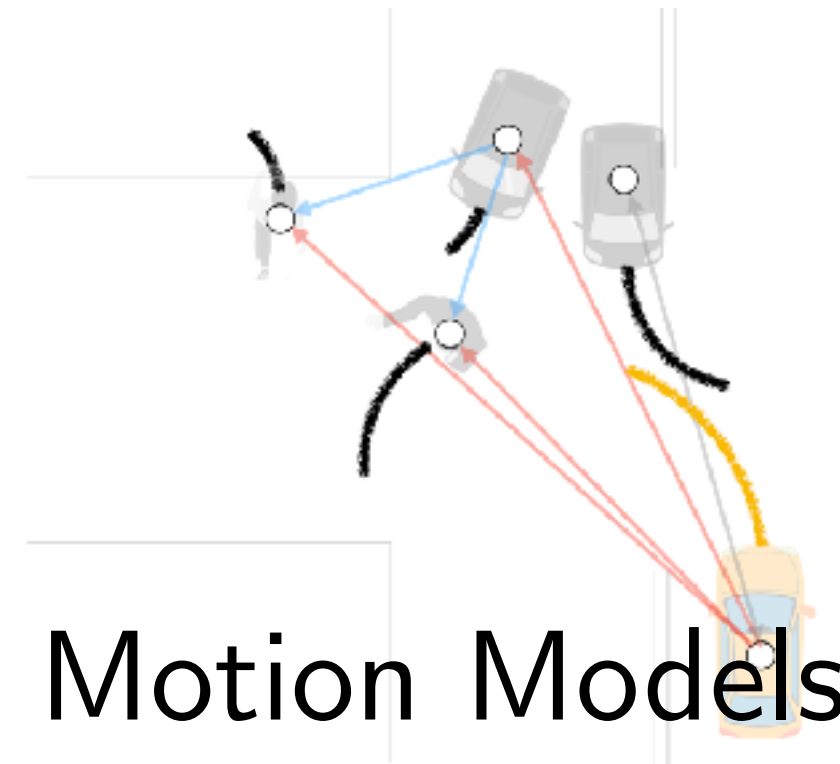
Unknown
dynamics



Models: From Simple to Complex



Physics Models



Motion Models



Open World Models

Simple

Complex

Known state

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

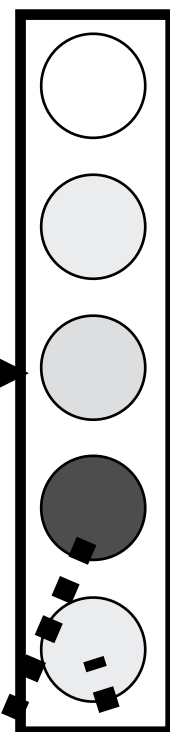
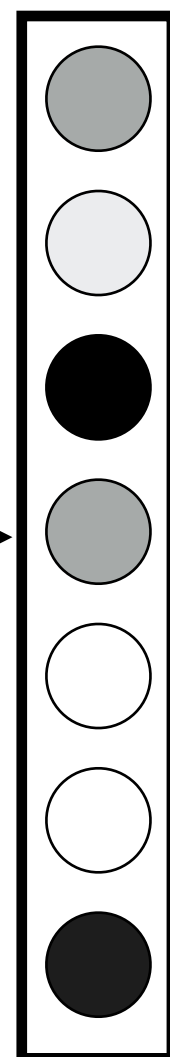
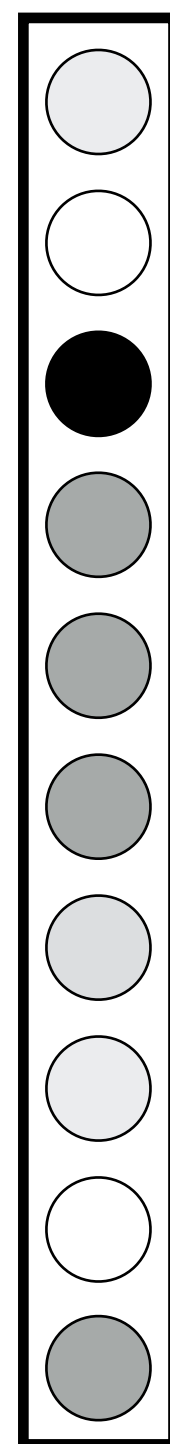
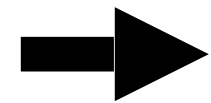
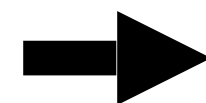
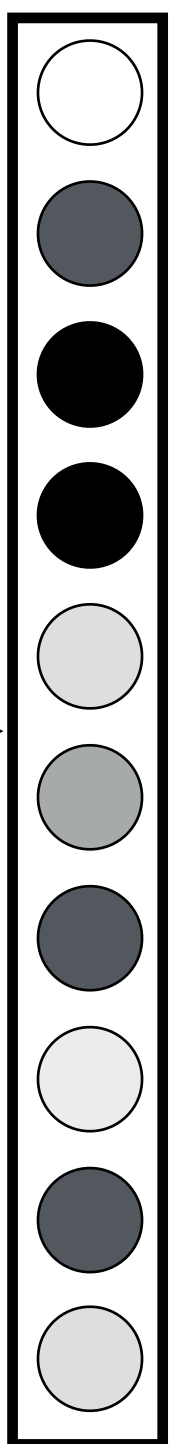
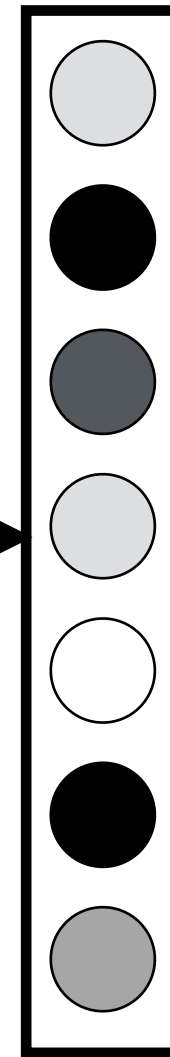
Challenge 2: Planning with complex dynamics

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

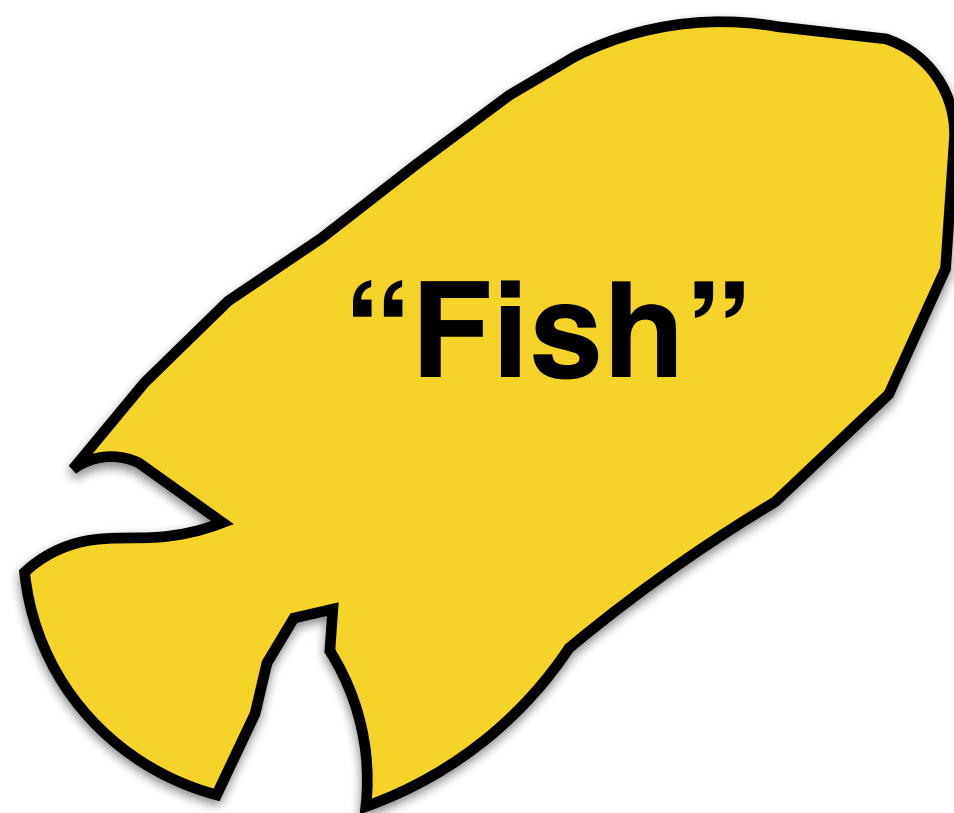
Idea: Use “auto-encoder” trick from
computer vision

\mathbf{X} 

Image

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

Reconstructed image



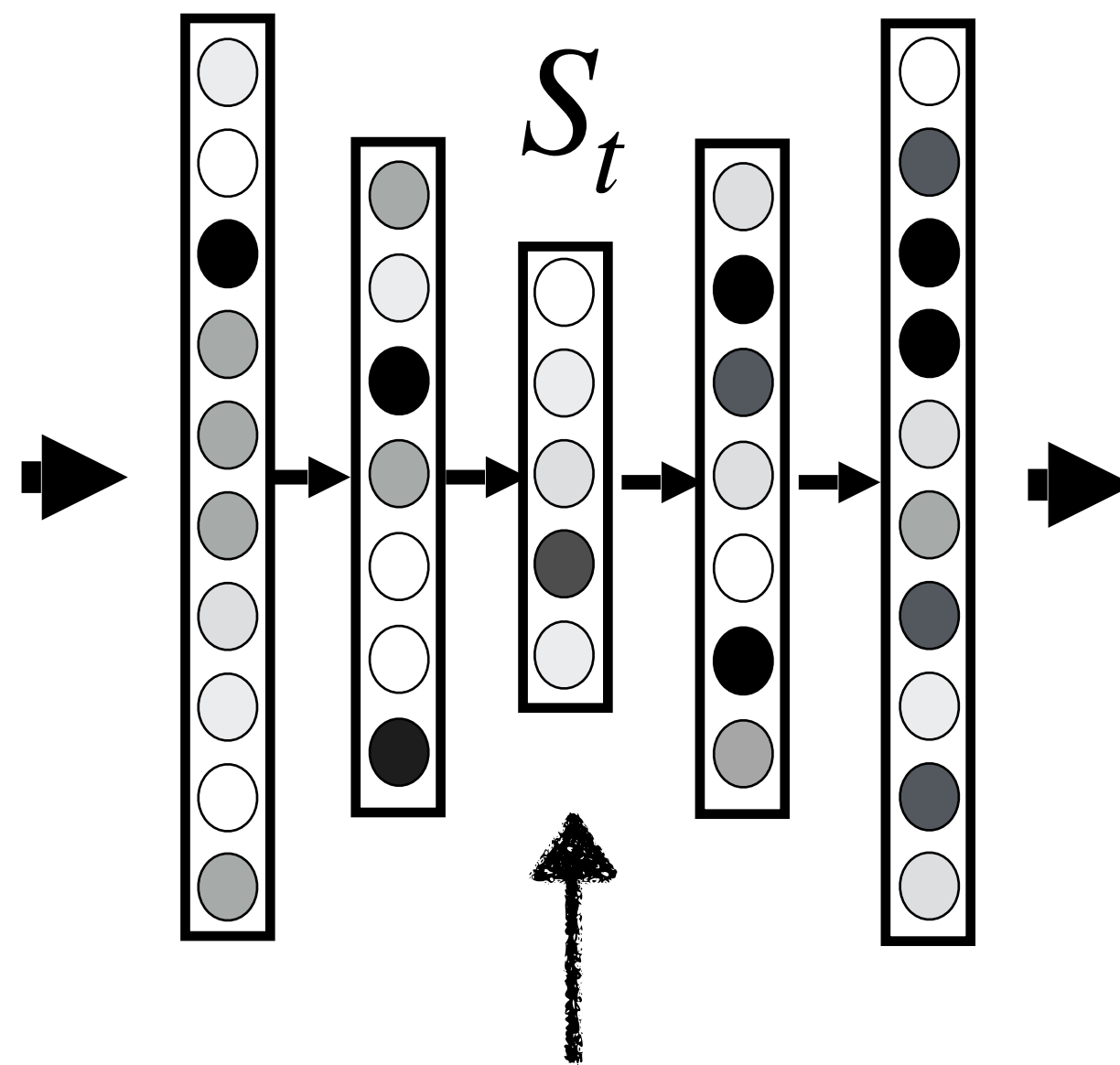
"Fish"



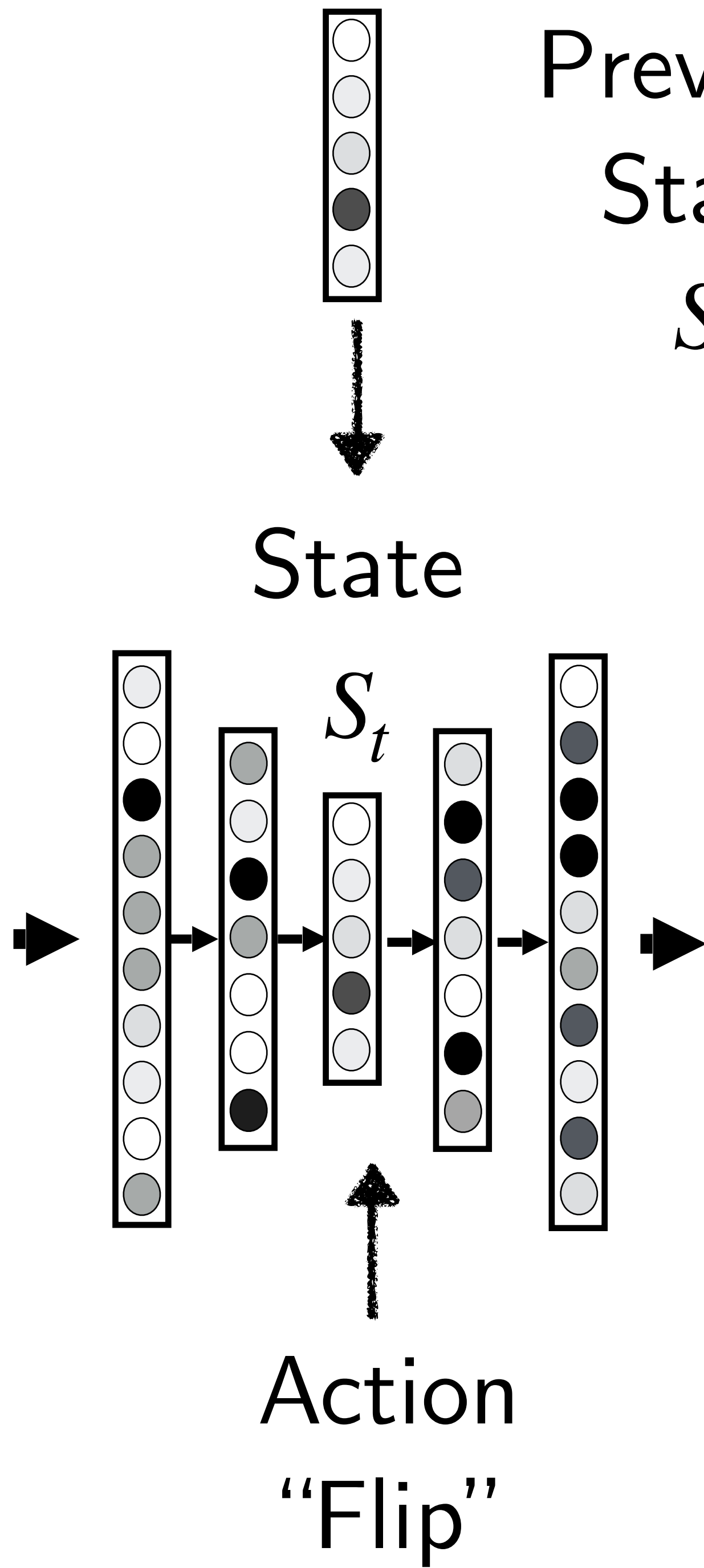
"Coral"



State

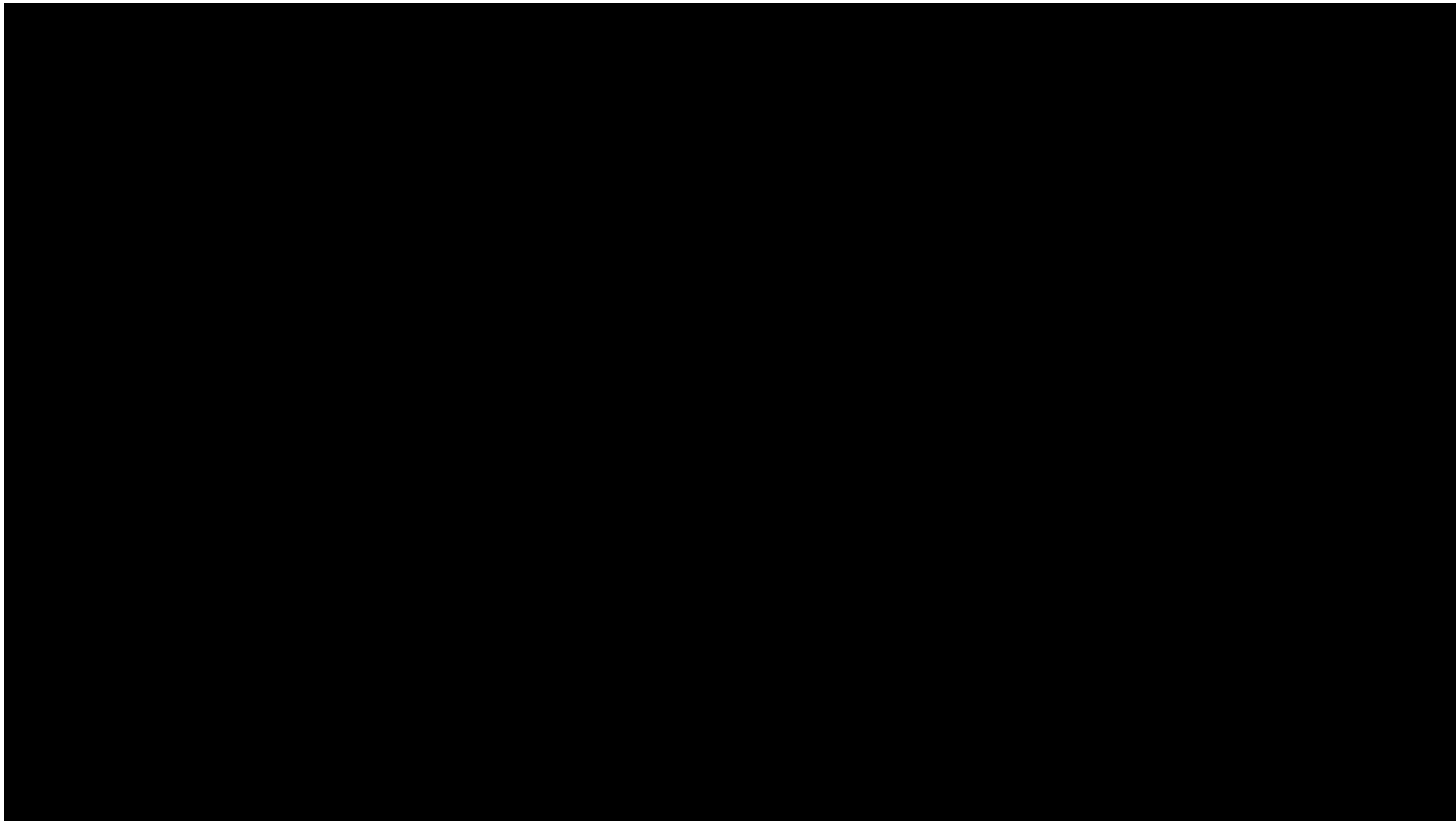


Action
"Flip"



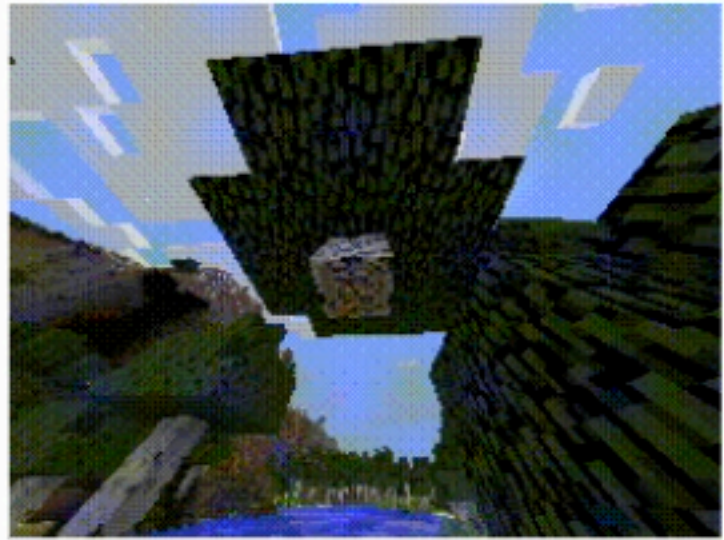
The DREAMER Algorithms

MineRL Diamond Challenge

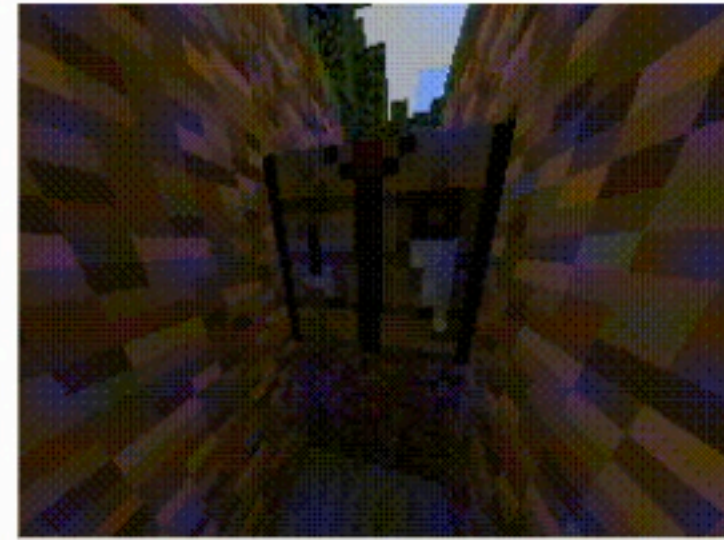


MineRL Diamond Challenge

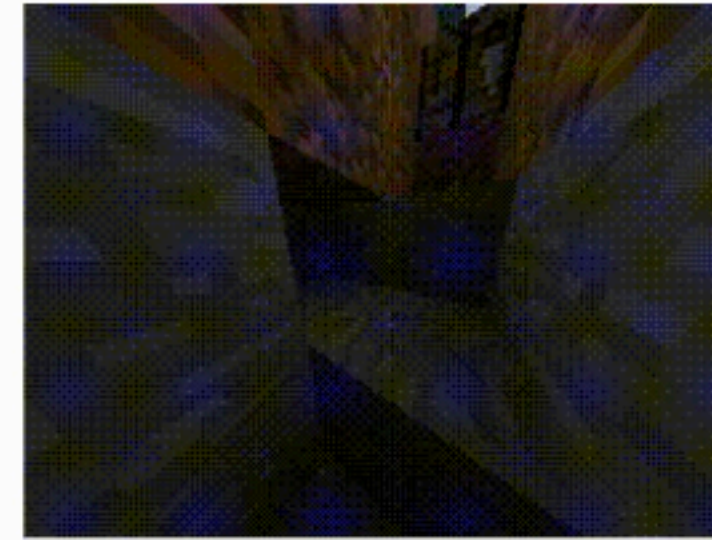
**Gather
Wood**



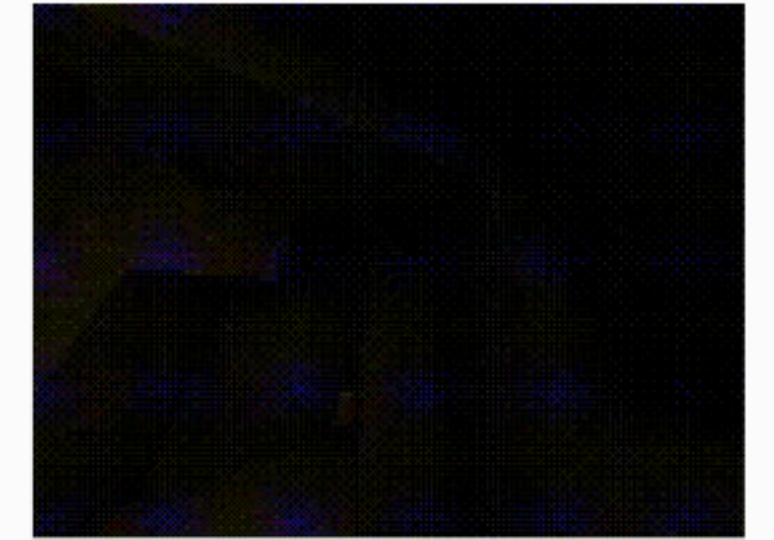
**Create
Wood Pickaxe**



**Mine Stone
and Create
Stone Pickaxe**



**Mine
Iron Ore**



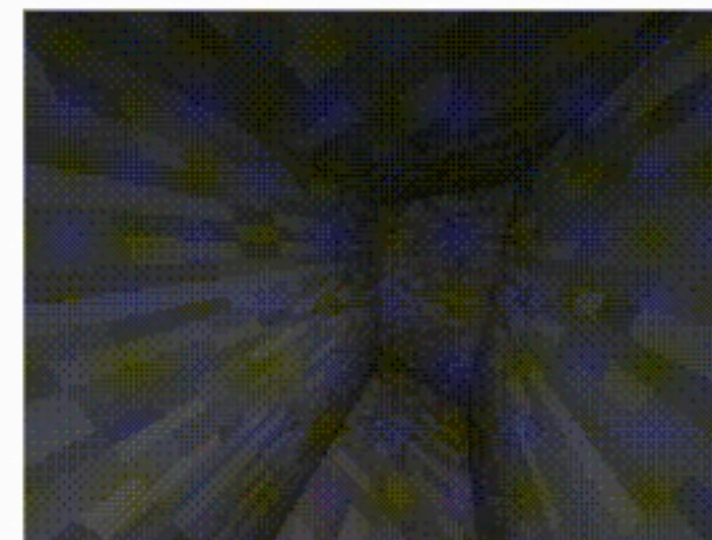
**Create
Furnace**



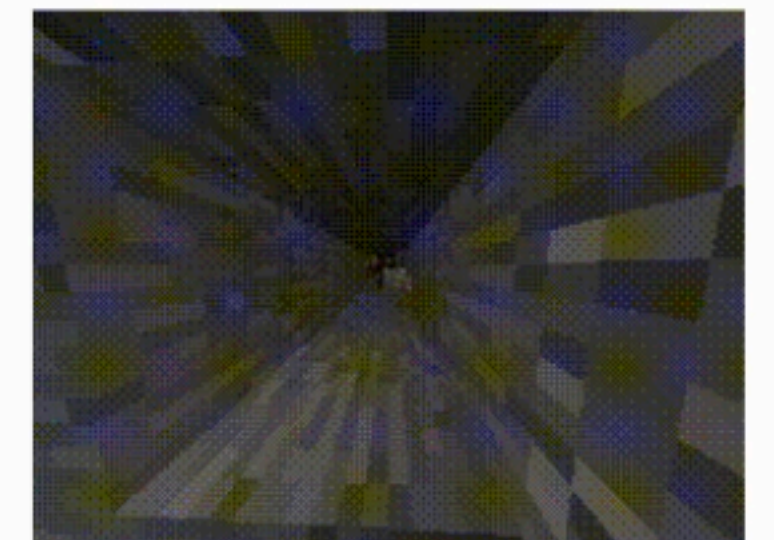
**Smelt Iron
and Create
Iron Pickaxe**



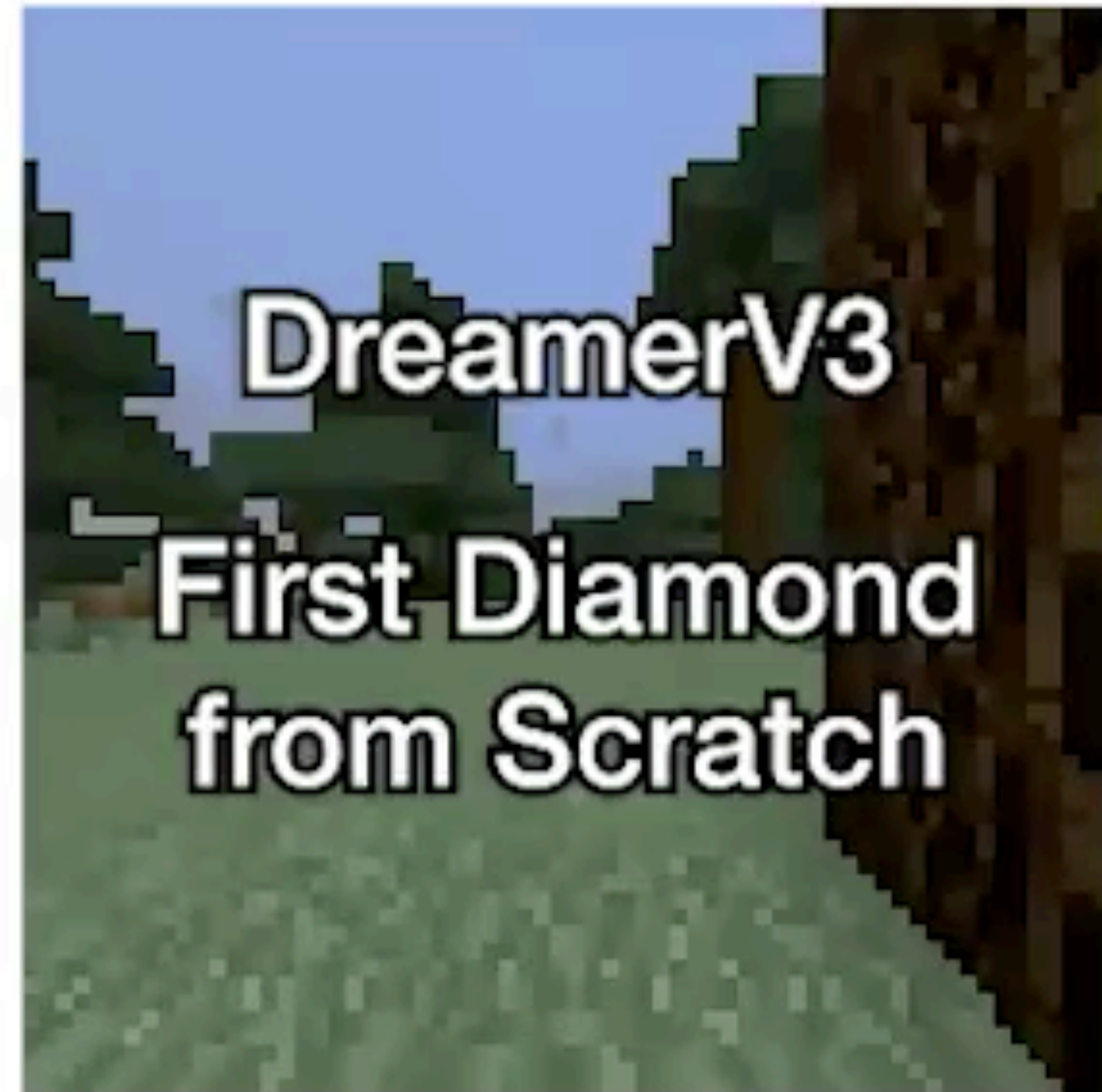
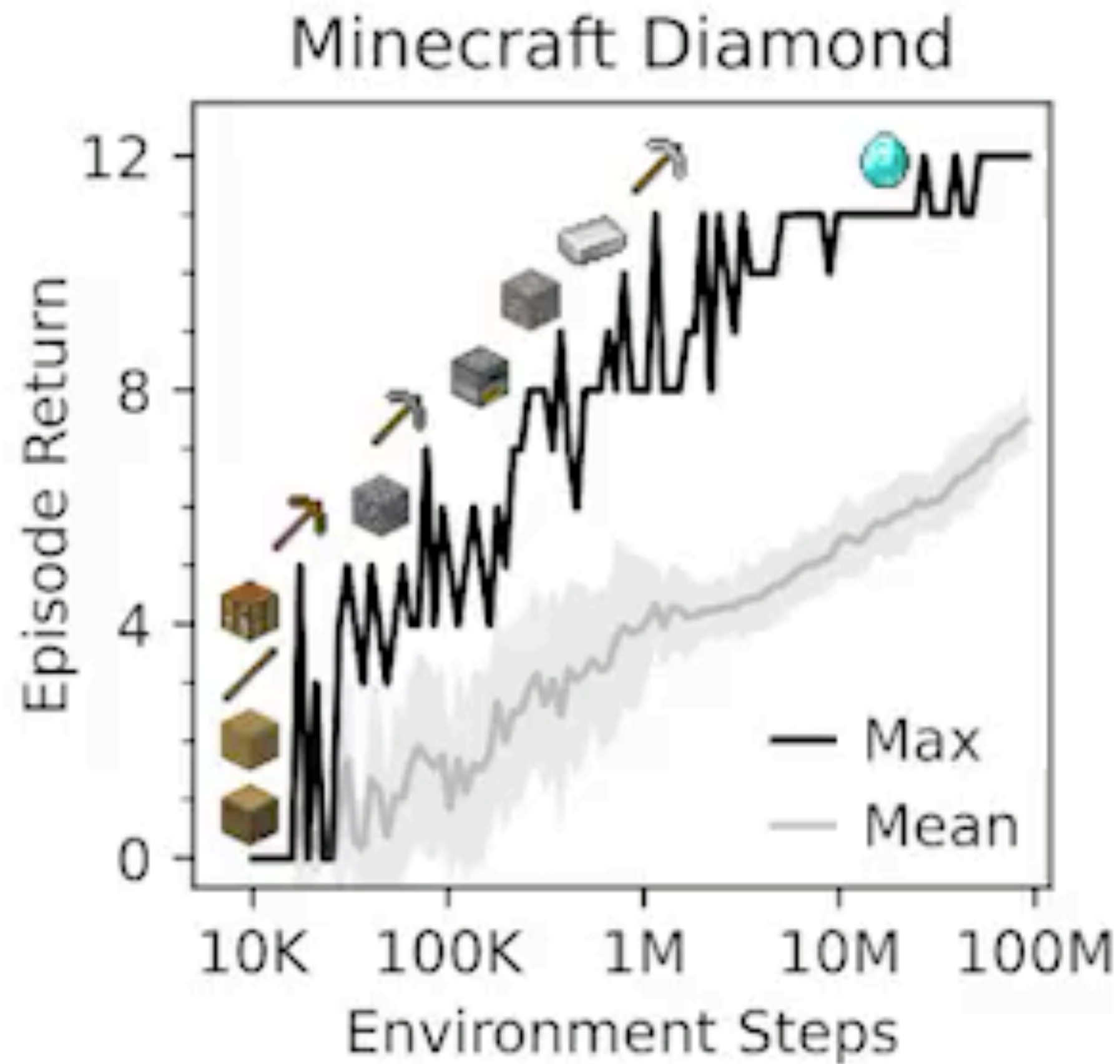
Search



**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

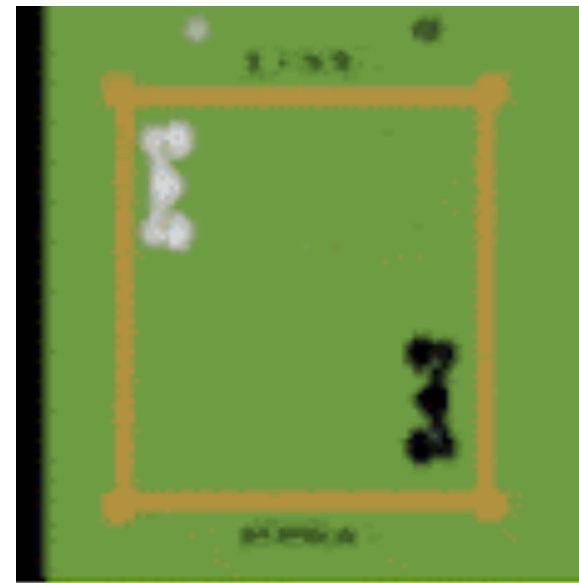
University of Toronto

Mohammad Norouzi

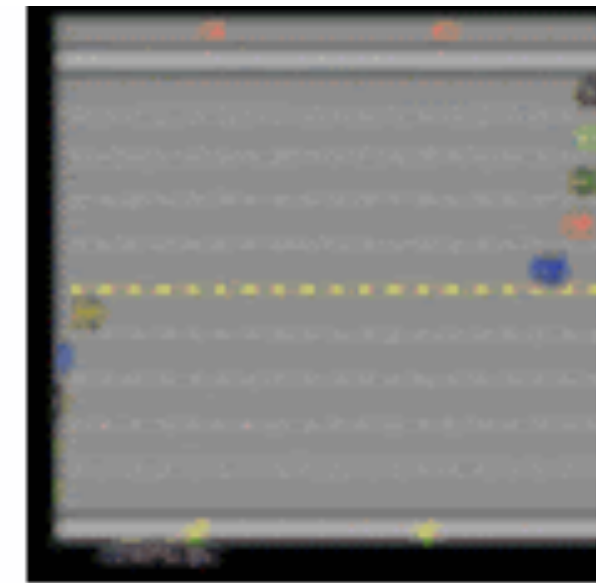
Google Brain

2020

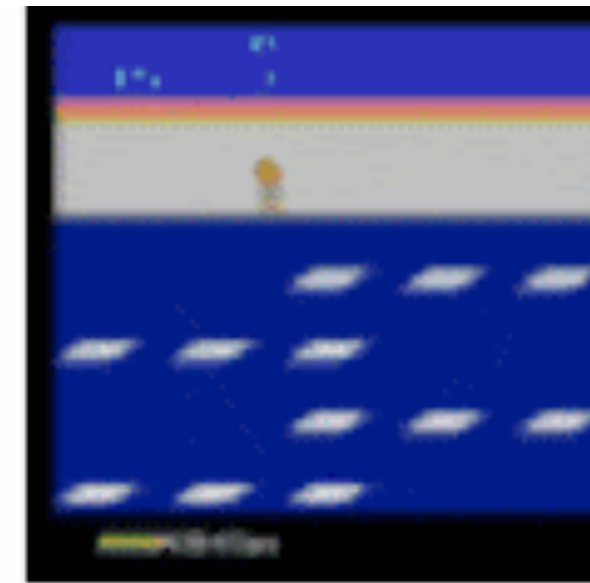
Look at the videos below



Boxing



Freeway



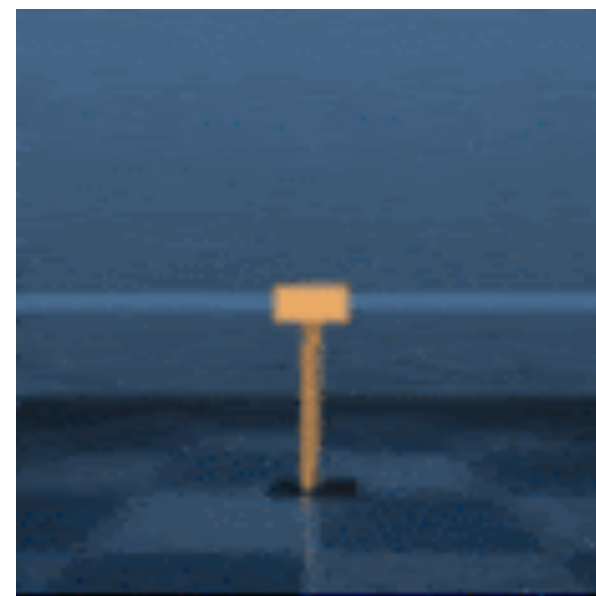
Frostbite



Collect Objects



Watermaze



Sparse Cartpole



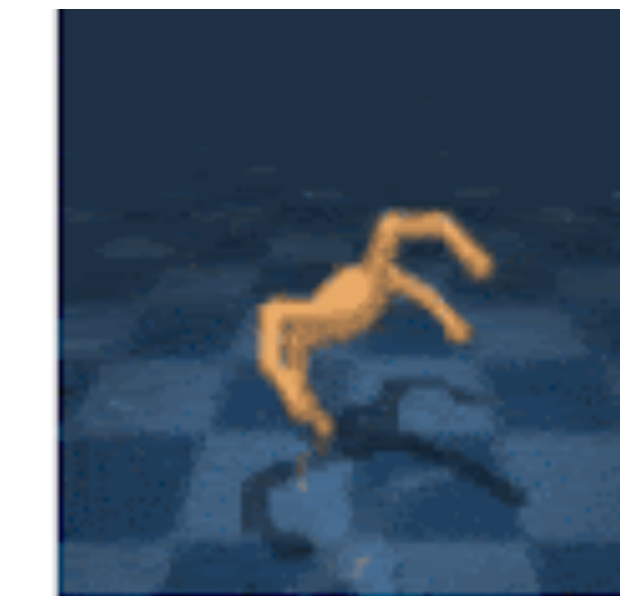
Acrobot Swingup



Hopper Hop



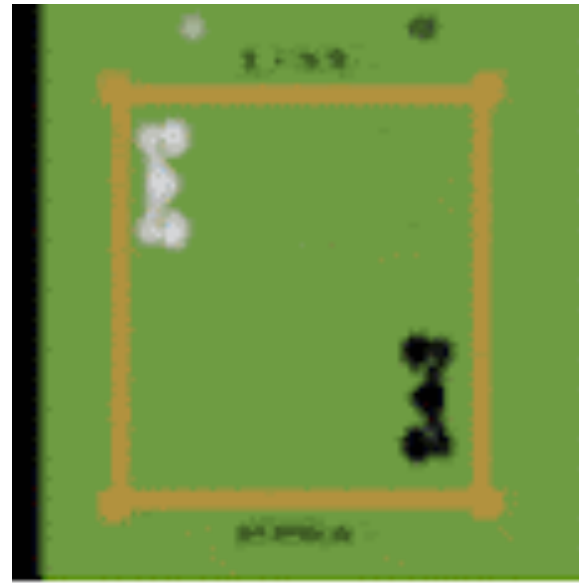
Walker Run



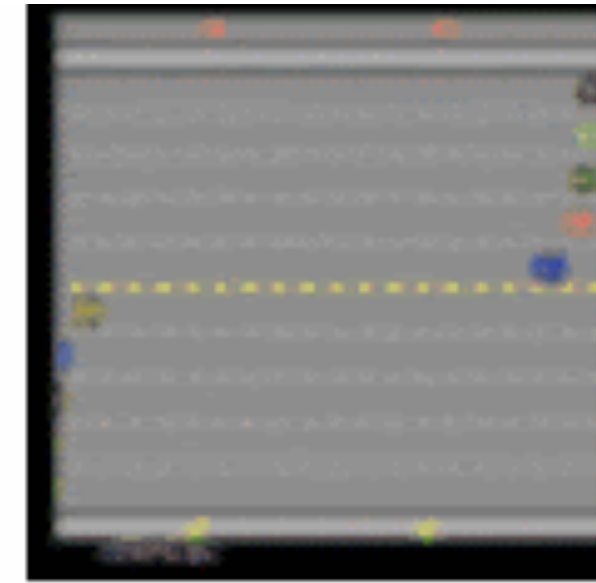
Quadruped Run

Is this from the actual simulator or predictions made by a model?

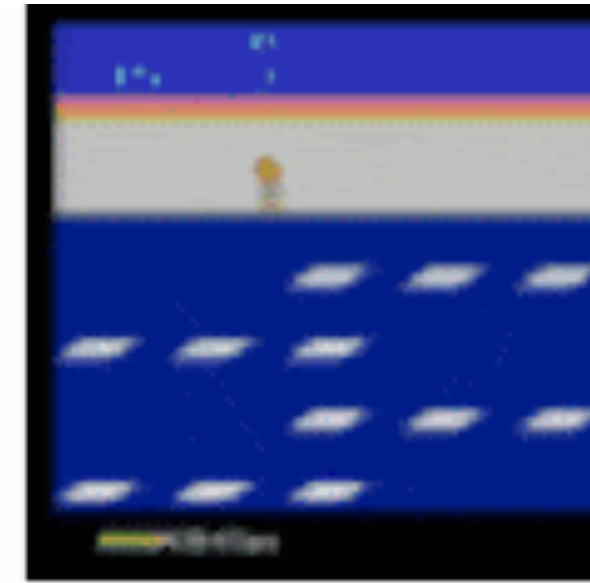
Look at the videos below



Boxing



Freeway



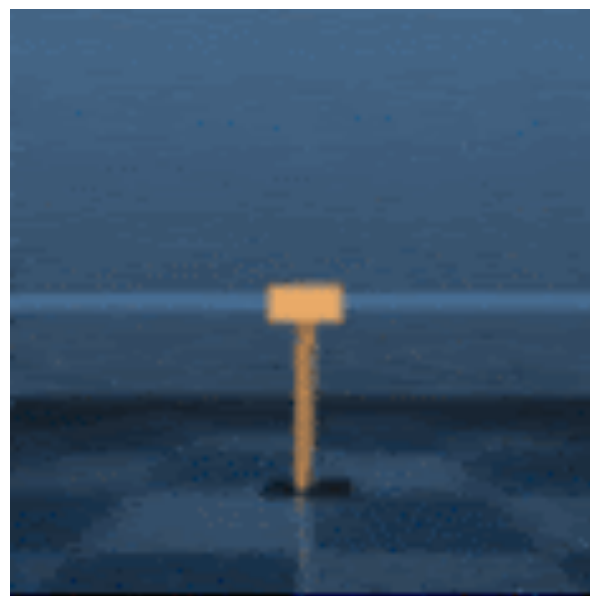
Frostbite



Collect Objects



Watermaze



Sparse Cartpole



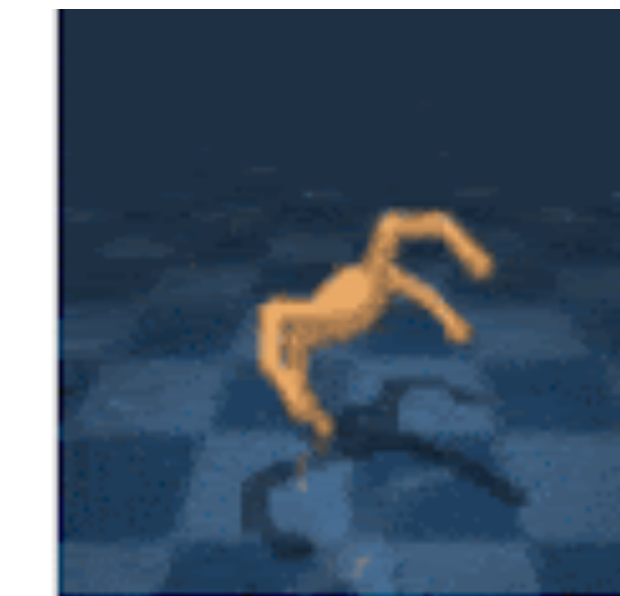
Acrobot Swingup



Hopper Hop



Walker Run

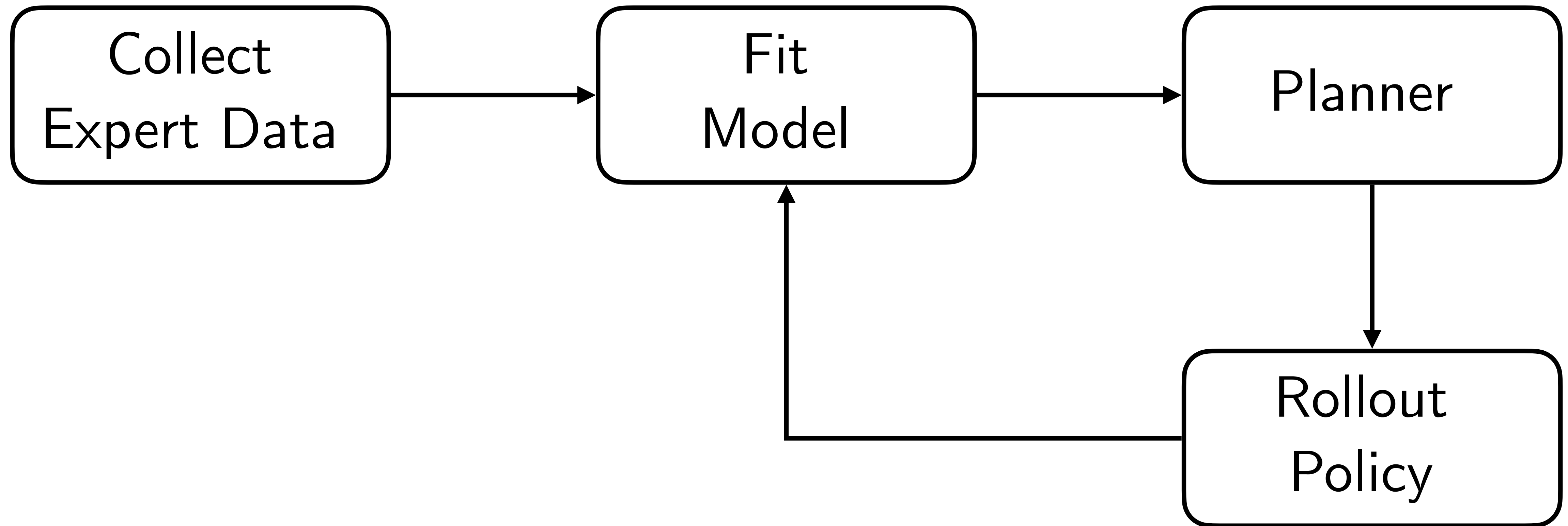


Quadruped Run

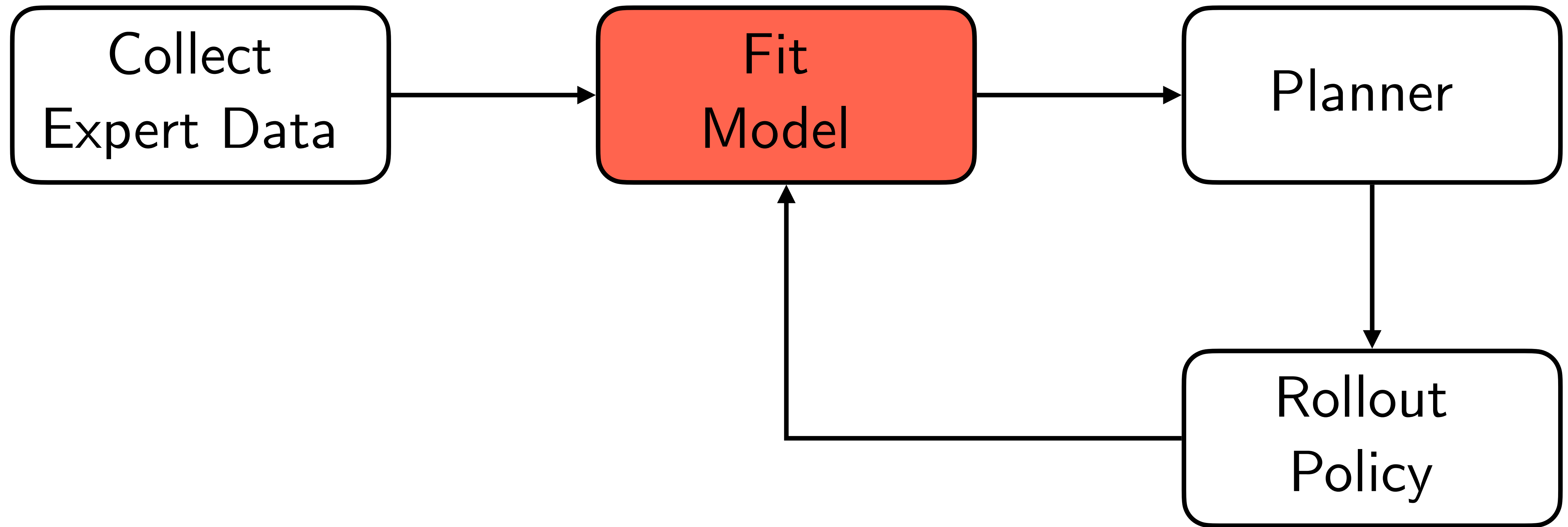
Predictions by a model!

Recap: Model-based RL

(Ross & Bagnell, 2012)



How does DREAMER fit a model?



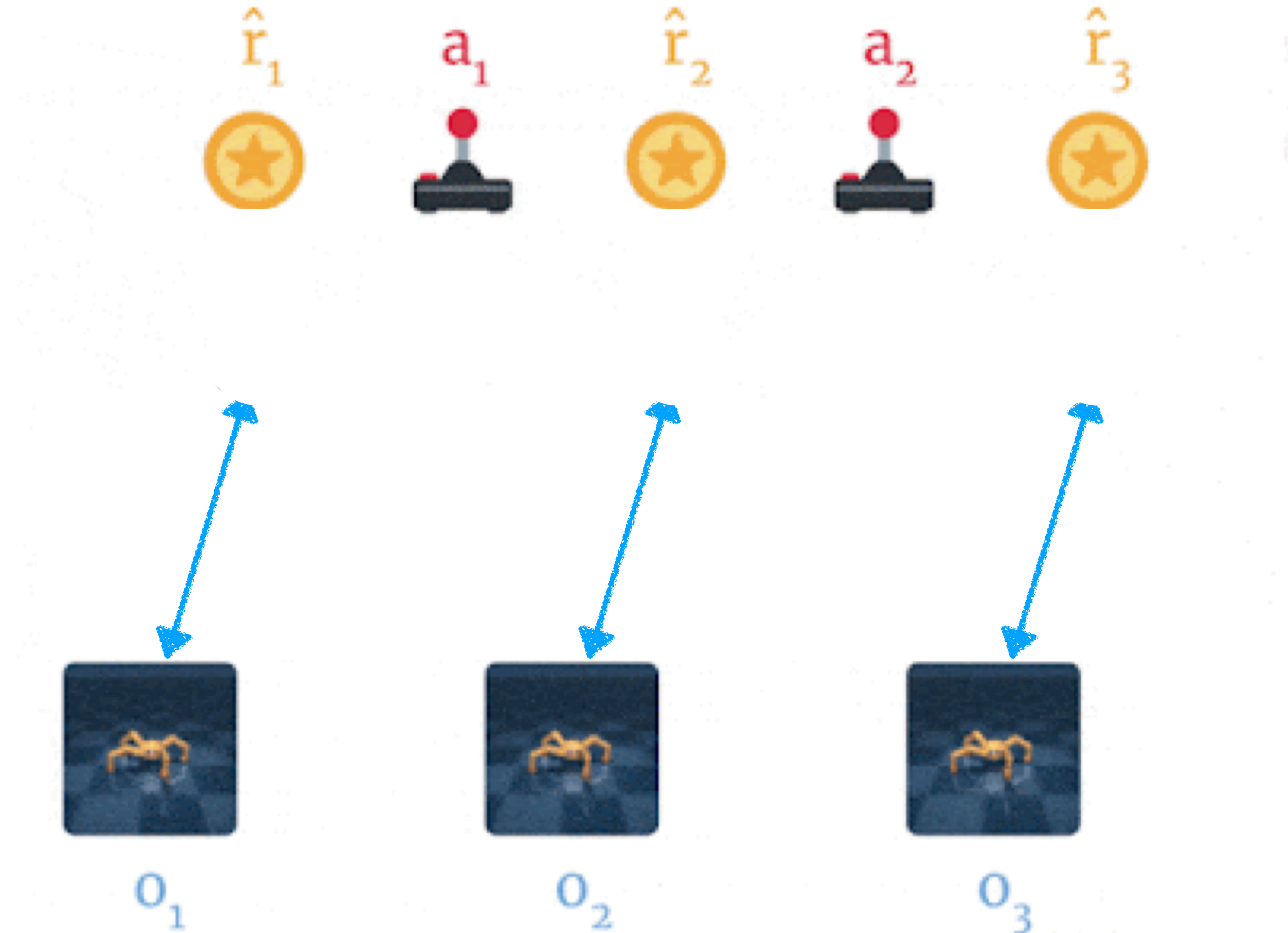
Goal: Fit a Model given data

Given Data:

Observations, rewards,
actions

Goal: Fit a Model given data

Given:
Observations, rewards,
actions



Predict:
States,
Dynamics Function,
Reward Function

Actions



Observations



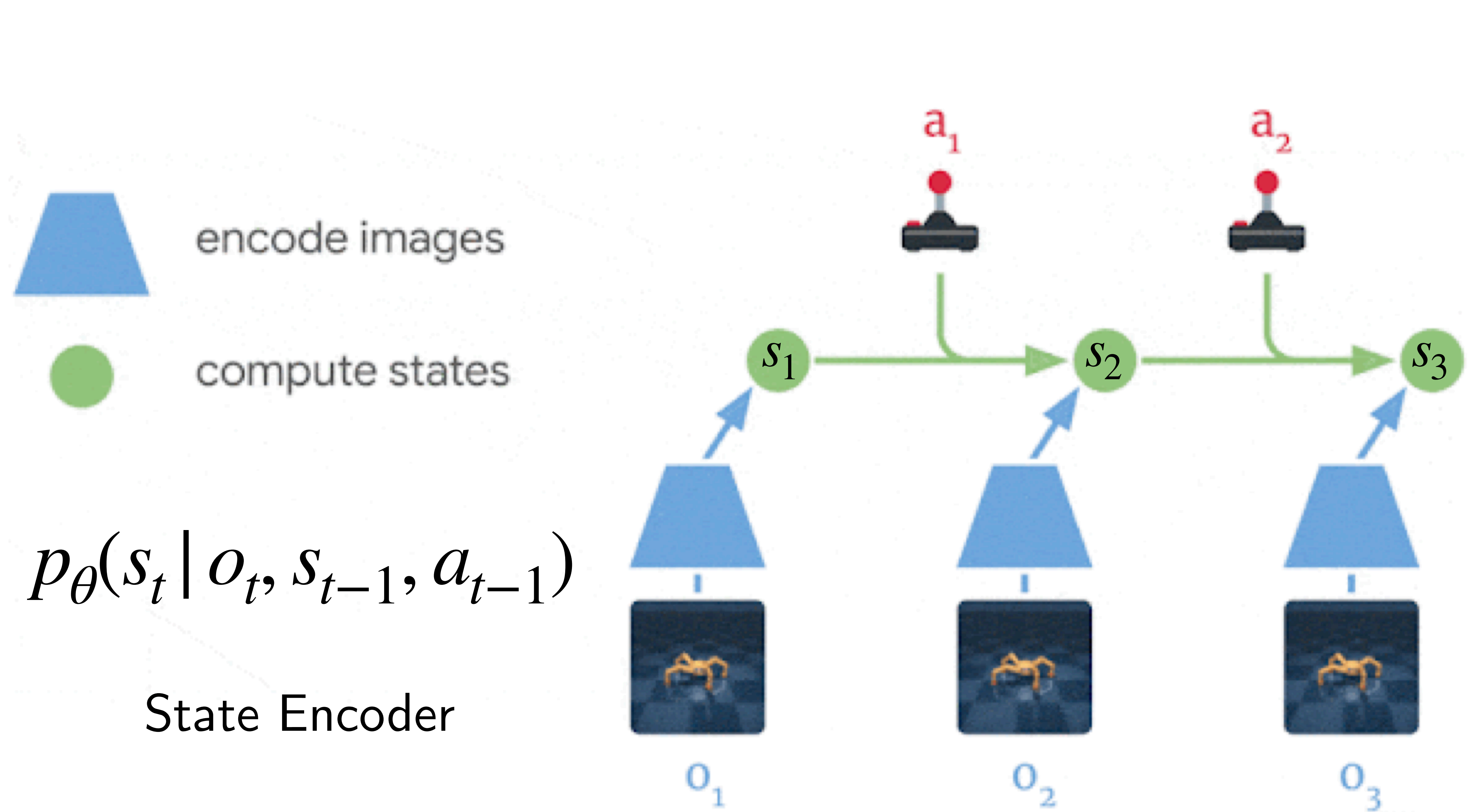
o_1



o_2



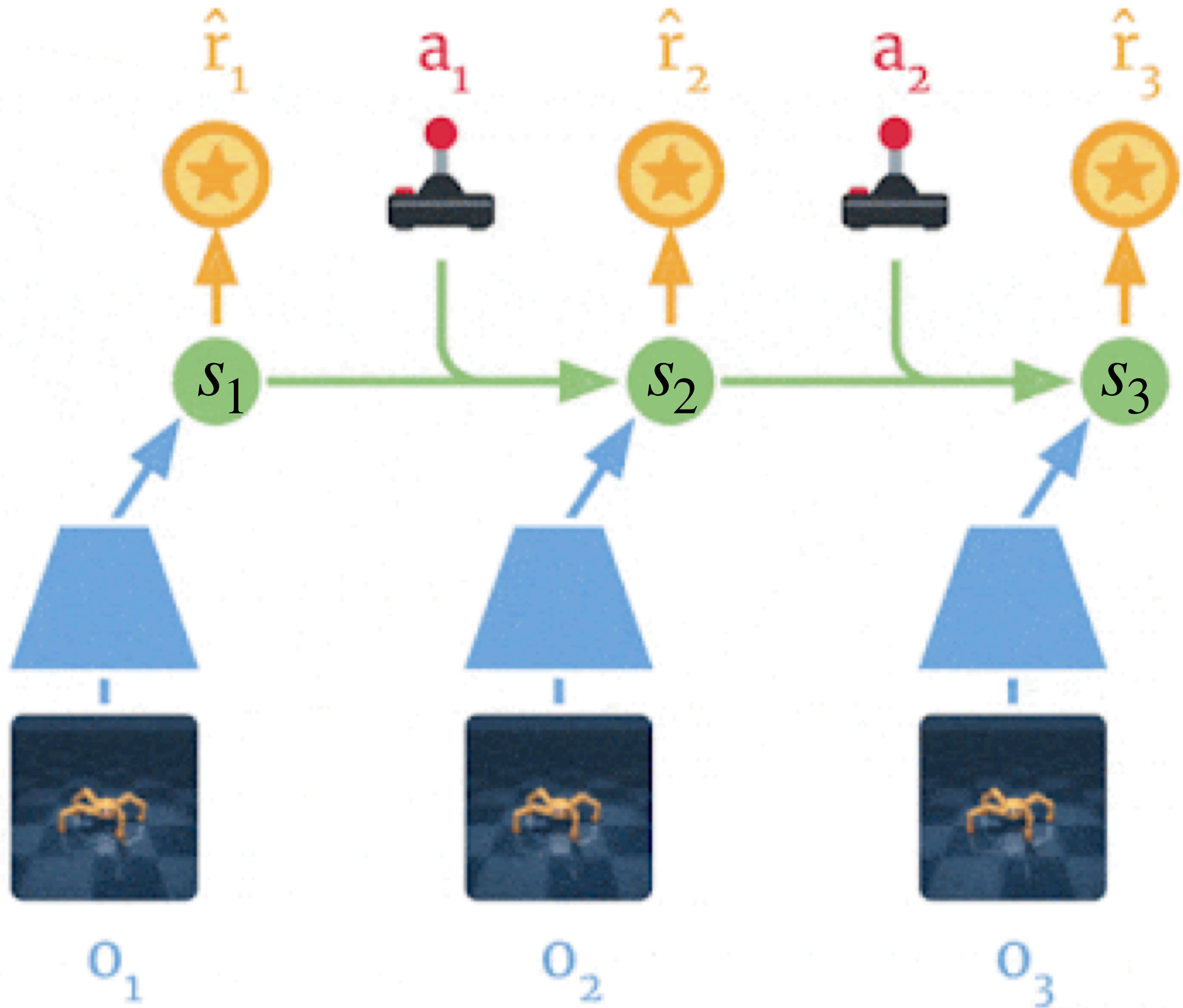
o_3



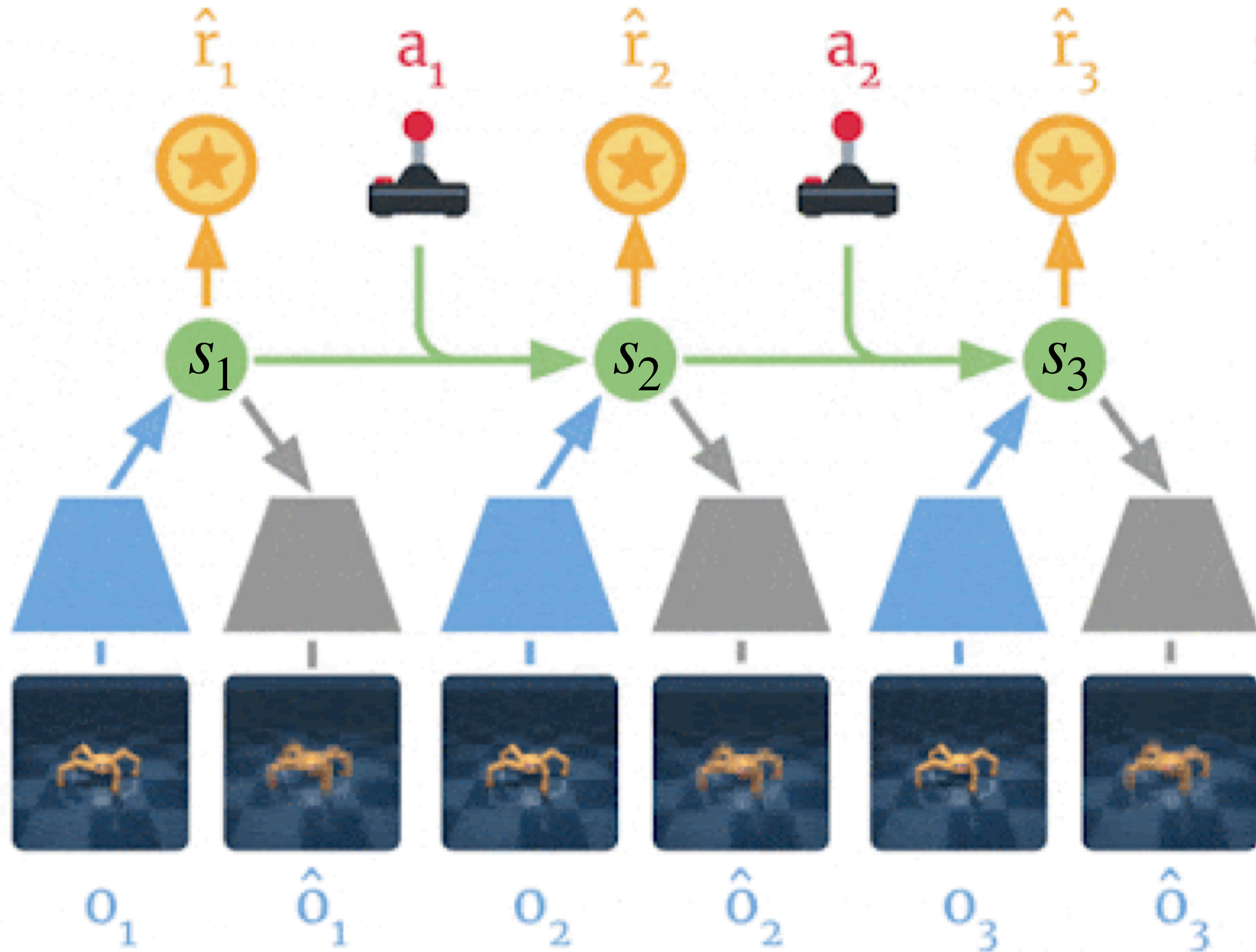
$$\mathcal{L} = (r_t - \hat{r}_t)^2$$

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

Reward Decoder



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

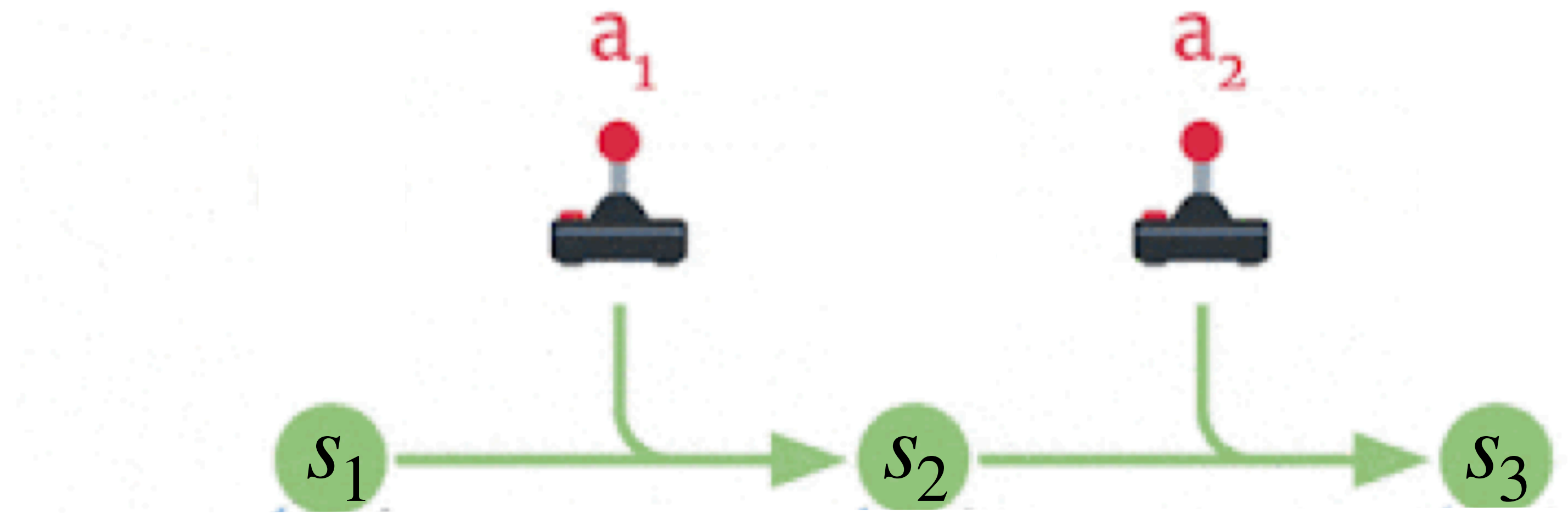


$$q_{\theta}(o_t | s_t)$$

Observation Decoder

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

Dynamics
Function



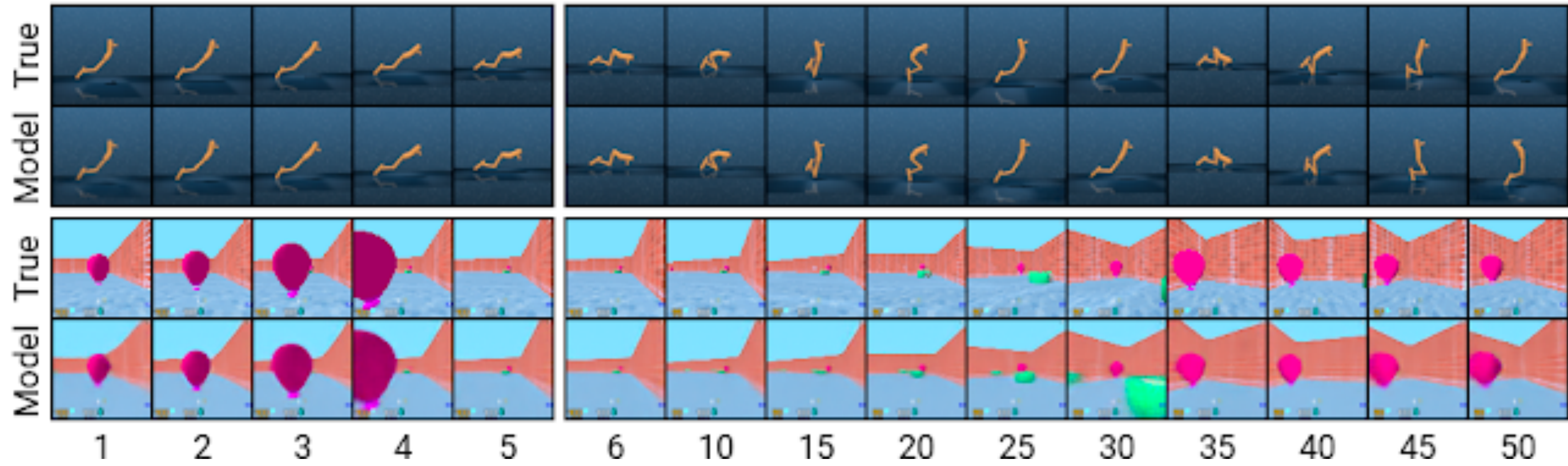
Results: Learning World Model



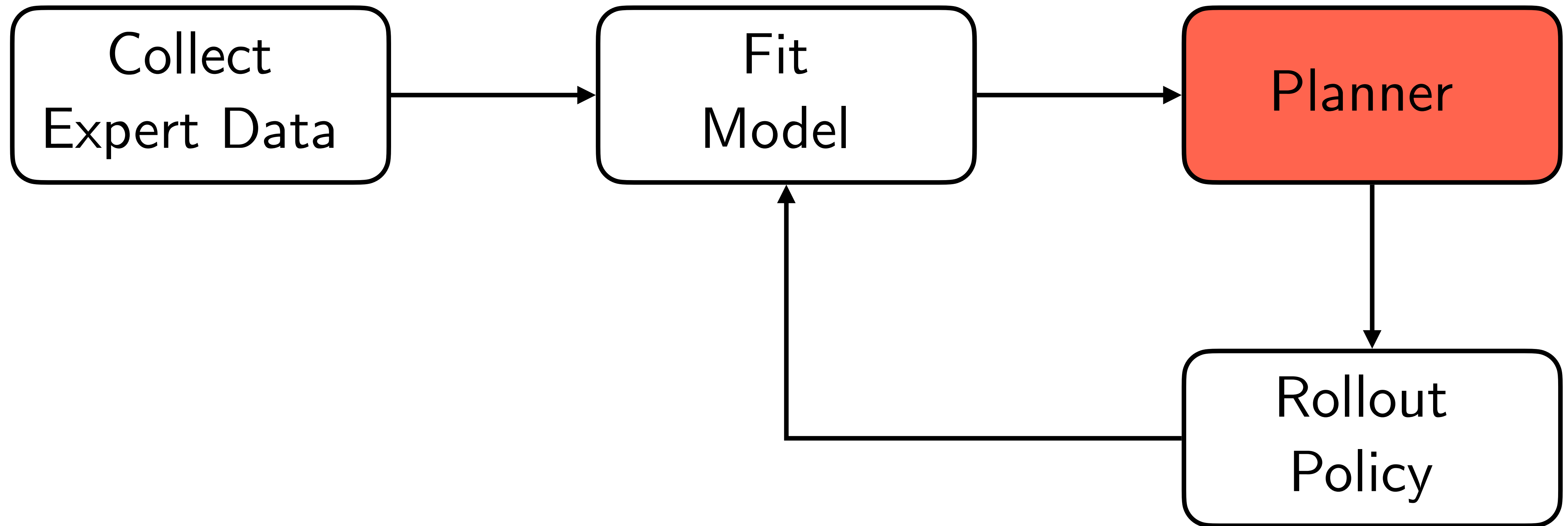
Results: Learning World Model

Input Images

Future Outcomes



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

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

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

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

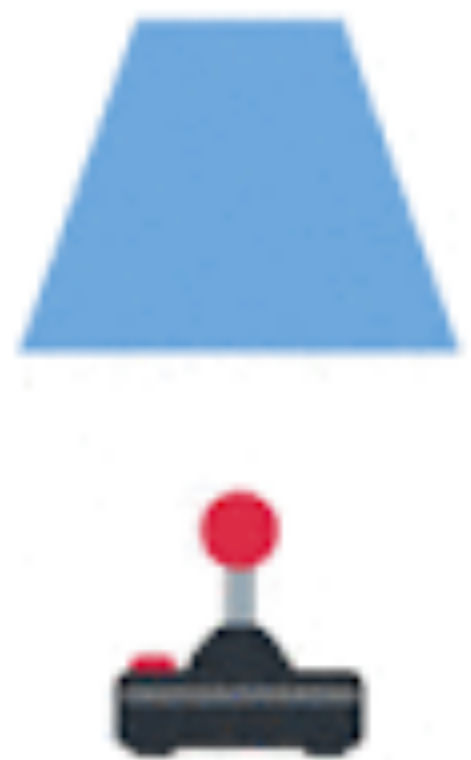


O_1



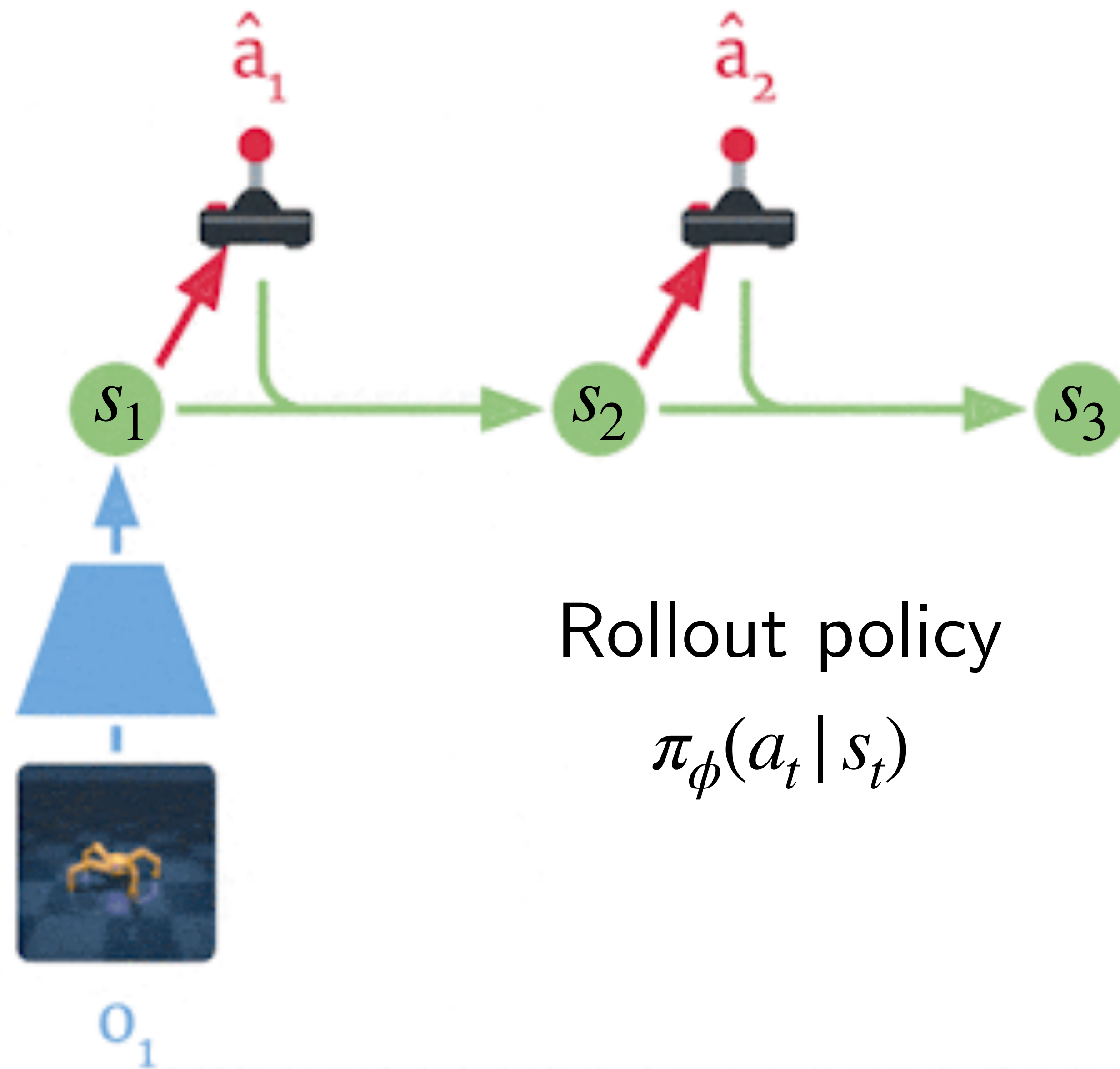
encode images





encode images

imagine ahead





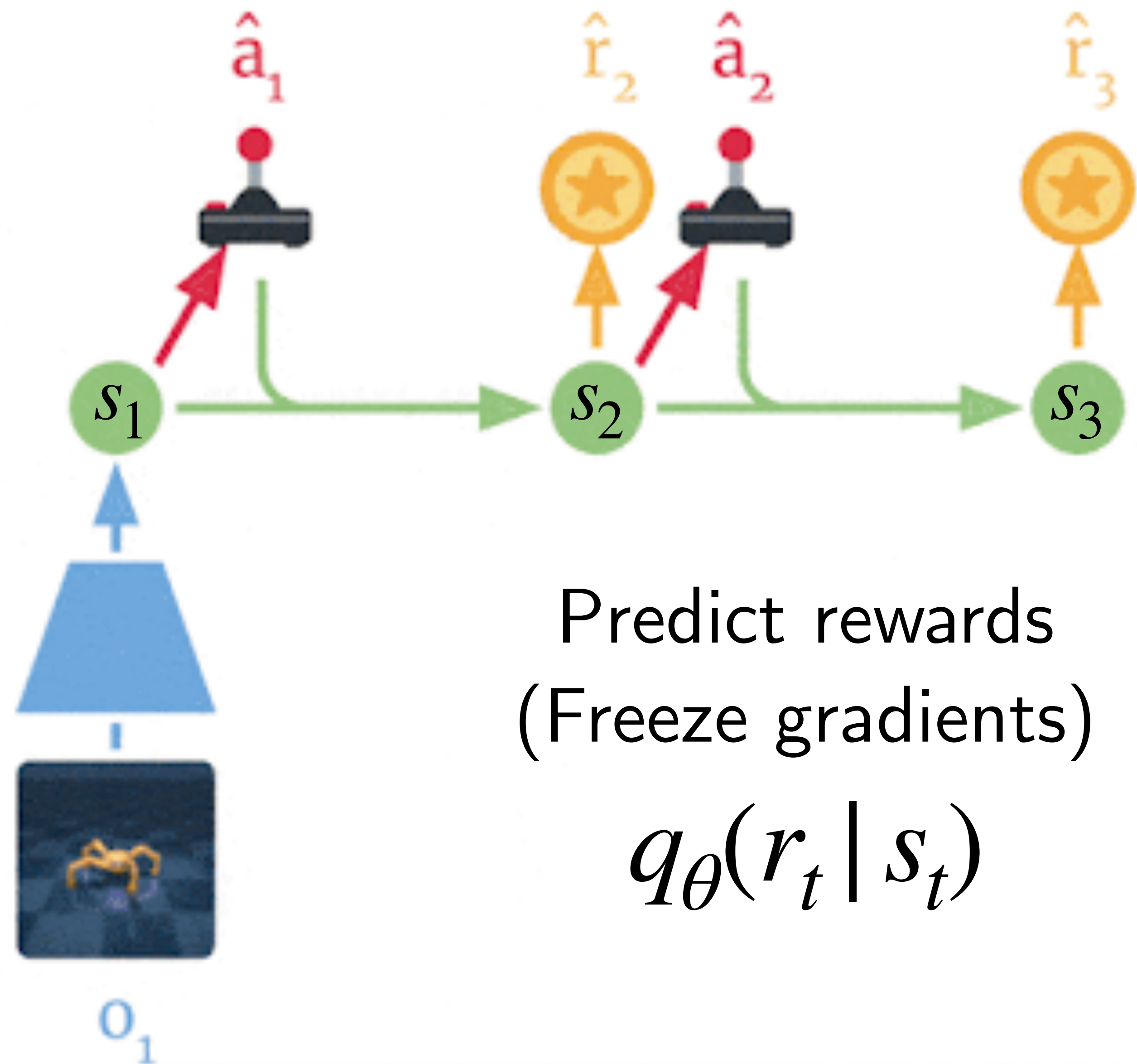
encode images



imagine ahead



predict rewards





encode images



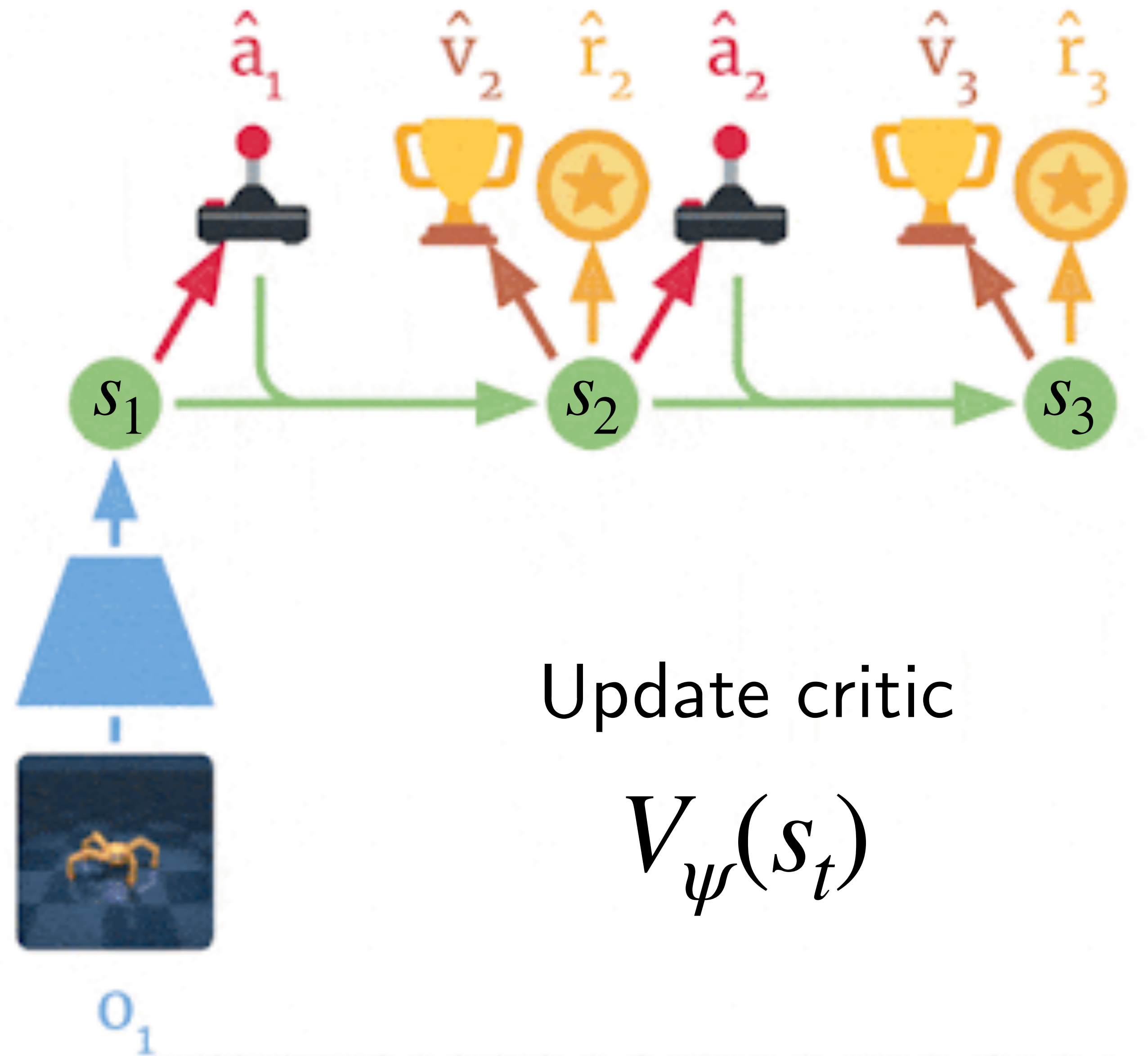
imagine ahead



predict rewards



predict values





encode images



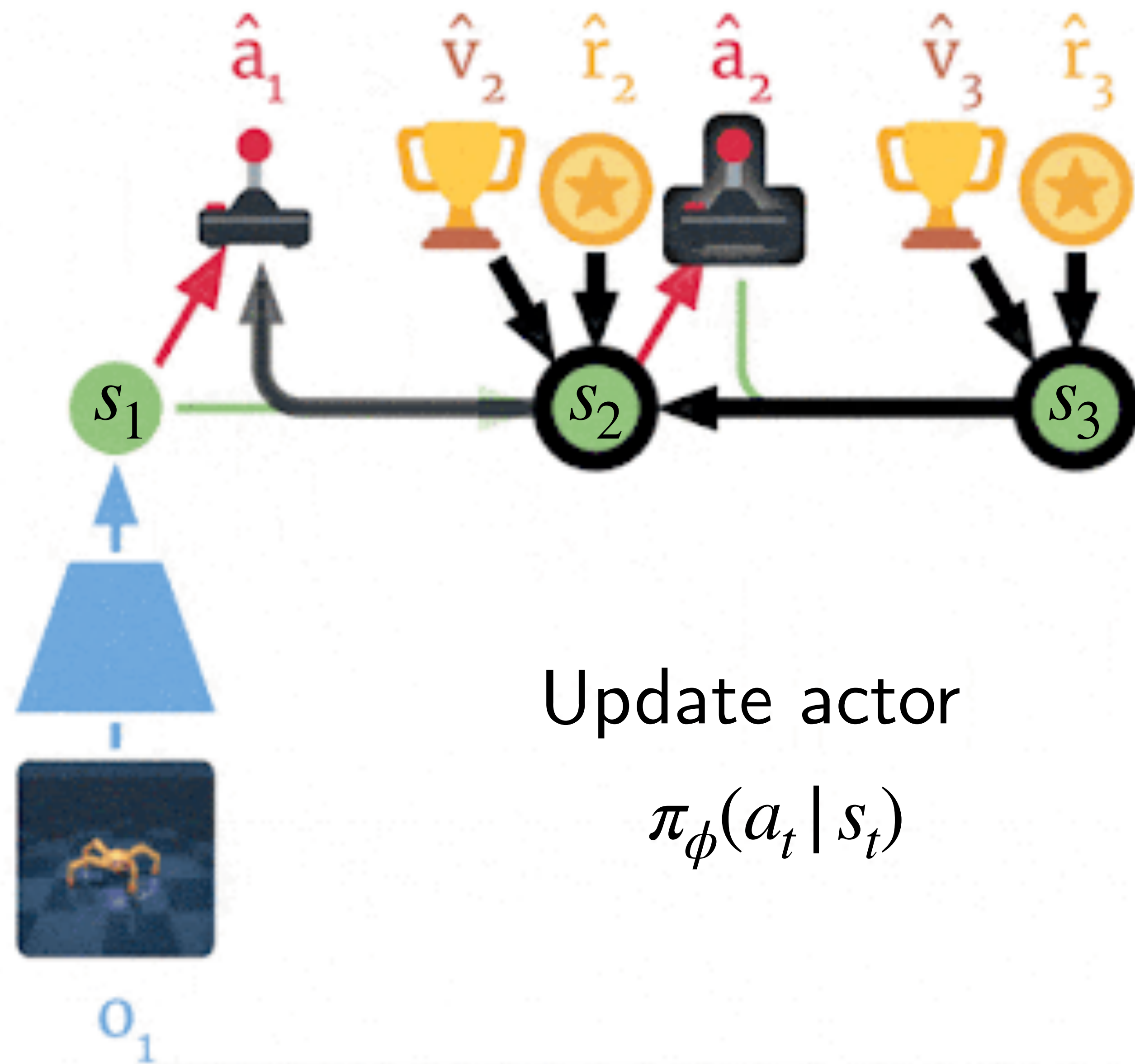
imagine ahead



predict rewards



predict values

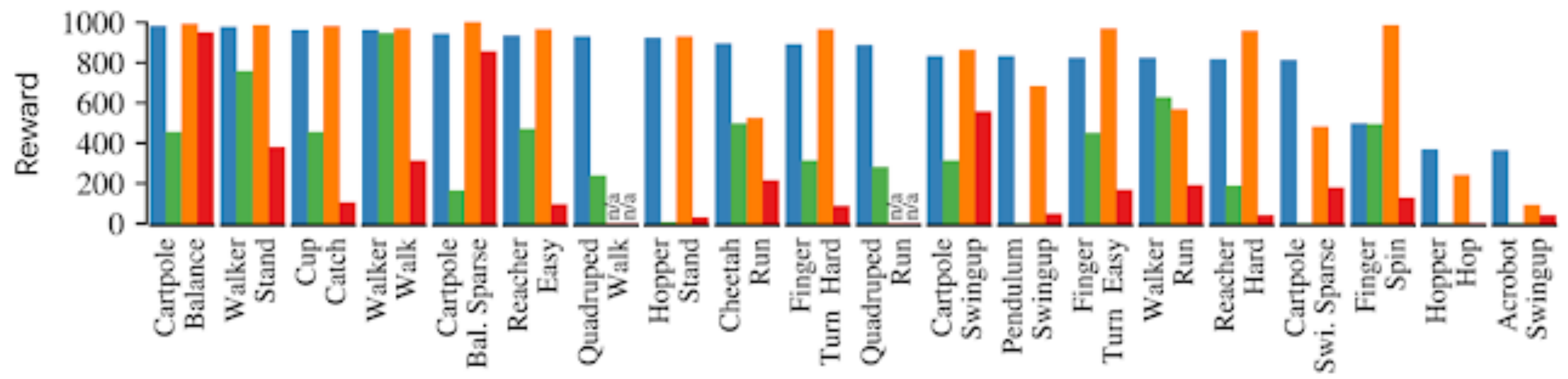


DREAMER: Results



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

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



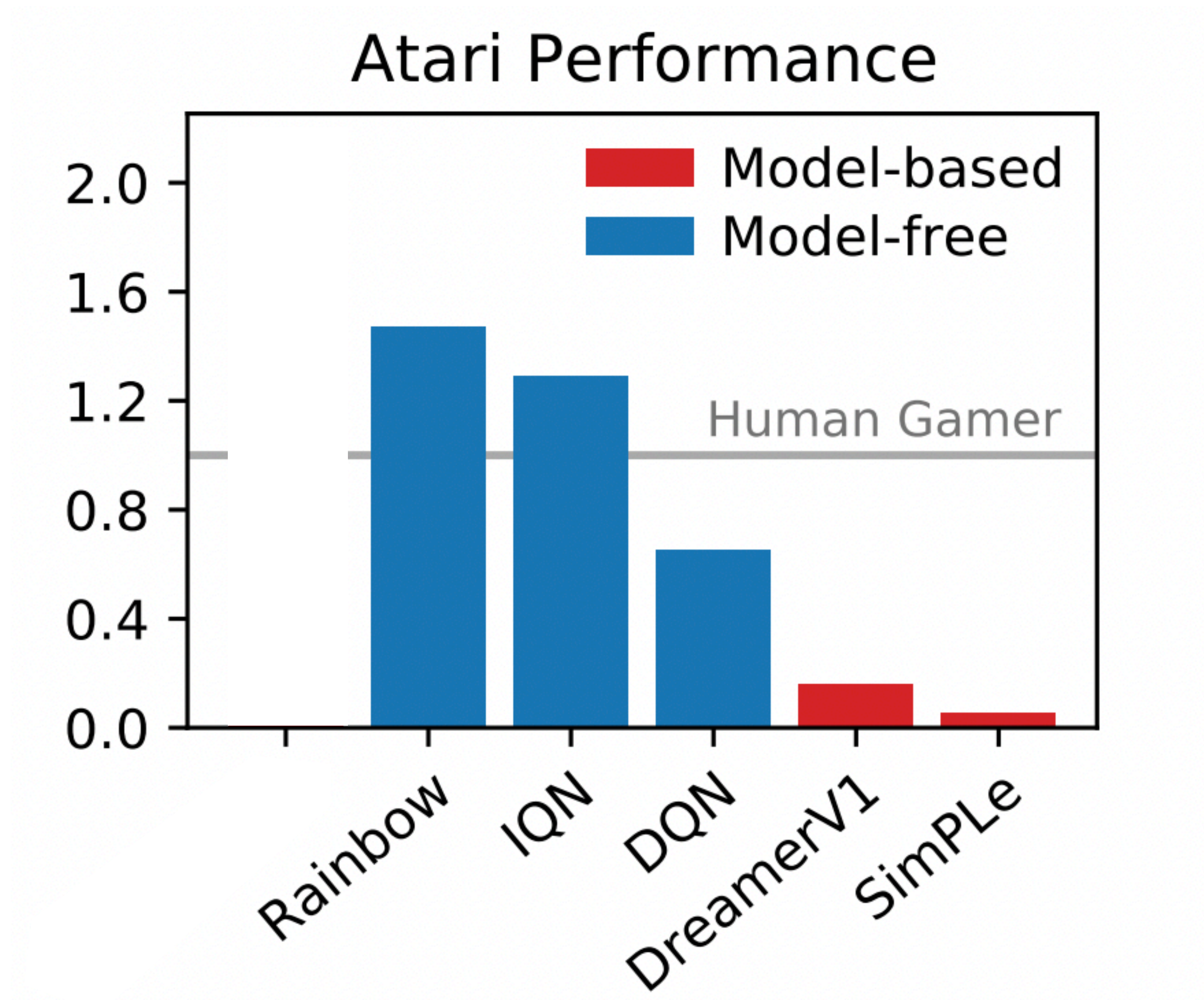
DREAMER is a template
for Model-based RL

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

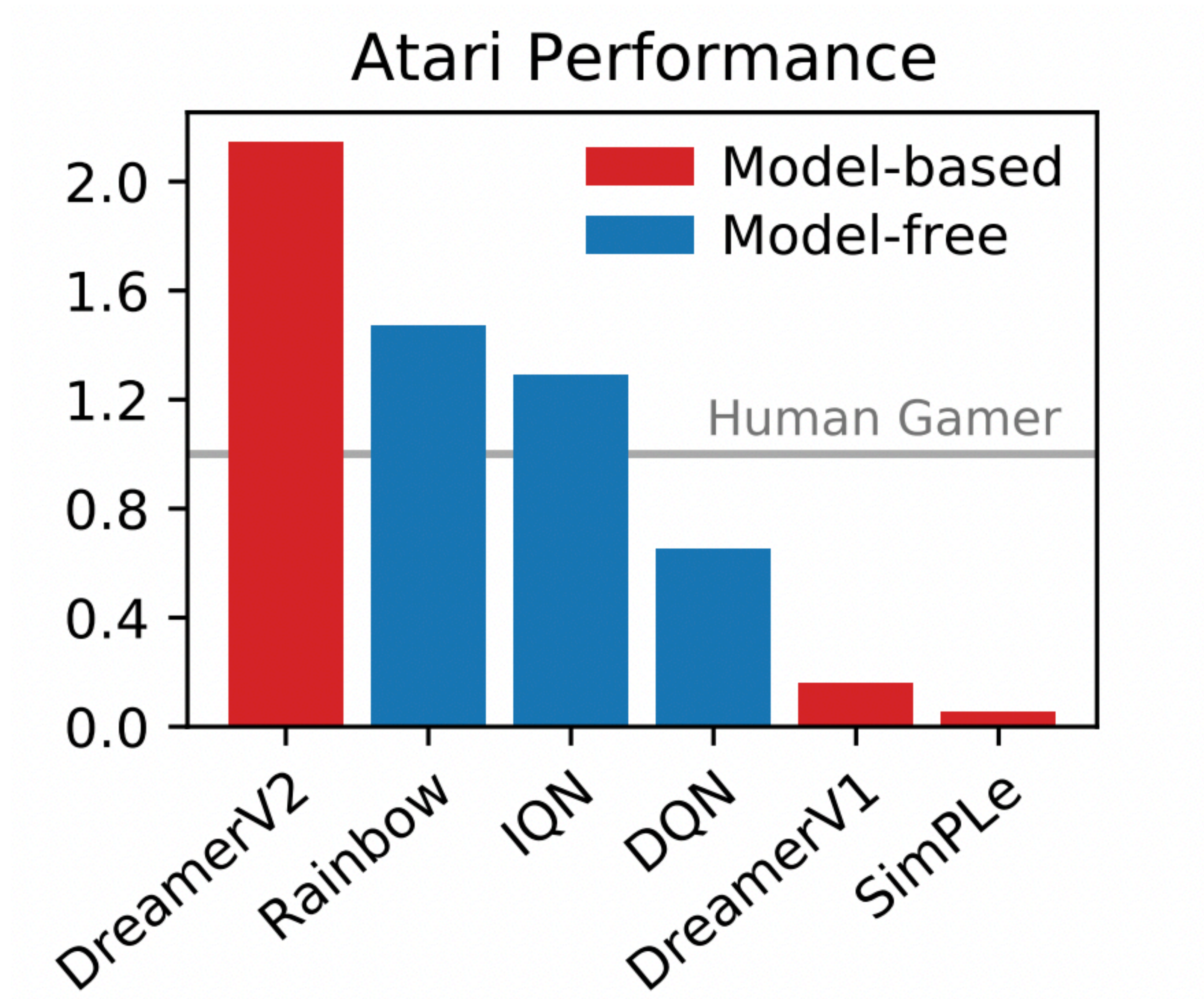
DREAMER V2:

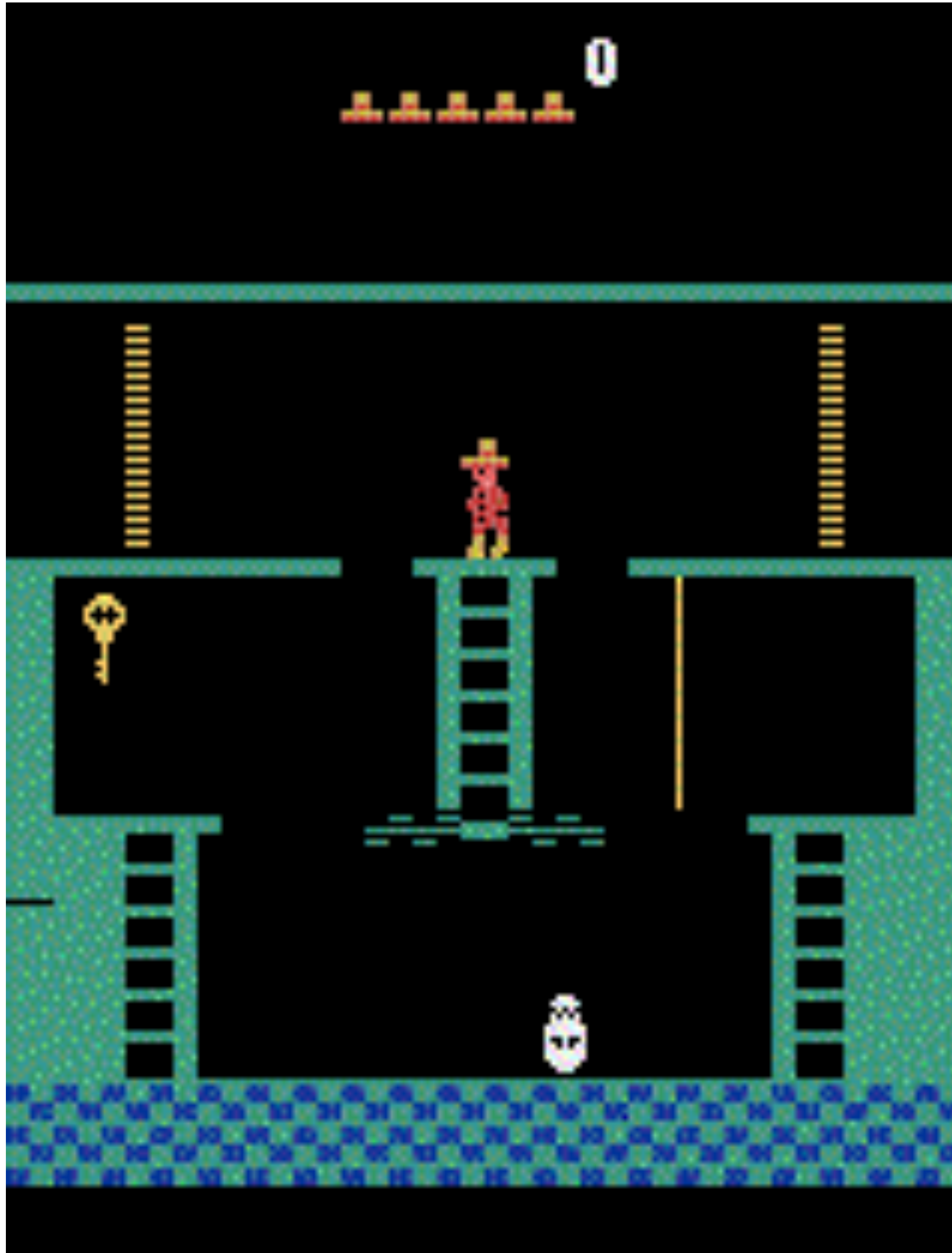
Tackling the world of Atari Games

Atari was hard for Model Based RL



DreamerV2 beats all model free!



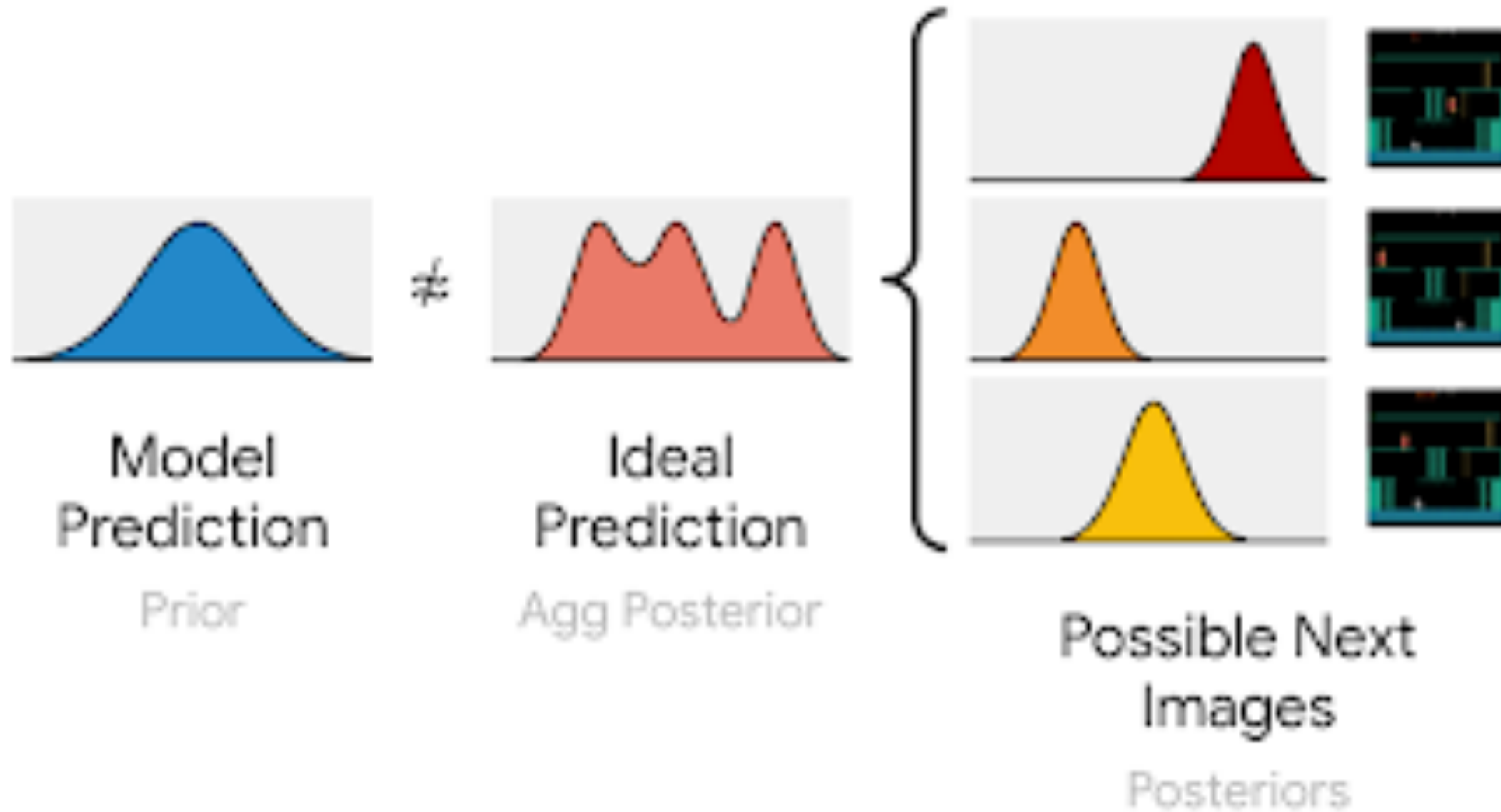


Montezuma's Revenge:

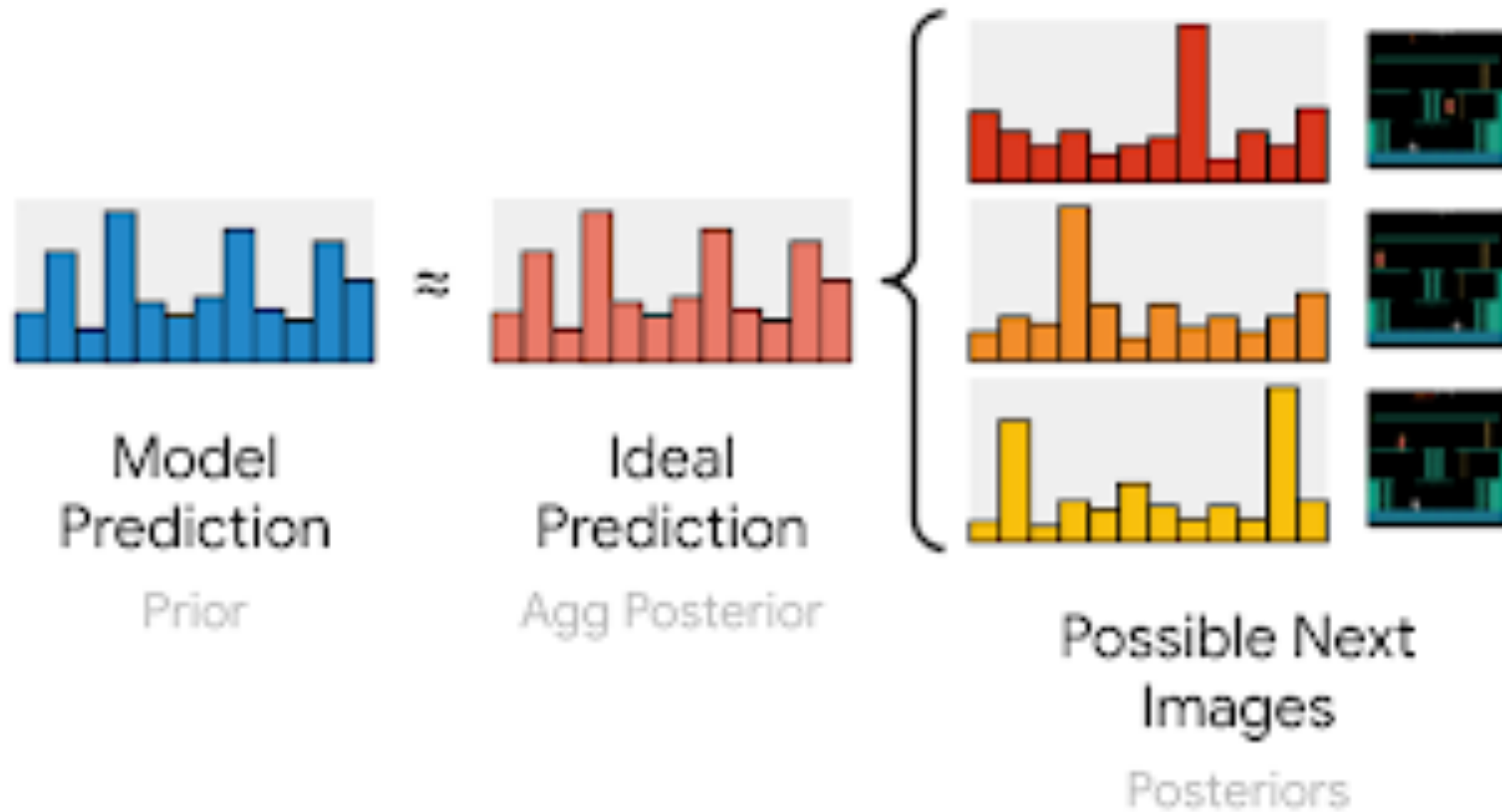
A really challenging
Atari Game!

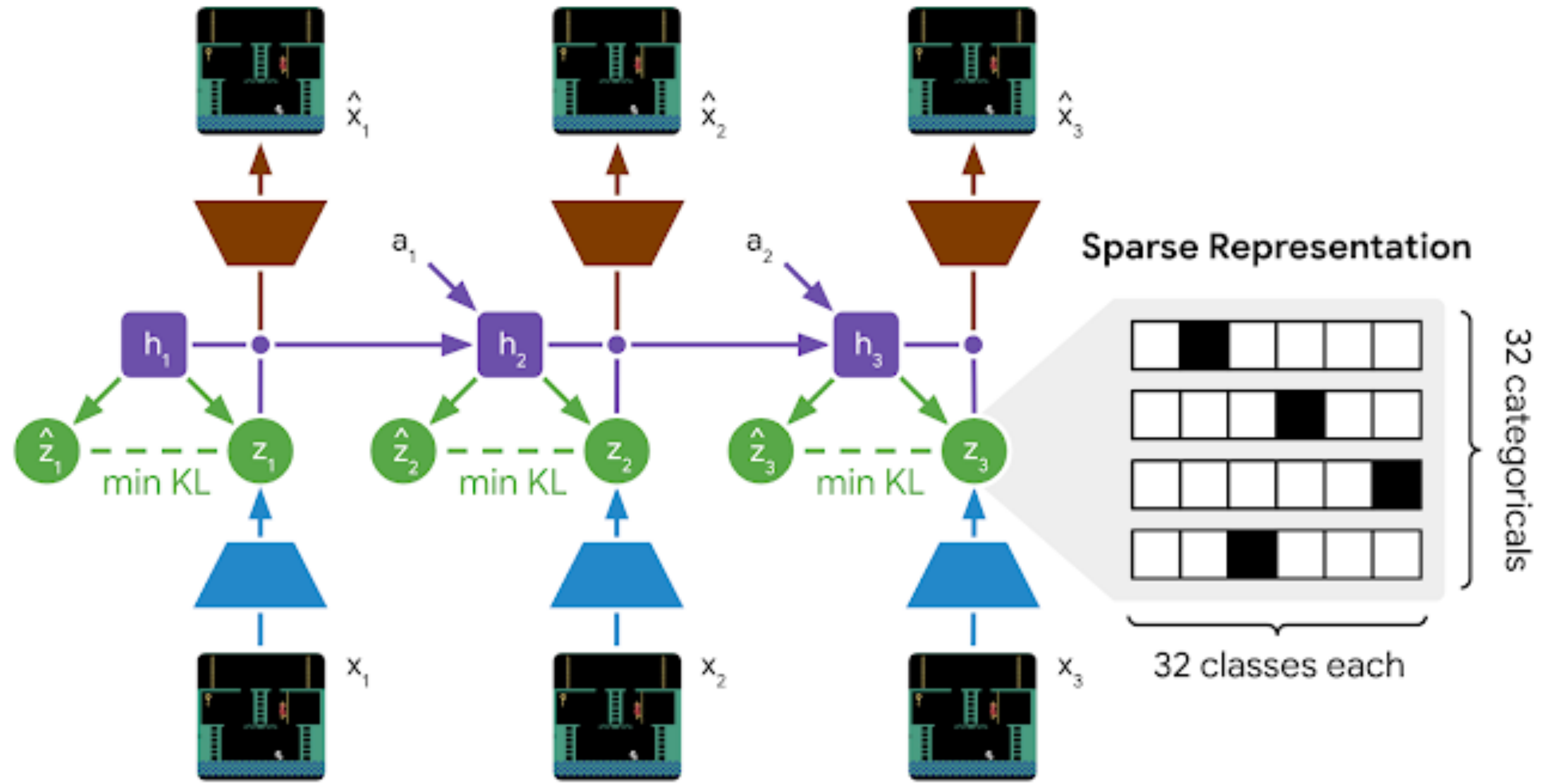
Challenge: Dreamer V1
predicts a single mode of
dynamics

Dreamer V1 predicts single mode dynamics



Idea: Predict multiple discrete modes!





True



Model

