

Exploration

CMPT 729 G100

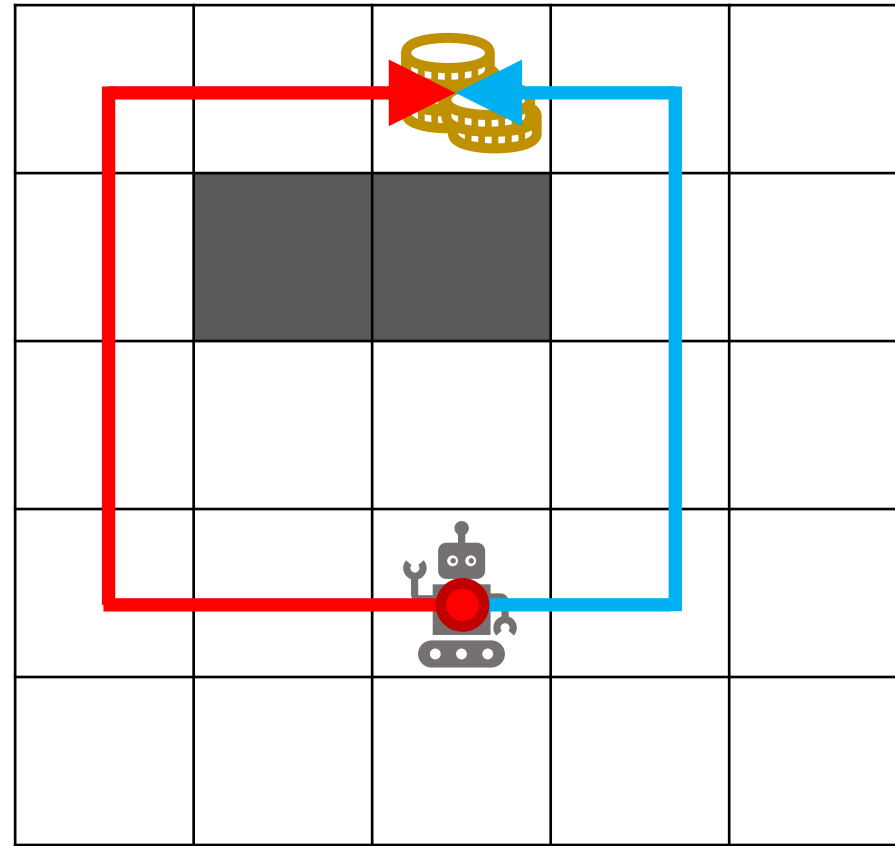
Jason Peng

Overview

- Exploration Exploitation Tradeoff
- Dense vs Sparse Rewards
- Intrinsic Motivation
- Count-Based Exploration
- Surprise Maximization

Exploration-Exploitation

Need to try new actions in case they are better



Exploration-Exploitation



Keep going to the same restaurant



Try new restaurant

Oil Drilling



Drill at best known location



Drill at a new location

Ad Recommendation



Show most successful ad



Show a different ad

Simple Exploration Strategies

ϵ -greedy exploration:

$$\pi(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 - \epsilon & \text{if } \mathbf{a} = \arg \max_{\mathbf{a}'} Q^k(\mathbf{s}, \mathbf{a}') \\ \epsilon & \text{otherwise} \end{cases}$$

Boltzmann exploration:

$$\pi(\mathbf{a}|\mathbf{s}) = \frac{1}{Z} \exp \left(\frac{1}{\beta} Q^k(\mathbf{s}, \mathbf{a}) \right)$$

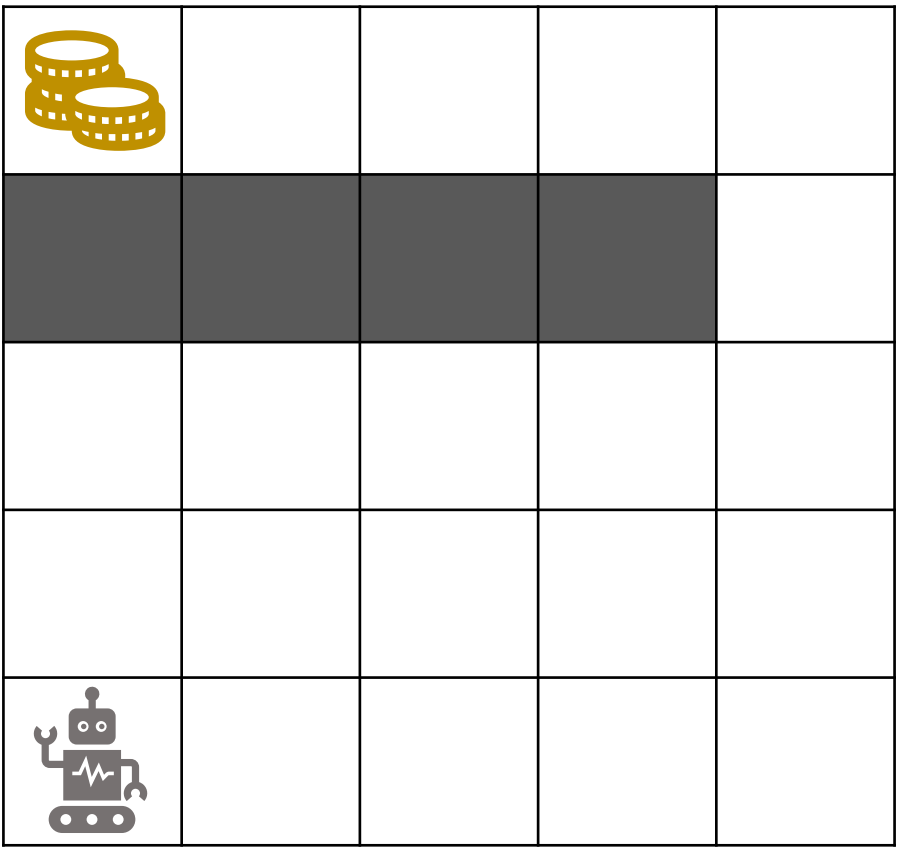
 Good action coverage

 Bad state coverage

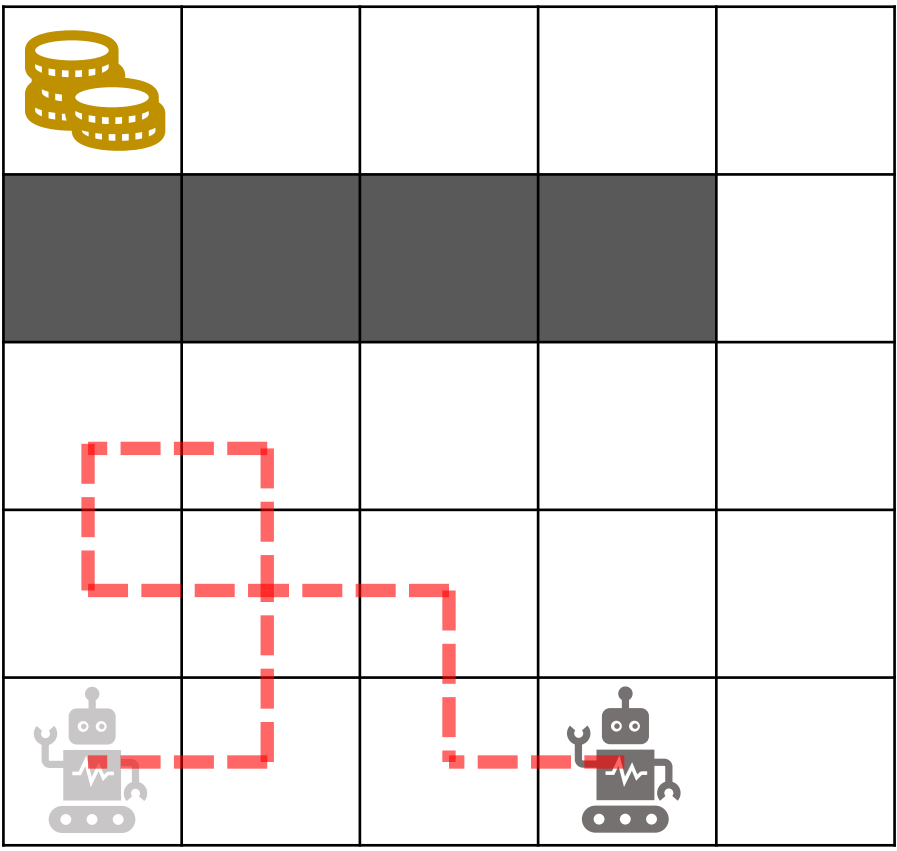
Gaussian policy:

$$\pi(\mathbf{a}|\mathbf{s}) = \frac{1}{Z} \exp \left(-\frac{1}{2} (\mathbf{a} - \mu_{\pi}(\mathbf{s}))^{\top} \Sigma(\mathbf{s})^{-1} (\mathbf{a} - \mu_{\pi}(\mathbf{s})) \right)$$

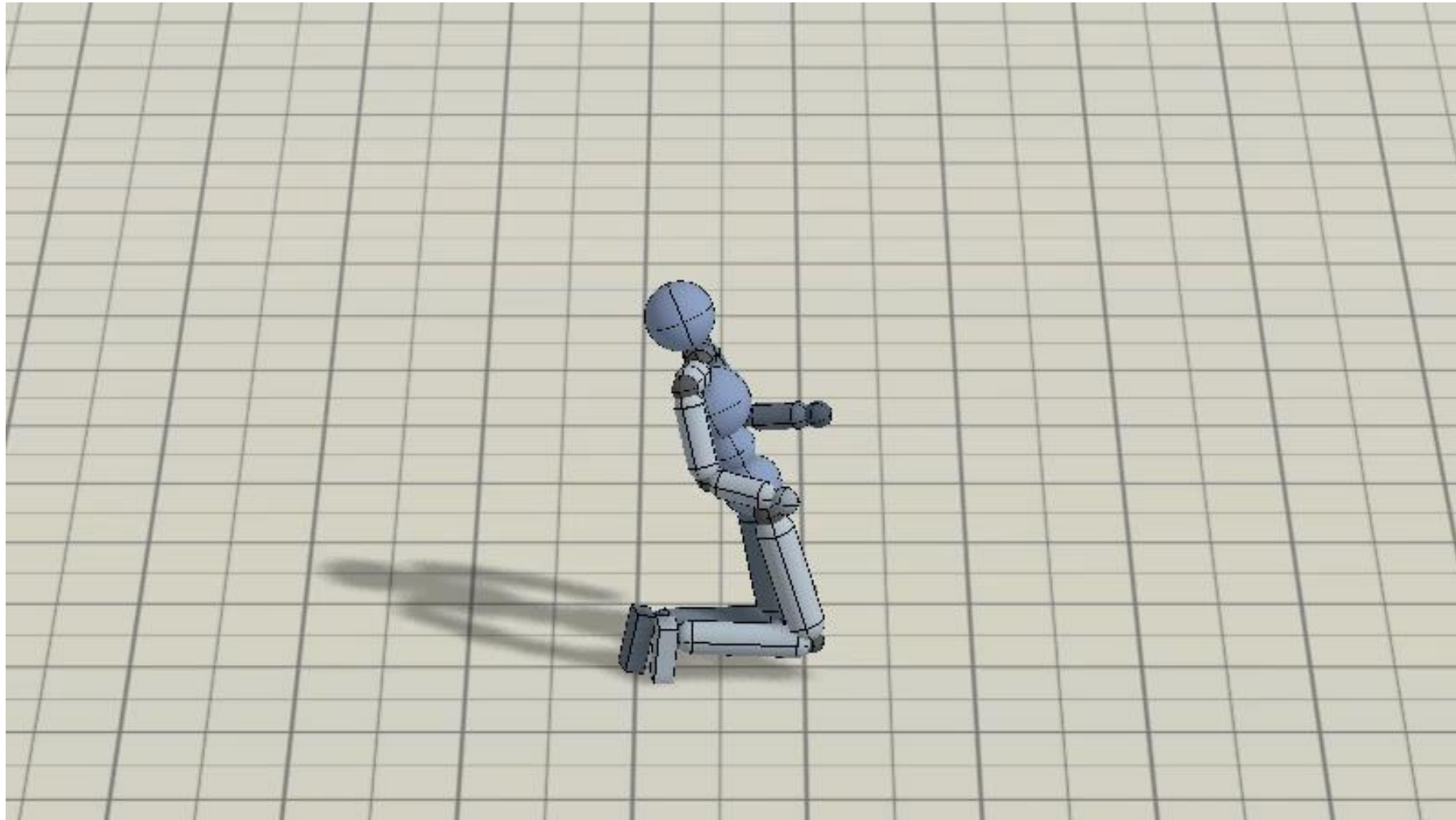
Simple Exploration Strategies



Simple Exploration Strategies



Simple Exploration Strategies



Reward Functions

- Reward function guides policy towards better actions
- Structure of reward function can have dramatic impact on exploration and performance
 - Well-shaped reward function: hard tasks can be made easy
 - Poorly shaped reward function: easy tasks can be made almost impossible

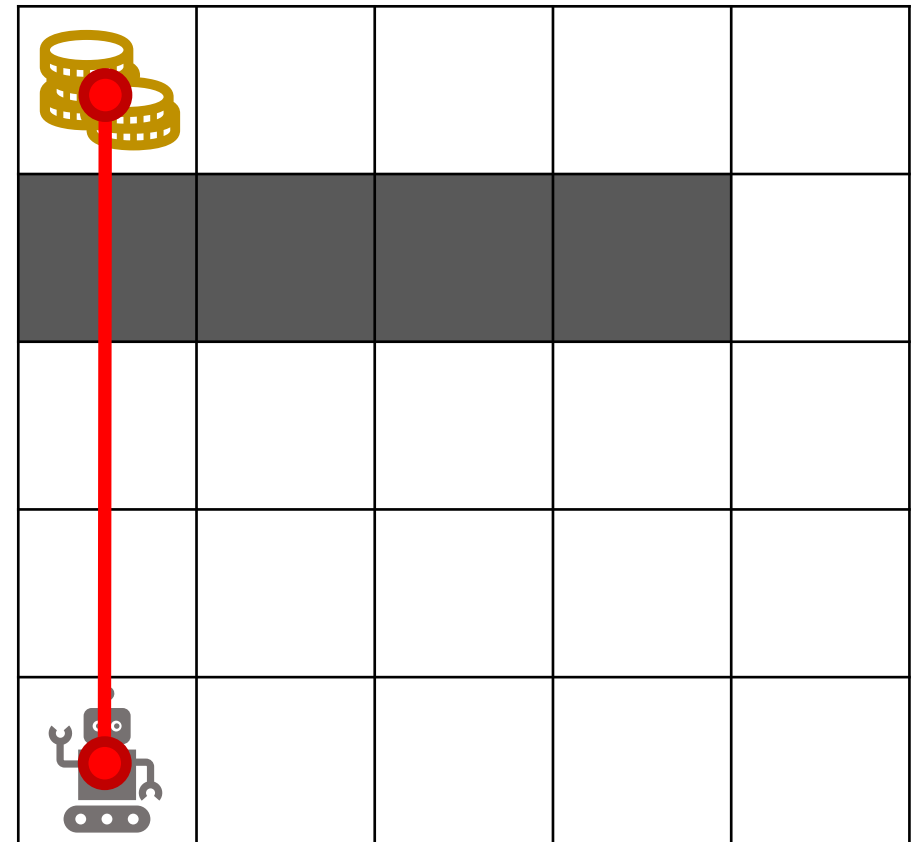
Reward Functions

- Dense reward
 - Non-zero reward at every timestep reflecting progress towards goal
- Sparse reward
 - Agent receives nonzero reward only when goal is completed

Dense Reward

Non-zero reward at every timestep
reflecting progress towards goal

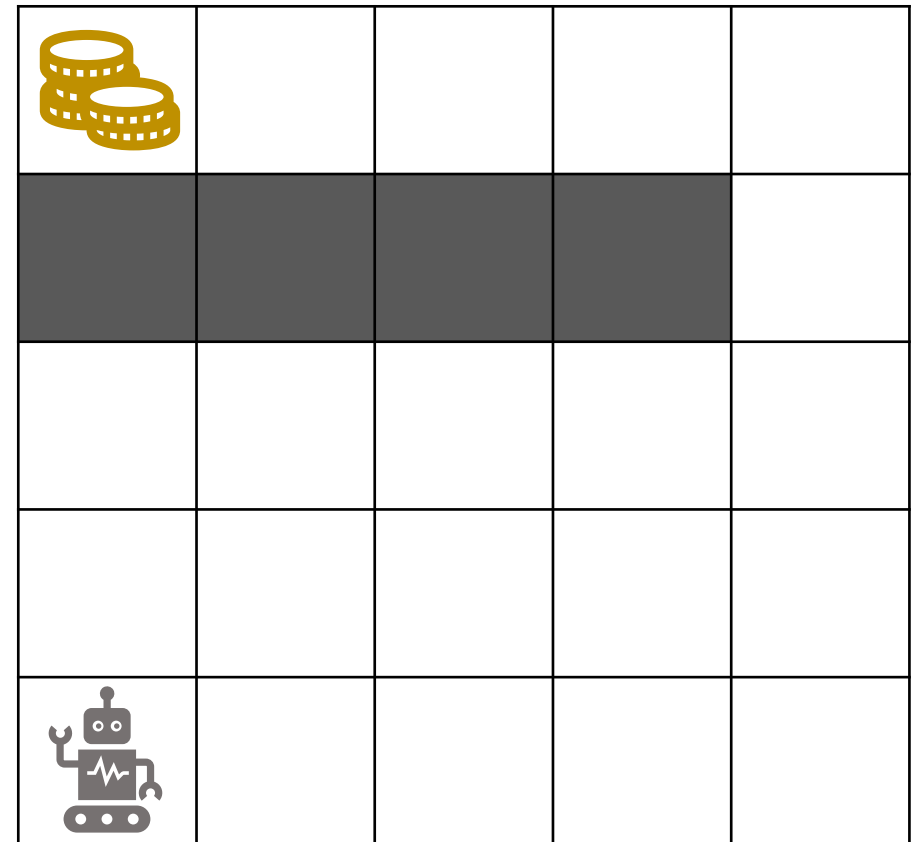
$$r_t = - \left\| \text{robot} - \text{goal} \right\|^2$$



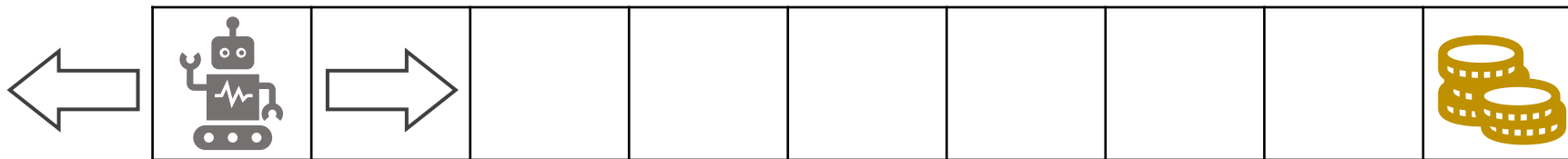
Sparse Reward

Agent receives nonzero reward only when goal is completed

$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{goal} \\ 0 & \text{otherwise} \end{cases}$$

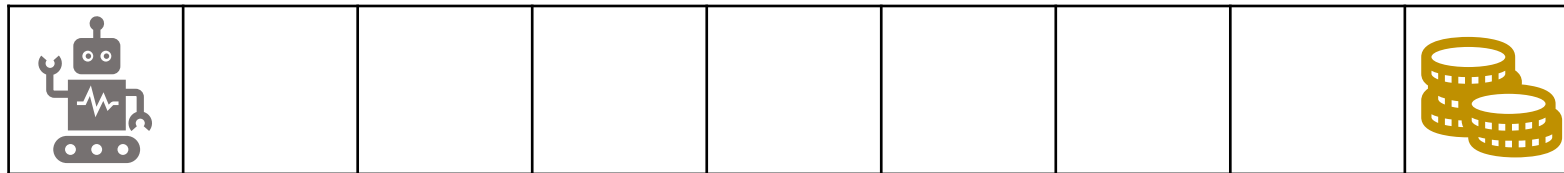


Reward Functions



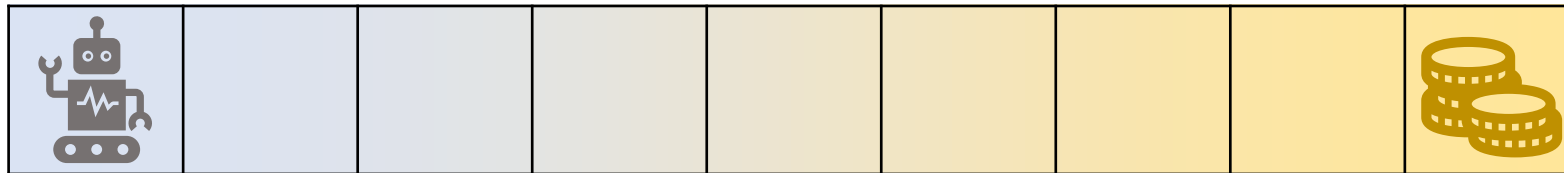
Dense Rewards

$$r_t = - \left\| \text{robot} - \text{coins} \right\|^2$$



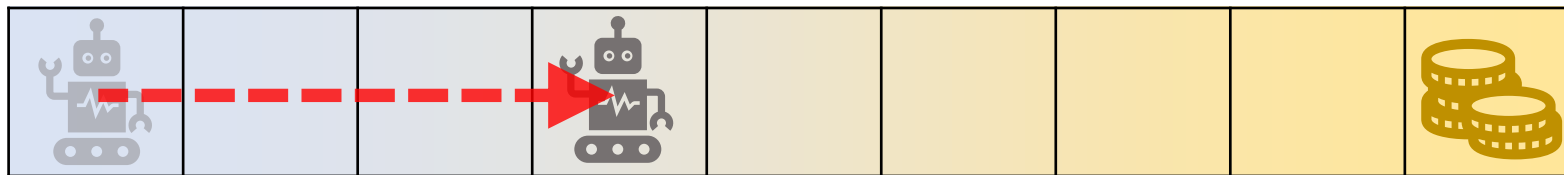
Dense Rewards

$$r_t = - \left\| \text{robot} - \text{coins} \right\|^2$$



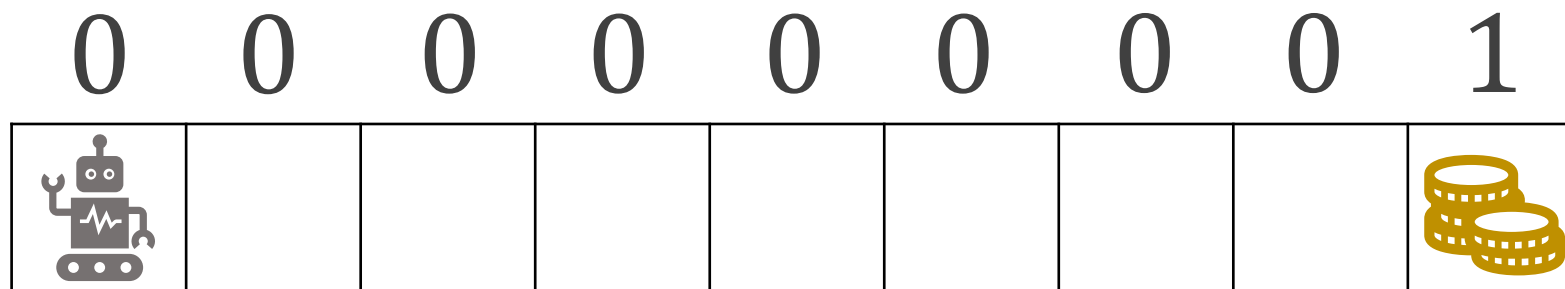
Dense Rewards

$$r_t = - \left\| \text{robot} - \text{coins} \right\|^2$$



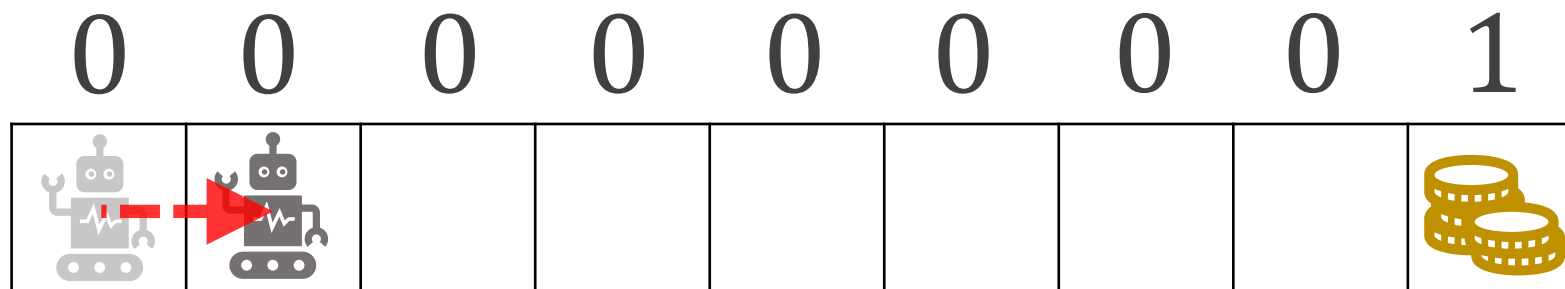
Sparse Rewards

$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{coins} \\ 0 & \text{otherwise} \end{cases}$$



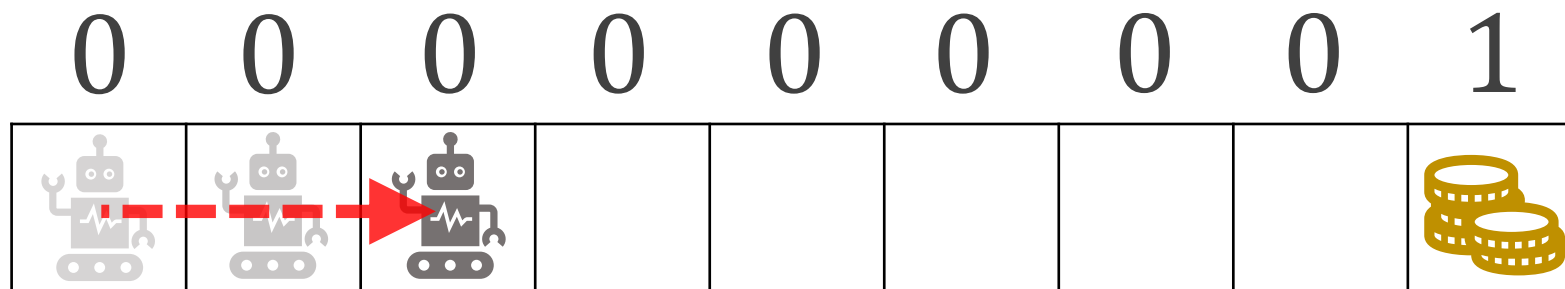
Sparse Rewards

$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{coins} \\ 0 & \text{otherwise} \end{cases}$$



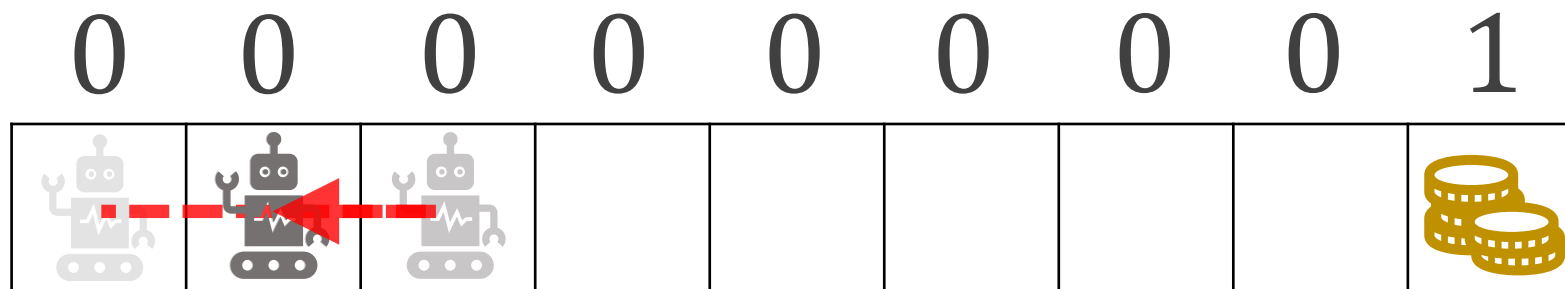
Sparse Rewards

$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{coins} \\ 0 & \text{otherwise} \end{cases}$$



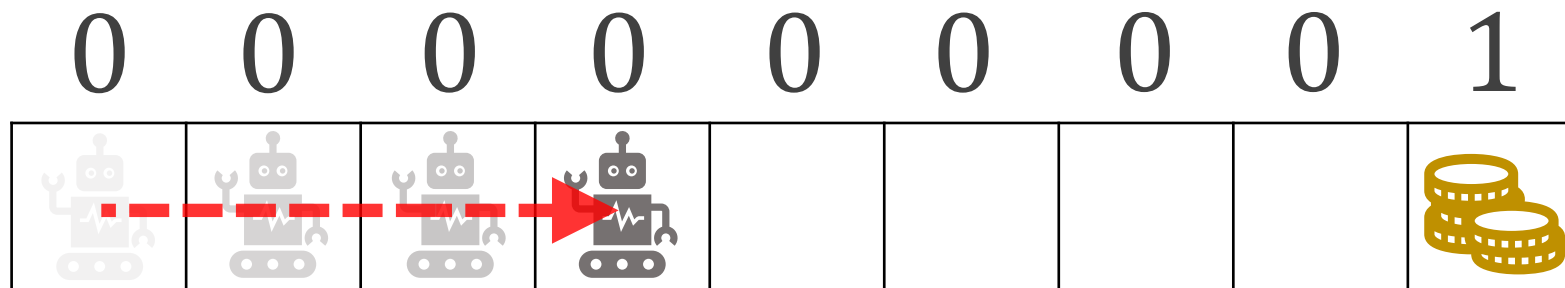
Sparse Rewards

$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{coins} \\ 0 & \text{otherwise} \end{cases}$$



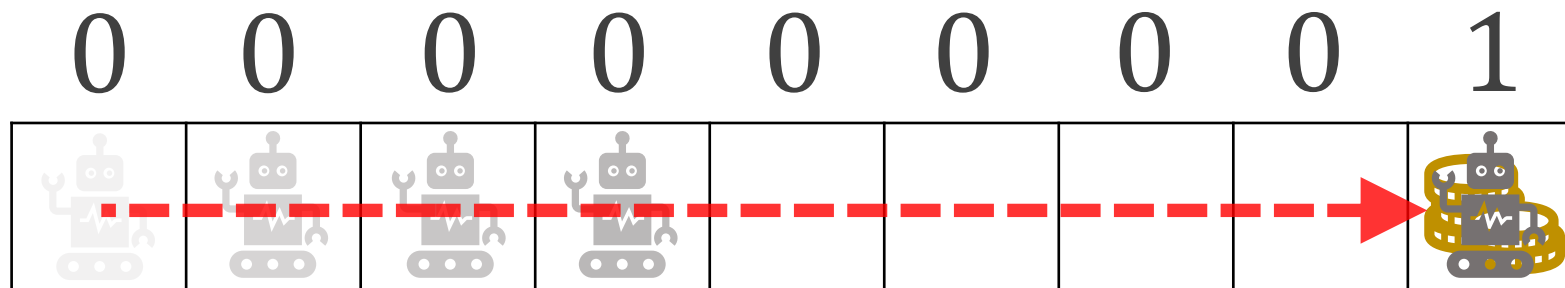
Sparse Rewards

$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{coins} \\ 0 & \text{otherwise} \end{cases}$$



Sparse Rewards

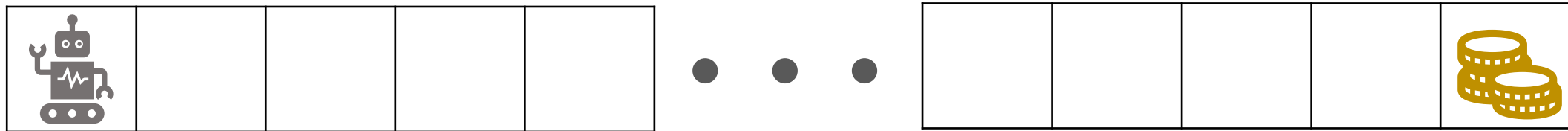
$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{coins} \\ 0 & \text{otherwise} \end{cases}$$



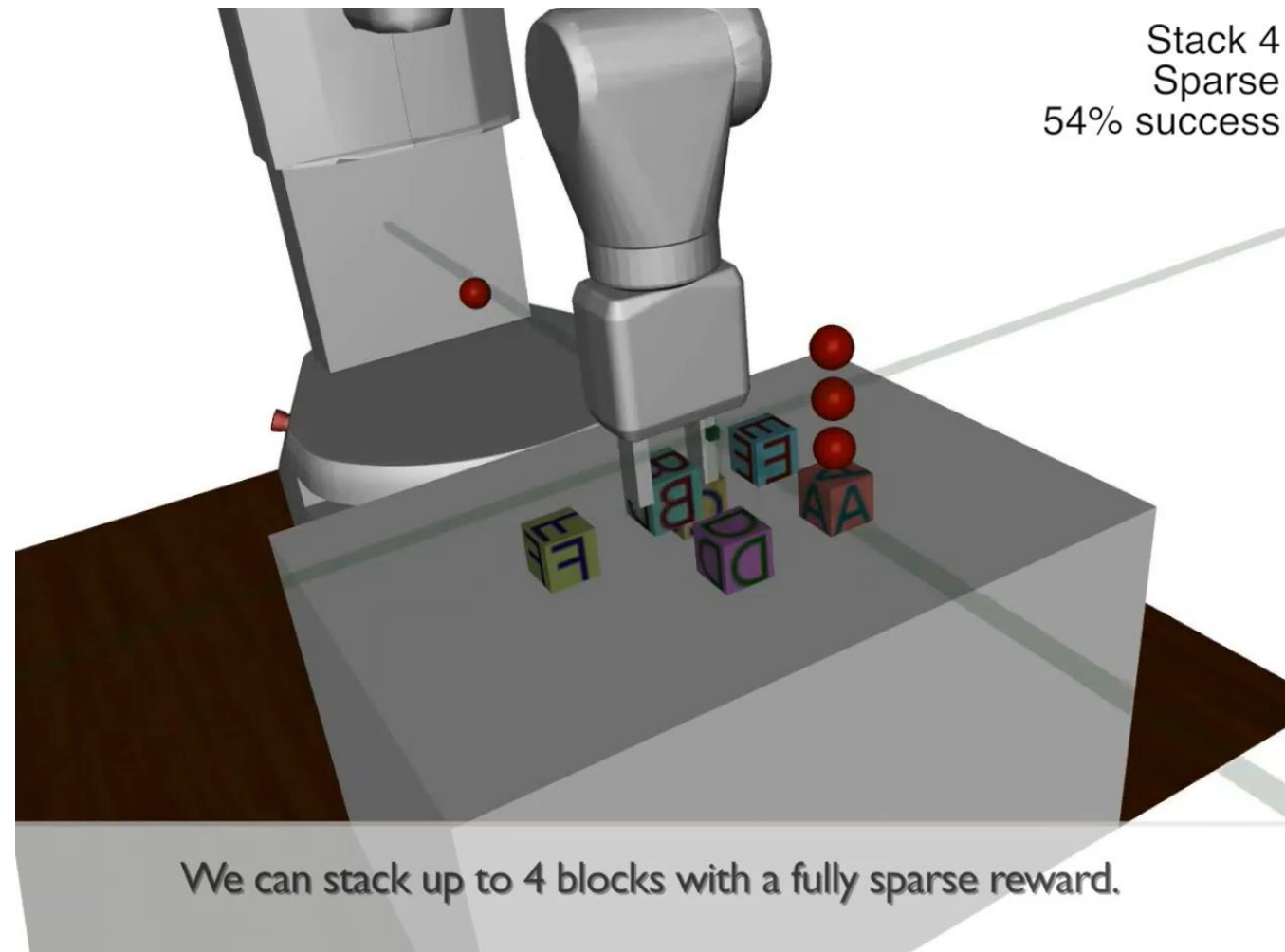
Sparse Rewards

$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{coins} \\ 0 & \text{otherwise} \end{cases}$$

Random exploration can be very inefficient for long horizon tasks



Long Horizon Tasks



Overcoming Exploration in Reinforcement Learning with Demonstrations
[Nair et al. 2018]

Dense vs Sparse Rewards

Dense reward

- ✓ Frequent feedback (faster learning)
- ✓ Easier exploration
- ✗ Shaping bias
- ✗ Harder to design

Sparse reward

- ✗ Infrequent feedback (slower learning)
- ✗ Harder exploration
- ✓ Less shaping bias
- ✓ Easier to design

Dense vs Sparse Rewards

Dense reward

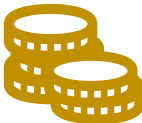
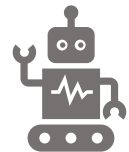
- ✓ Frequent feedback (faster learning)
- ✓ Better exploration
- ✗ Shaping bias
- ✗ Harder to design

Sparse reward

- ✗ Infrequent feedback (slower learning)
- ✗ Harder exploration
- ✓ Less shaping bias
- ✓ Easier to design

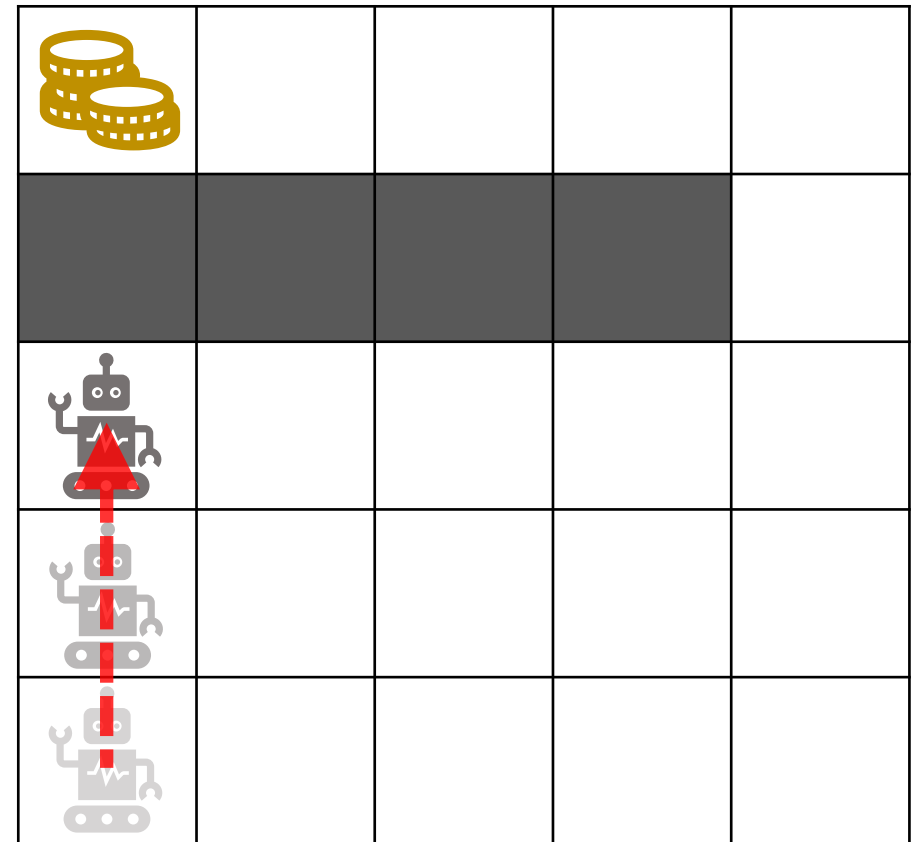
Shaping Bias

$$r_t = - \left\| \text{robot} - \text{coins} \right\|^2$$

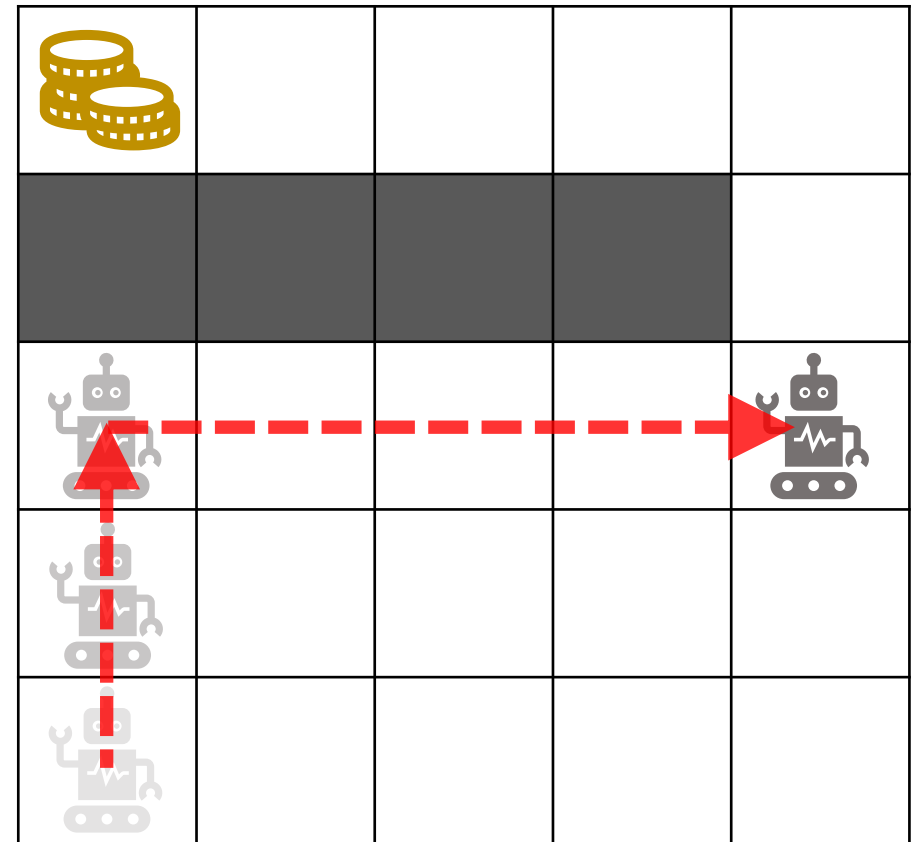
Shaping Bias

$$r_t = - \left\| \text{robot} - \text{coins} \right\|^2$$



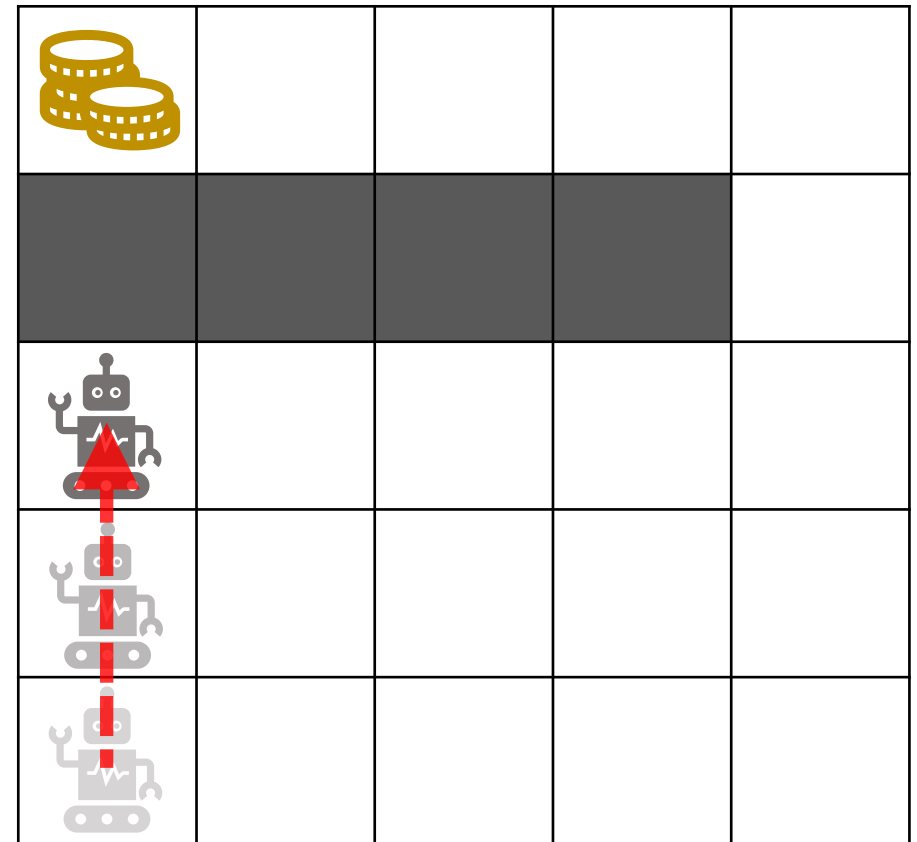
Shaping Bias

$$r_t = - \left\| \text{robot} - \text{coins} \right\|^2$$



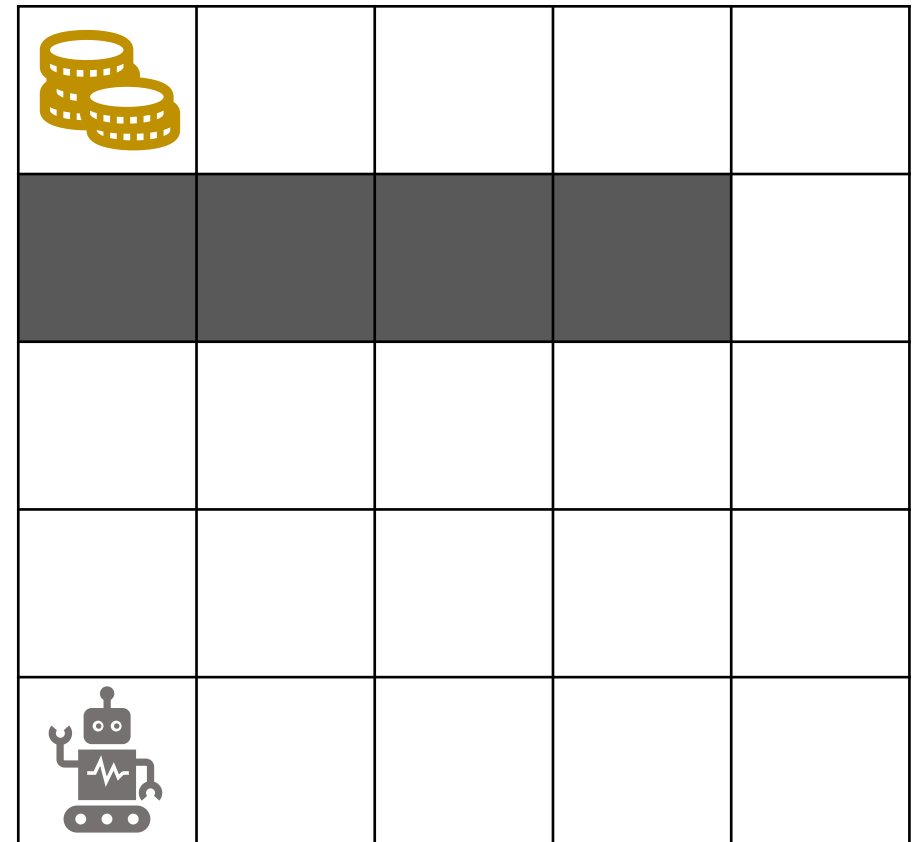
Shaping Bias

$$r_t = - \left\| \text{robot} - \text{coins} \right\|^2$$



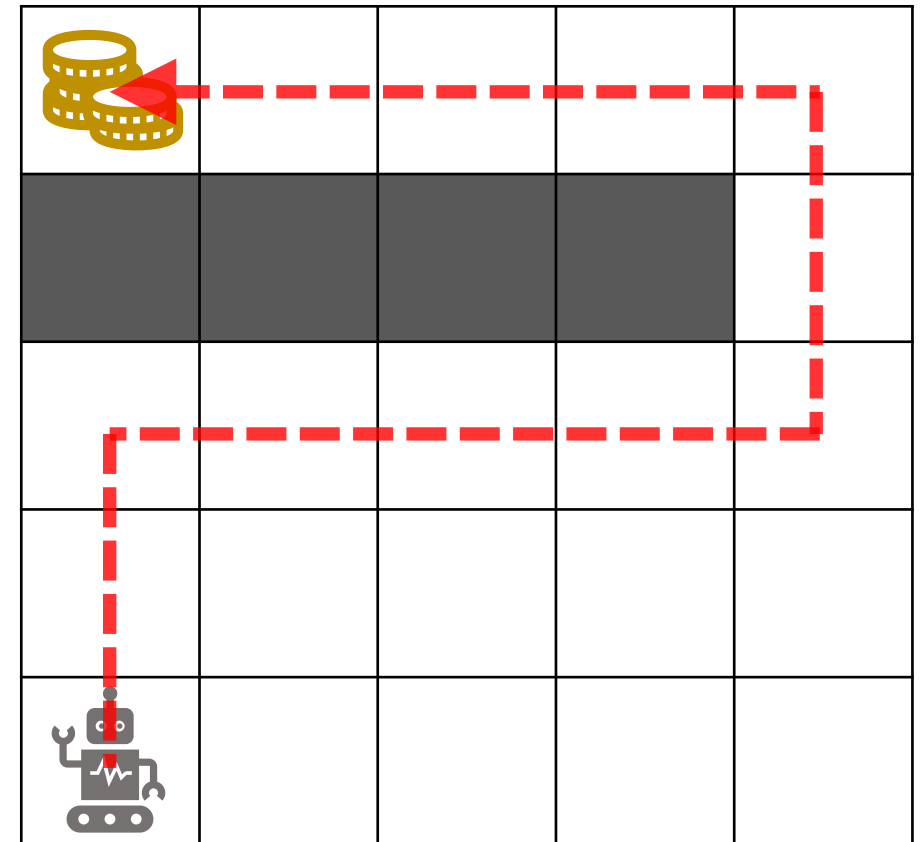
Shaping Bias

$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{coins} \\ 0 & \text{otherwise} \end{cases}$$



Shaping Bias

$$r_t = \begin{cases} 1 & \text{if } \text{robot} \text{ at } \text{coins} \\ 0 & \text{otherwise} \end{cases}$$

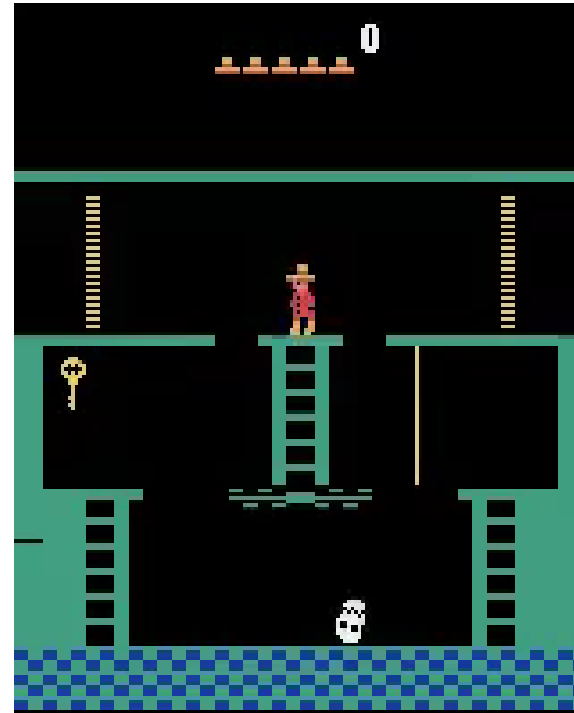


Atari

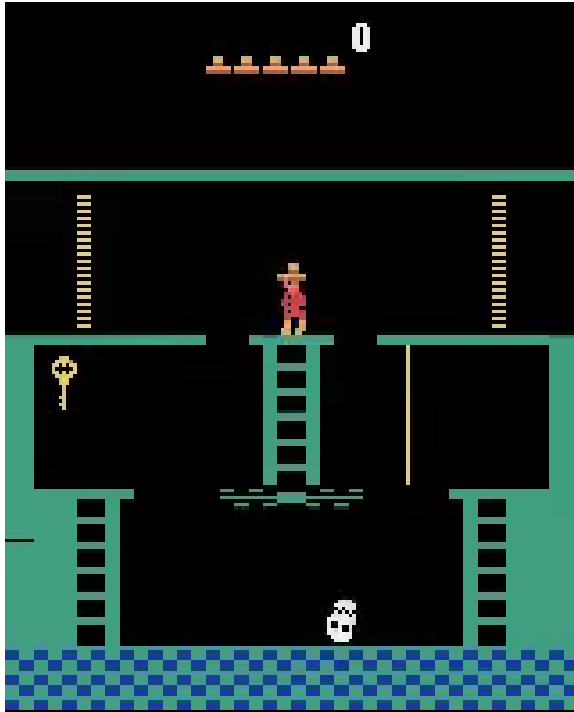
Fairly Easy



Almost Impossible

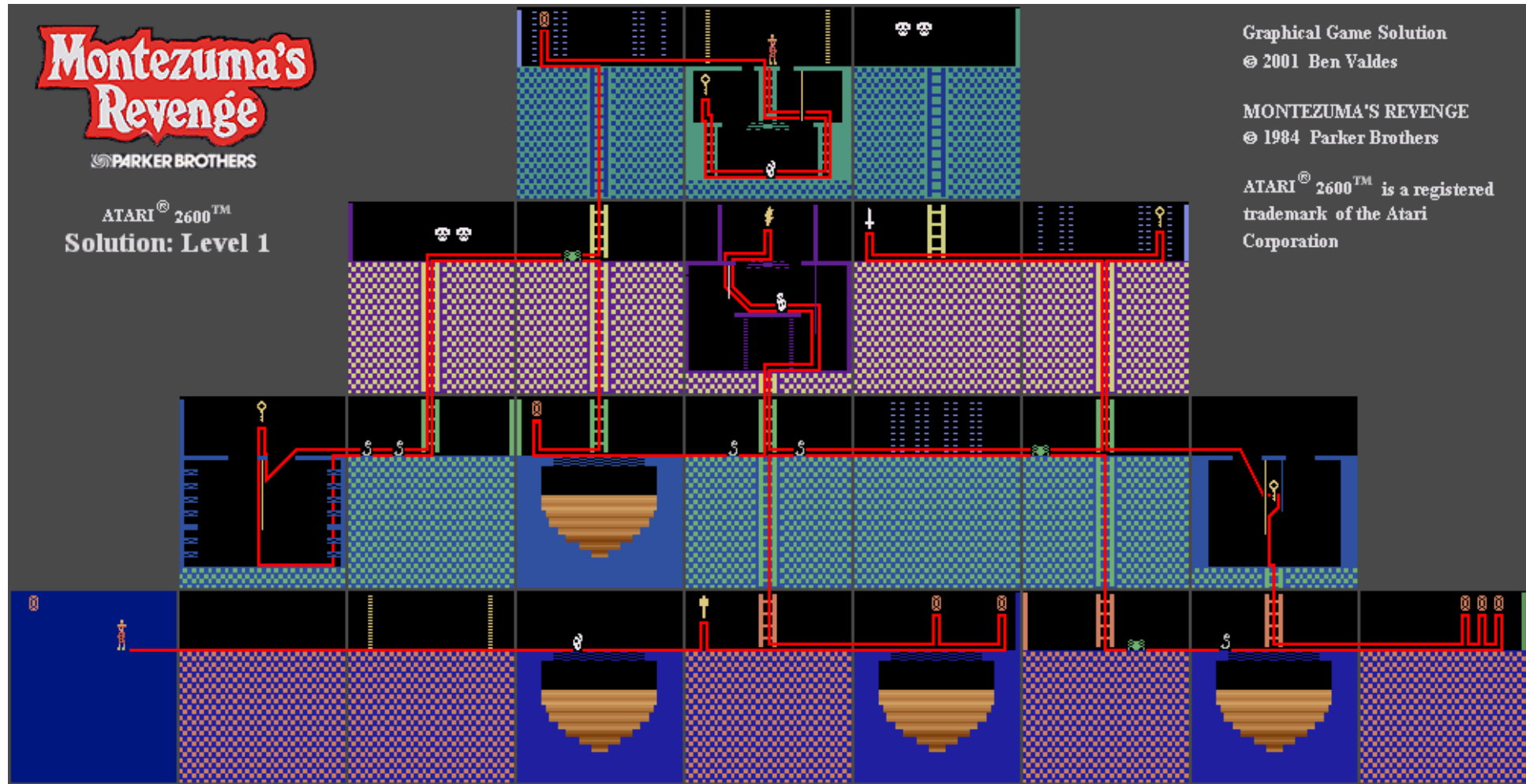


Montezuma's Revenge



- Very sparse reward
- +1: getting key
- +1: opening door
- RL algorithm has no idea what keys and doors are

Montezuma's Revenge

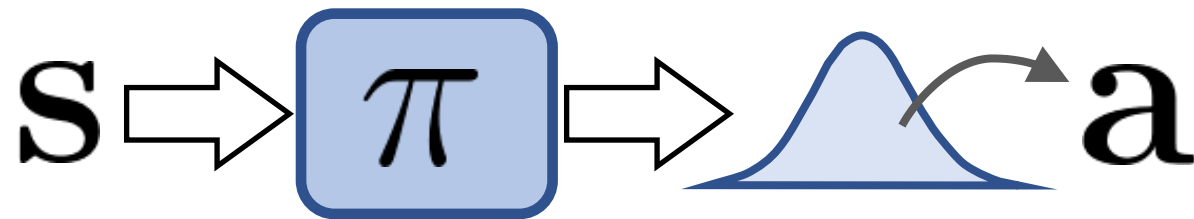


[Graphic Created by [Ben Valdes](#)]

Better Exploration Strategies

- Agent needs to visit new states and try new actions to find optimal strategies
- Encourage coverage of *both* states and actions

relatively easy

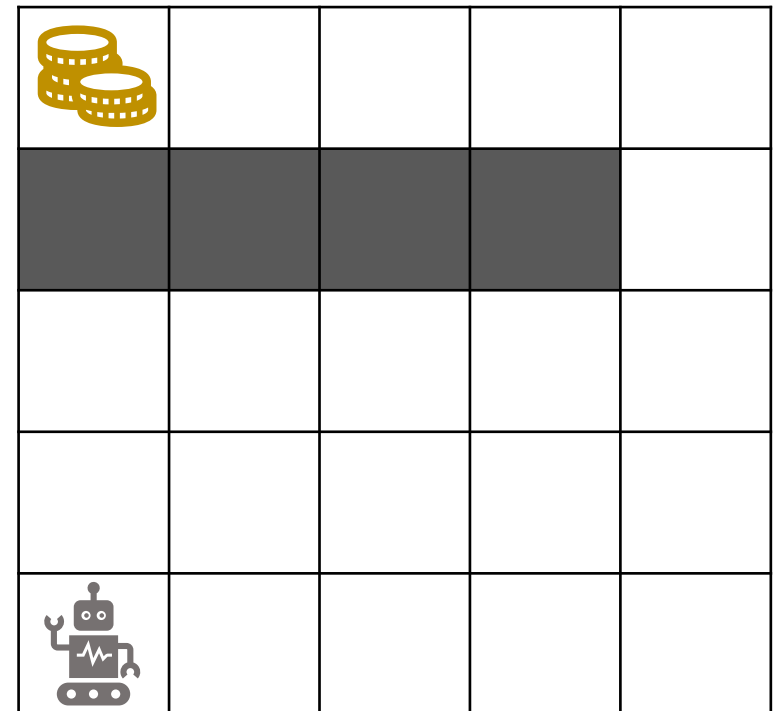


Better Exploration Strategies

- Agent needs to visit new states and try new actions to find optimal strategies
- Encourage coverage of *both* states and actions

much harder

encourage agent to
visit new states



Intrinsic Motivation

r_t

Intrinsic Motivation

$$\hat{r}_t = \underbrace{r_t^{\text{ext}}}_{\text{Extrinsic Reward}} + \beta \underbrace{r_t^{\text{int}}}_{\text{Intrinsic Reward}}$$

Extrinsic reward

- from environment
- encourage agent to perform a given task

Intrinsic reward

- from the agent itself
- encourage agent to explore new states

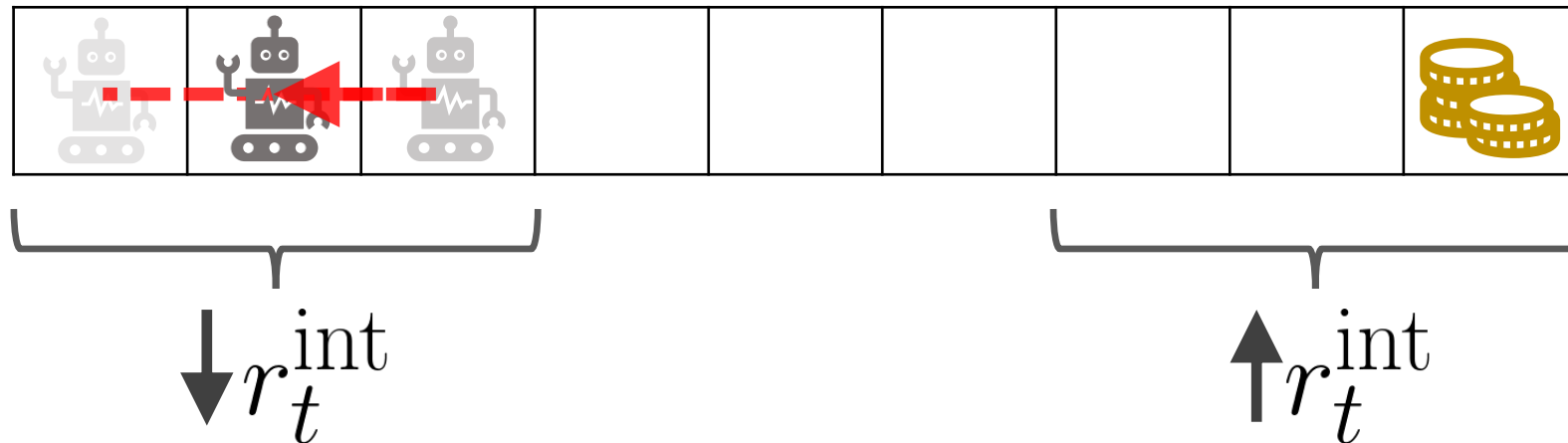
Intrinsic Reward

- Nonstationary reward
- Low for frequently visited states $\downarrow r_t^{\text{int}}$
- High for rarely visited states $\uparrow r_t^{\text{int}}$

Count-Based Exploration

Keep count $N(s)$ on how many times agent visited a particular state

dense reward $r_t^{\text{int}} = \frac{1}{1 + N(s_t)}$

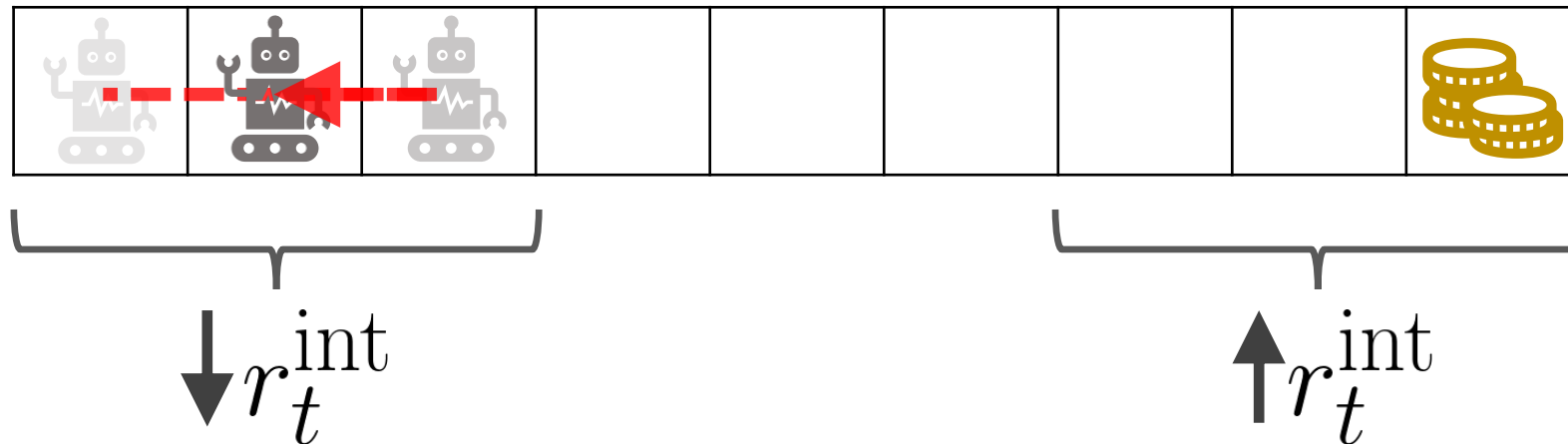


Count-Based Exploration

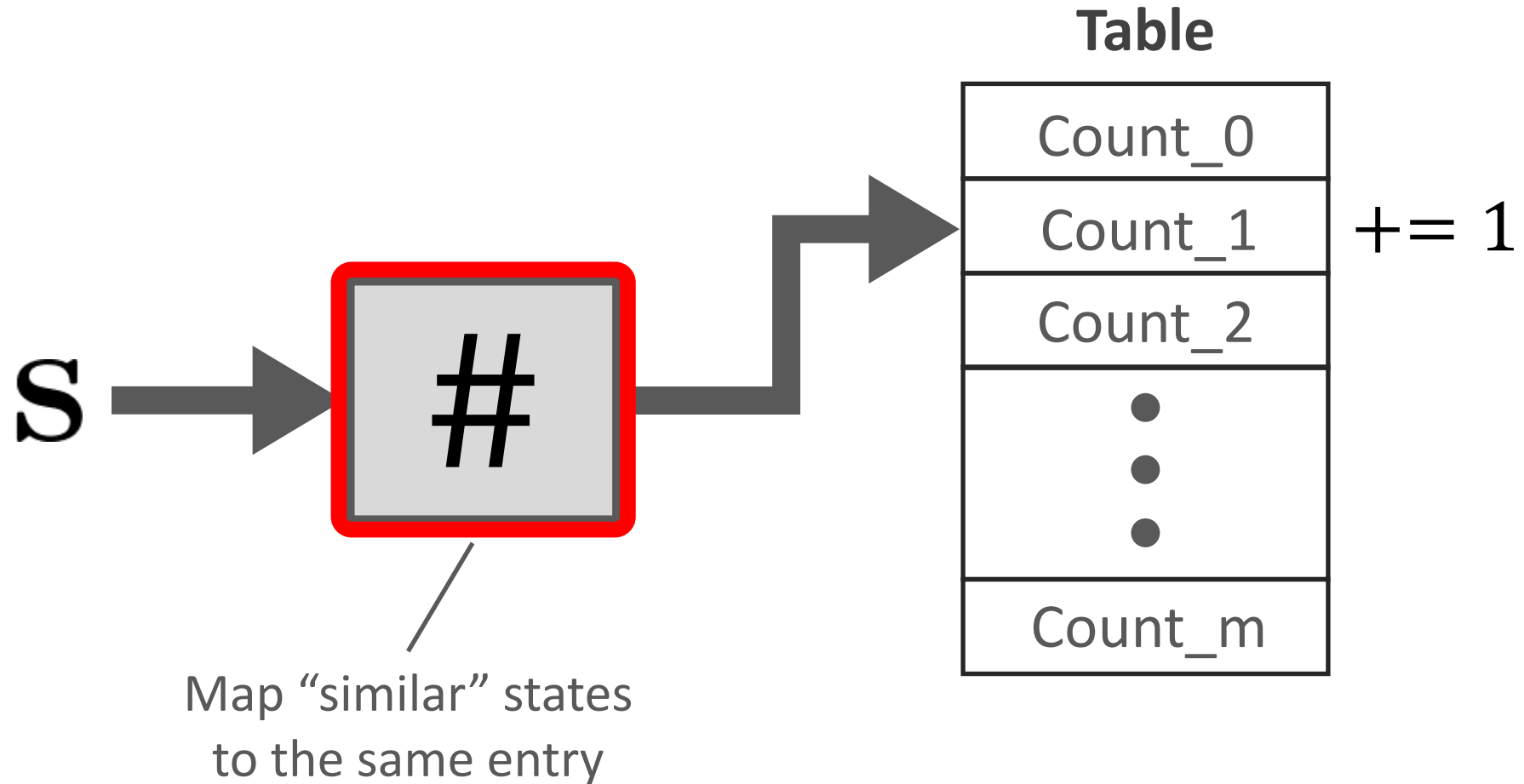
Keep count $N(s)$ on how many times agent visited a particular state

$$r_t^{\text{int}} = \frac{1}{1 + N(s_t)}$$

What about large/
continuous states?

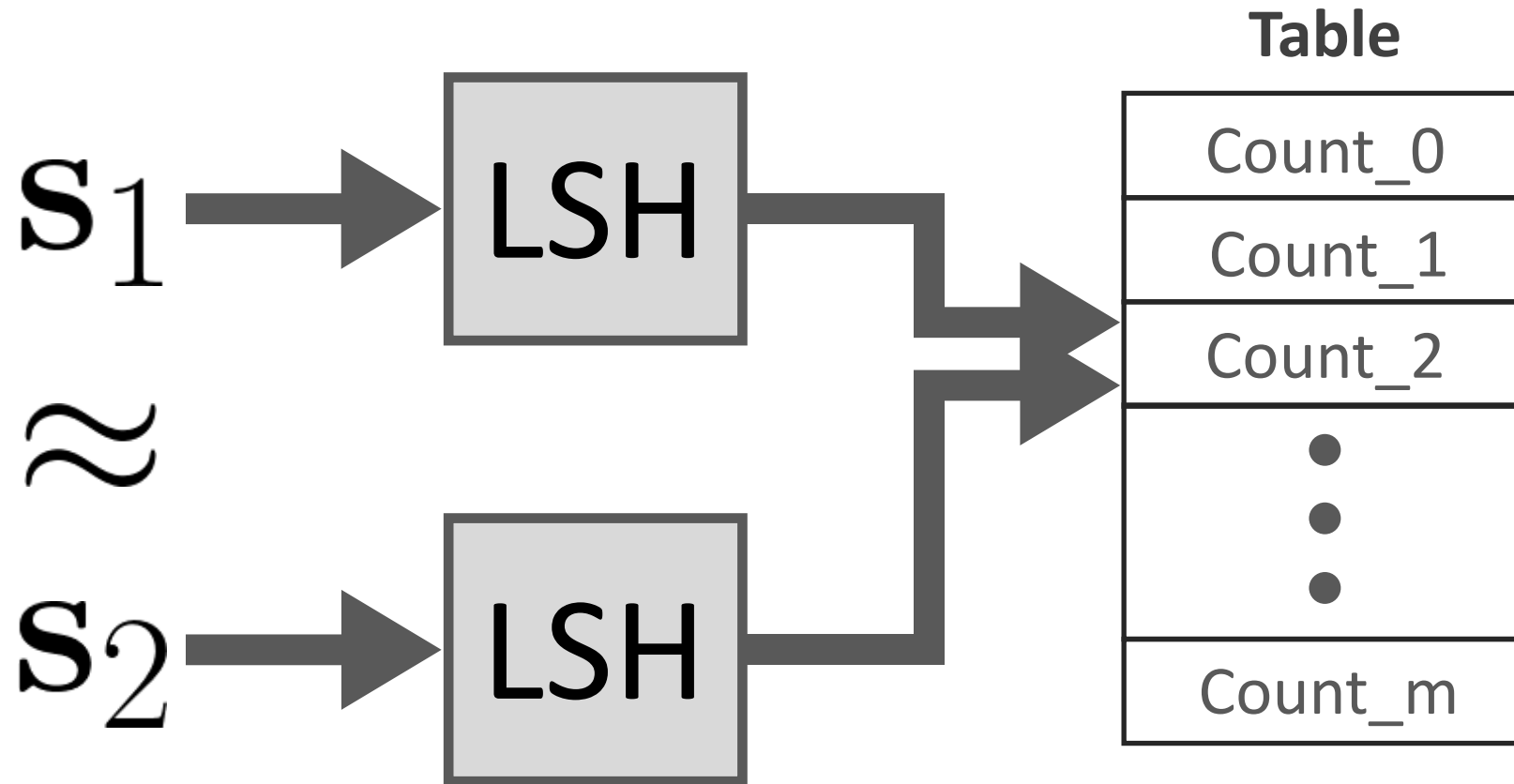


State Hashing

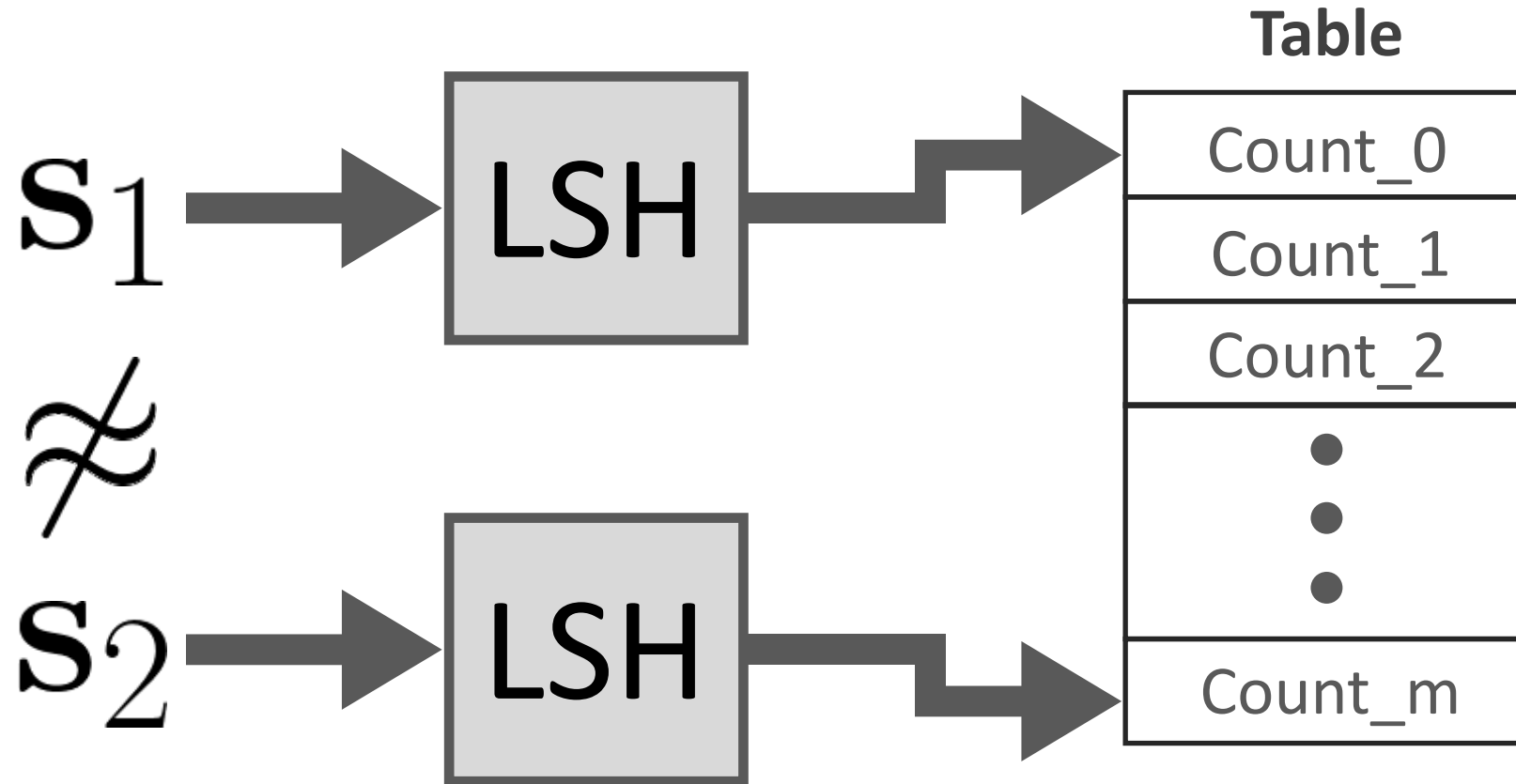


#Exploration: A Study of Count-Based Exploration for Deep Reinforcement Learning
[Tang et al. 2017]

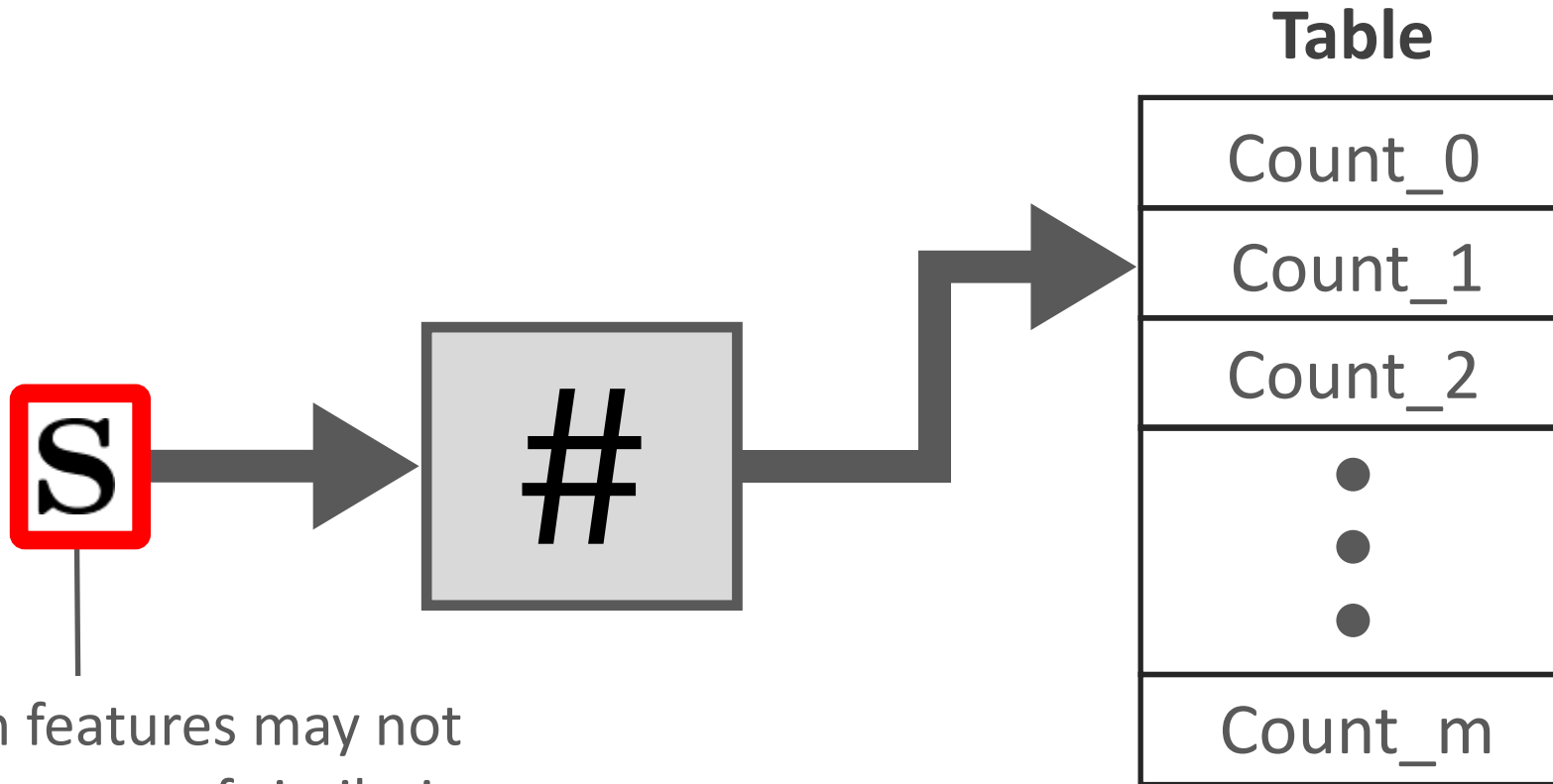
Locality-Sensitive Hashing



Locality-Sensitive Hashing



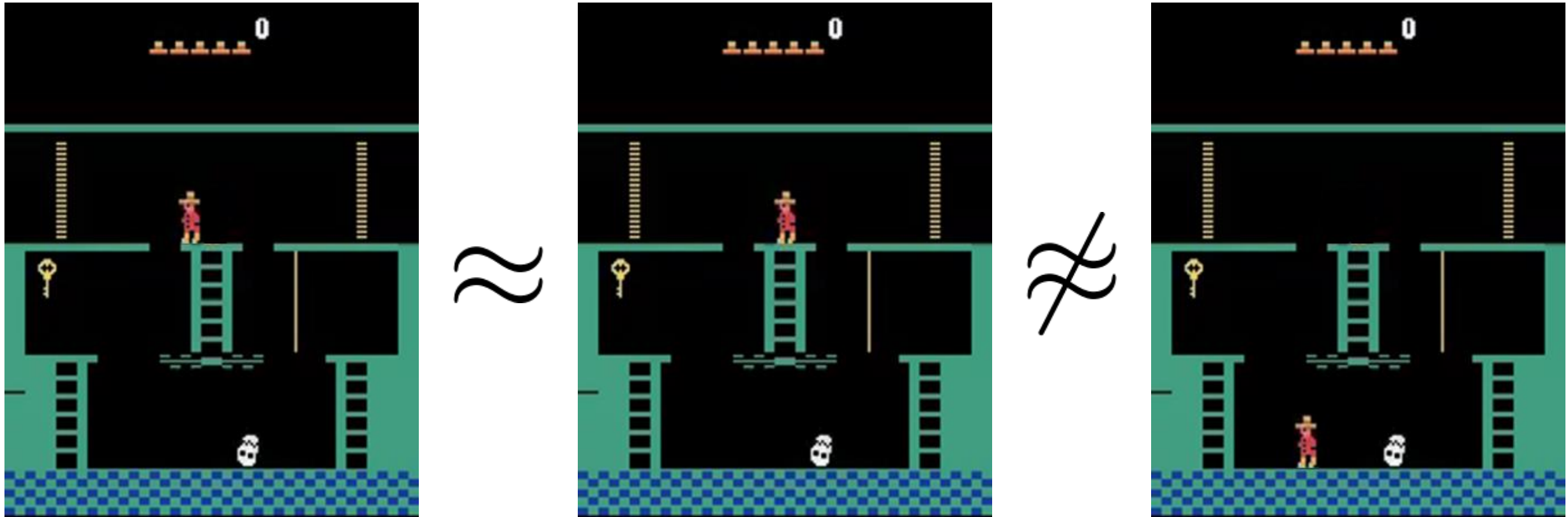
Locality-Sensitive Hashing



Observation features may not be a good measure of similarity

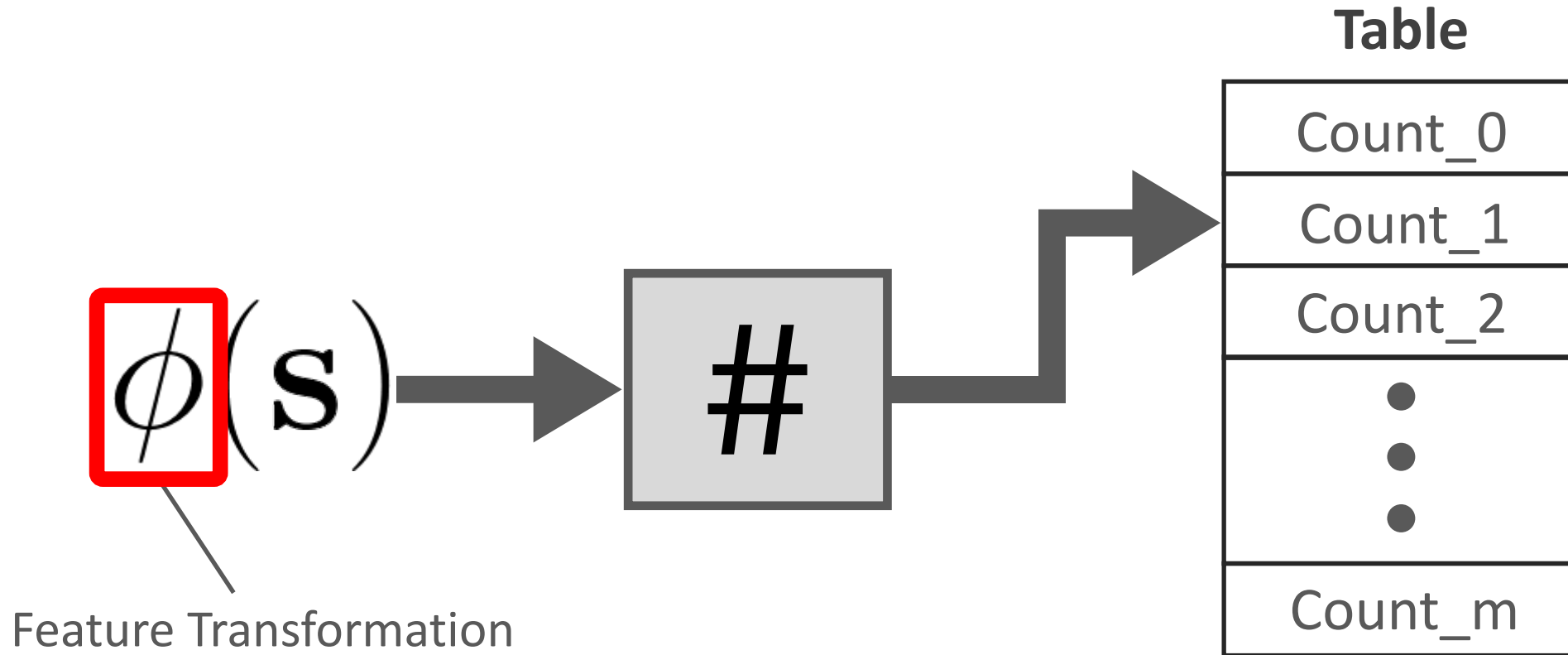
#Exploration: A Study of Count-Based Exploration for Deep Reinforcement Learning
[Tang et al. 2017]

Locality-Sensitive Hashing



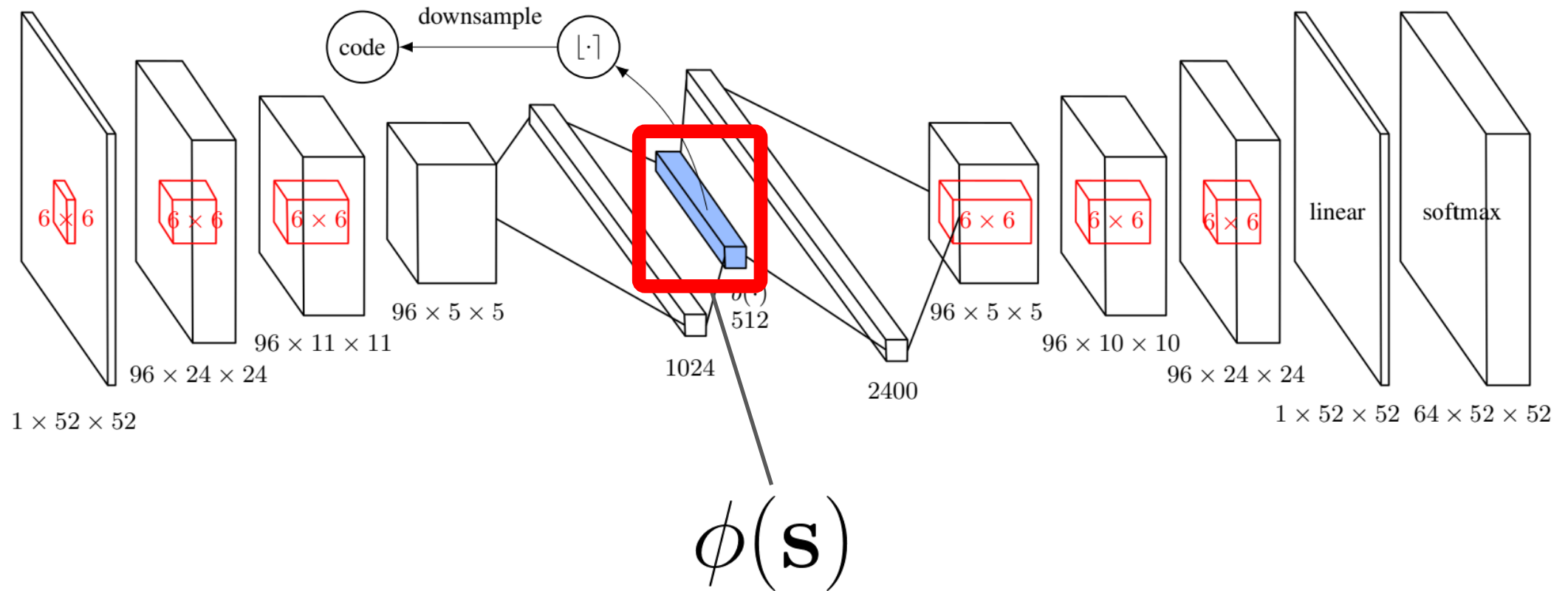
#Exploration: A Study of Count-Based Exploration for Deep Reinforcement Learning
[Tang et al. 2017]

Locality-Sensitive Hashing



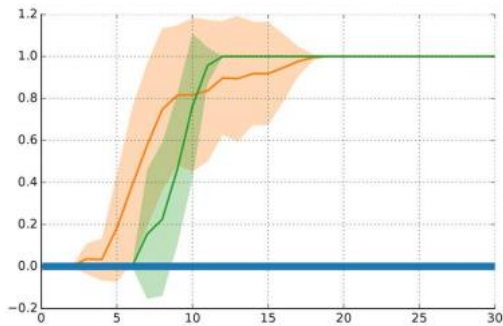
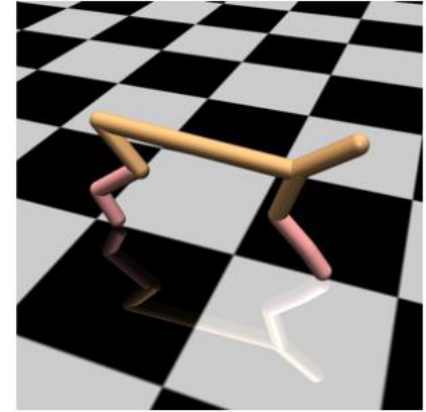
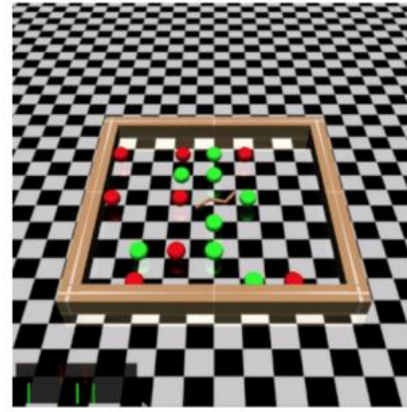
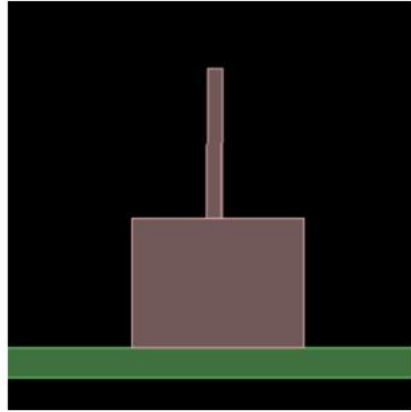
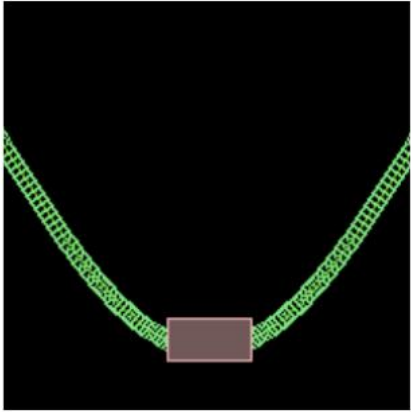
#Exploration: A Study of Count-Based Exploration for Deep Reinforcement Learning
[Tang et al. 2017]

Feature Embedding

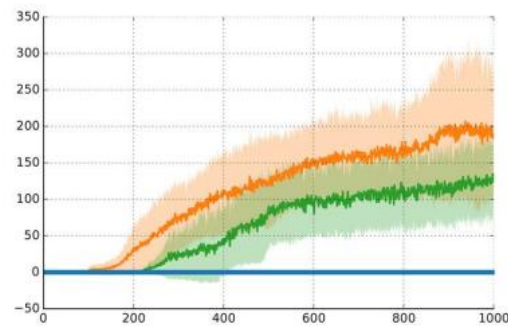


#Exploration: A Study of Count-Based Exploration for Deep Reinforcement Learning
[Tang et al. 2017]

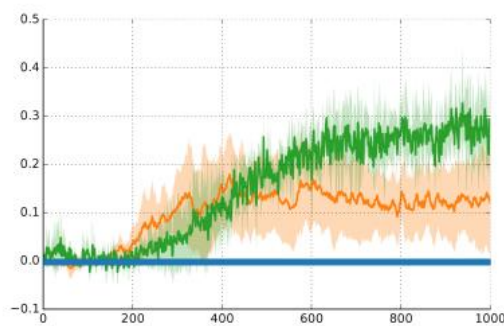
State Hashing



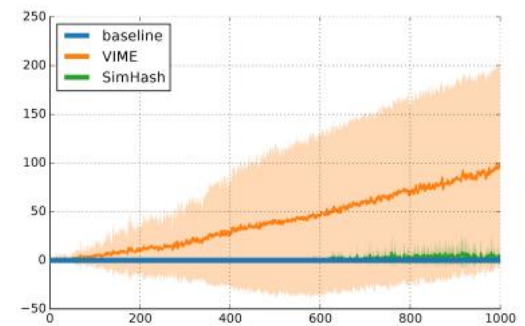
(a) MountainCar



(b) CartPoleSwingup



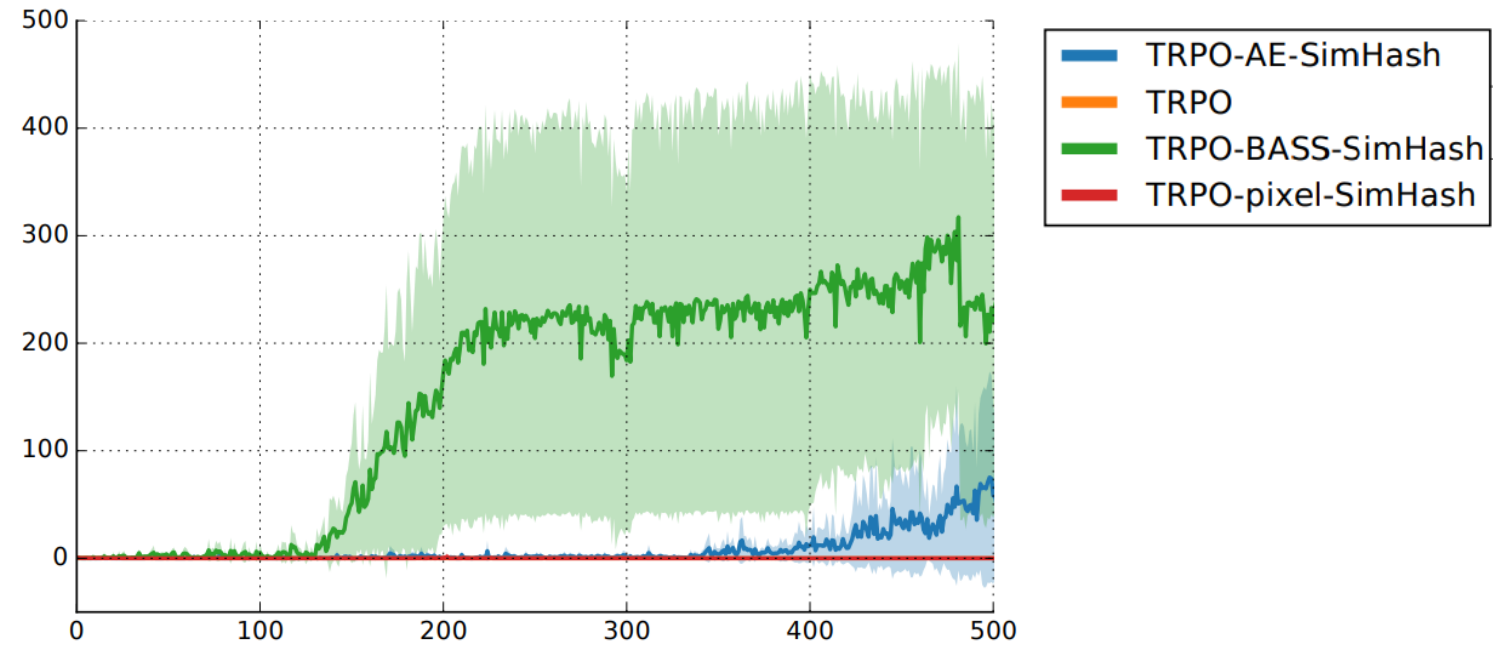
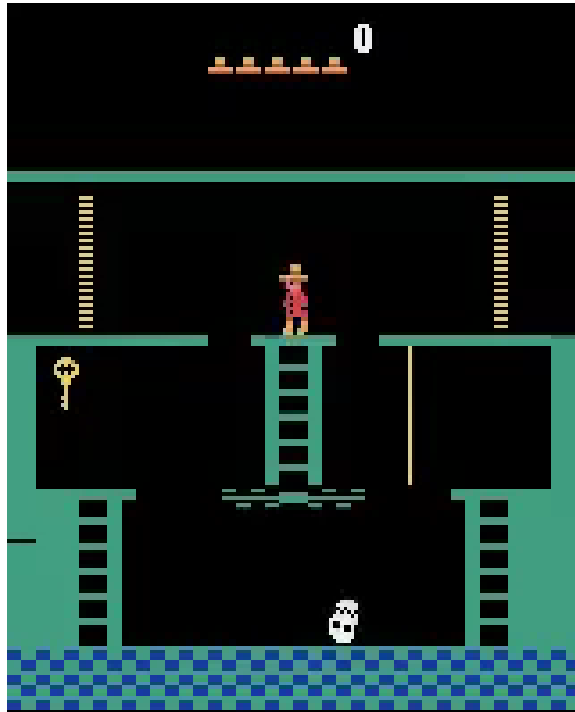
(c) SwimmerGather



(d) HalfCheetah

#Exploration: A Study of Count-Based Exploration for Deep Reinforcement Learning
[Tang et al. 2017]

State Hashing



(d) Montezuma's Revenge

#Exploration: A Study of Count-Based Exploration for Deep Reinforcement Learning
[Tang et al. 2017]

State Hashing (Drawbacks)

- Learning an effective representation for hashing can be difficult
- Prone to aliasing
- Distribution of states changes during training (feature transform needs to be updated)
- Hard to pick hash table size a-priori

Density Estimation

$$p(\mathbf{s}) = \frac{N(\mathbf{s})}{n}$$

probability / density

count

total timesteps

Density Estimation

$$p(\mathbf{s}) = \frac{N(\mathbf{s})}{n}$$

Idea: use density to estimate count

Density Estimation

Agent visits a state \mathbf{s}

Before

$$p(\mathbf{s}) = \frac{N(\mathbf{s})}{n}$$

After

$$p'(\mathbf{s}) = \frac{N(\mathbf{s}) + 1}{n + 1}$$

2 equations and 2 unknowns ($N(\mathbf{s}), n$)

Density Estimation

Solve for $N(\mathbf{s})$

$$p(\mathbf{s}) = \frac{N(\mathbf{s})}{n}$$


$$p'(\mathbf{s}) = \frac{N(\mathbf{s}) + 1}{n + 1}$$

Density Estimation

Solve for $N(\mathbf{s})$

$$p(\mathbf{s}) = \frac{N(\mathbf{s})}{n}$$

$$p'(\mathbf{s}) = \frac{N(\mathbf{s}) + 1}{n + 1}$$

$$N(\mathbf{s}) = np(\mathbf{s}) \leftarrow$$


Density Estimation

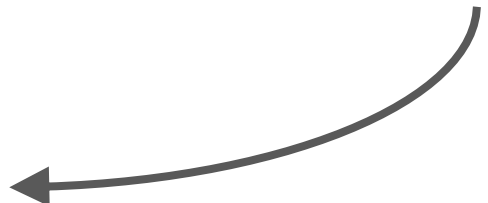
Solve for $N(\mathbf{s})$

$$p(\mathbf{s}) = \frac{N(\mathbf{s})}{n}$$

$$p'(\mathbf{s}) = \frac{N(\mathbf{s}) + 1}{n + 1}$$

$$N(\mathbf{s}) = np(\mathbf{s})$$

$$np'(\mathbf{s}) + p'(\mathbf{s}) = N(\mathbf{s}) + 1$$



Density Estimation

Solve for $N(\mathbf{s})$

$$p(\mathbf{s}) = \frac{N(\mathbf{s})}{n} \qquad p'(\mathbf{s}) = \frac{N(\mathbf{s}) + 1}{n + 1}$$


$$N(\mathbf{s}) = np(\mathbf{s})$$

$$np'(\mathbf{s}) + p'(\mathbf{s}) = \underline{N(\mathbf{s})} + 1 \\ = np(\mathbf{s})$$

Density Estimation

Solve for $N(\mathbf{s})$

$$p(\mathbf{s}) = \frac{N(\mathbf{s})}{n} \qquad p'(\mathbf{s}) = \frac{N(\mathbf{s}) + 1}{n + 1}$$


$$N(\mathbf{s}) = np(\mathbf{s})$$

$$np'(\mathbf{s}) + p'(\mathbf{s}) = N(\mathbf{s}) + 1$$

$$np'(\mathbf{s}) + p'(\mathbf{s}) = np(\mathbf{s}) + 1$$

$$n(p'(\mathbf{s}) - p(\mathbf{s})) = 1 - p'(\mathbf{s})$$


$$n = \frac{1 - p'(\mathbf{s})}{p'(\mathbf{s}) - p(\mathbf{s})}$$

Density Estimation

Solve for $N(\mathbf{s})$

$$p(\mathbf{s}) = \frac{N(\mathbf{s})}{n}$$

$$p'(\mathbf{s}) = \frac{N(\mathbf{s}) + 1}{n + 1}$$


$$N(\mathbf{s}) = np(\mathbf{s})$$

$$np'(\mathbf{s}) + p'(\mathbf{s}) = N(\mathbf{s}) + 1$$

$$np'(\mathbf{s}) + p'(\mathbf{s}) = np(\mathbf{s}) + 1$$

$$n(p'(\mathbf{s}) - p(\mathbf{s})) = 1 - p'(\mathbf{s})$$

$$n = \frac{1 - p'(\mathbf{s})}{p'(\mathbf{s}) - p(\mathbf{s})}$$

$$N(\mathbf{s}) = \frac{1 - p'(\mathbf{s})}{p'(\mathbf{s}) - p(\mathbf{s})}p(\mathbf{s})$$

Density Estimation

$$N(\mathbf{s}) = \frac{1 - p'(\mathbf{s})}{p'(\mathbf{s}) - p(\mathbf{s})} p(\mathbf{s})$$

Pseudo-Count

$$\hat{N}(\mathbf{s}) = \frac{1 - p'(\mathbf{s})}{p'(\mathbf{s}) - p(\mathbf{s})} p(\mathbf{s})$$

“pseudo-count”

Pseudo-Count

$$\hat{N}(\mathbf{s}) = \frac{1 - p'(\mathbf{s})}{p'(\mathbf{s}) - p(\mathbf{s})} p(\mathbf{s})$$

How?

Pseudo-Count

$$\hat{N}(\mathbf{s}) = \frac{1 - p'(\mathbf{s})}{p'(\mathbf{s}) - p(\mathbf{s})} p(\mathbf{s})$$

Fit density model

- E.g. flow models, autoregressive models, CTS (Bellemare et al. 2016), etc.

$$\rho(\mathbf{s}) \approx p(\mathbf{s}) \quad \rho'(\mathbf{s}) \approx p'(\mathbf{s})$$

Exploration Pseudo-Count

ALGORITHM: Exploration with Pseudo Counts

- 1: $\mathcal{D} \leftarrow$ initialize dataset
 - 2: Fit density model $\rho(\mathbf{s})$ to \mathcal{D}

 - 3: **for** every timestep t **do**
 - 4: Observe state \mathbf{s}_t
 - 5: Sample action from policy $\mathbf{a}_t \sim \pi(\mathbf{a}_t|\mathbf{s}_t)$
 - 6: Apply \mathbf{a}_t and observe new state \mathbf{s}_{t+1} and extrinsic reward r_t^{ext}
 - 7: Store \mathbf{s}_{t+1} in \mathcal{D}

 - 8: Fit density model $\rho'(\mathbf{s})$ to \mathcal{D}
 - 9: Use $\rho(\mathbf{s}_{t+1})$ and $\rho'(\mathbf{s}_{t+1})$ to estimate pseudo-count $\hat{N}(\mathbf{s}_{t+1})$
 - 10: Calculate intrinsic reward r_t^{int} with $N(\mathbf{s}_{t+1})$
 - 11: Calculate reward $r_t = r_t^{\text{ext}} + \beta r_t^{\text{int}}$

 - 12: $\rho \leftarrow \rho'$
 - 13: **end for**
-

$$r_t^{\text{int}} = \frac{1}{1 + N(\mathbf{s}_{t+1})}$$

Exploration Pseudo-Count

ALGORITHM: Exploration with Pseudo Counts

- 1: $\mathcal{D} \leftarrow$ initialize dataset
 - 2: Fit density model $\rho(\mathbf{s})$ to \mathcal{D}

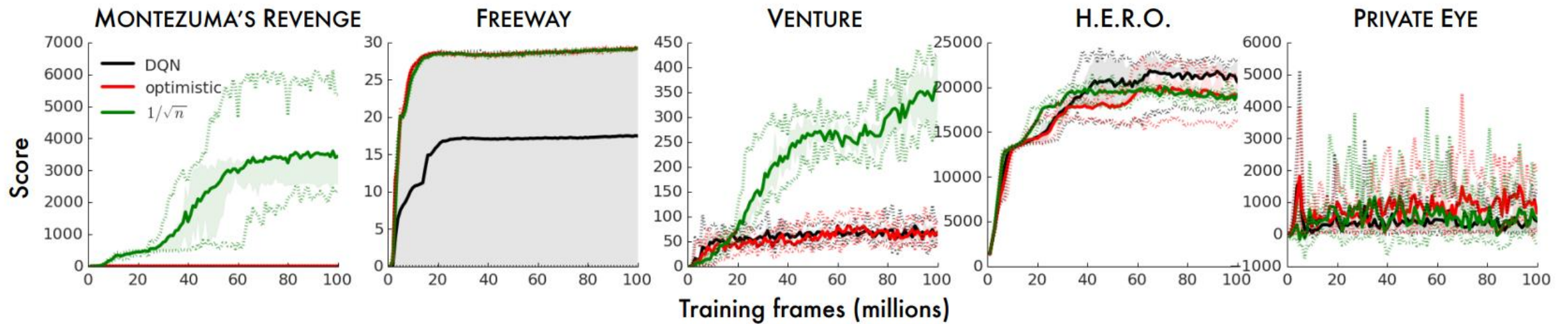
 - 3: **for** every timestep t **do**
 - 4: Observe state \mathbf{s}_t
 - 5: Sample action from policy $\mathbf{a}_t \sim \pi(\mathbf{a}_t|\mathbf{s}_t)$
 - 6: Apply \mathbf{a}_t and observe new state \mathbf{s}_{t+1} and extrinsic reward r_t^{ext}
 - 7: Store \mathbf{s}_{t+1} in \mathcal{D}

 - 8: Fit density model $\rho'(\mathbf{s})$ to \mathcal{D}
 - 9: Use $\rho(\mathbf{s}_{t+1})$ and $\rho'(\mathbf{s}_{t+1})$ to estimate pseudo-count $\hat{N}(\mathbf{s}_{t+1})$
 - 10: Calculate intrinsic reward r_t^{int} with $N(\mathbf{s}_{t+1})$
 - 11: Calculate reward $r_t = r_t^{\text{ext}} + \beta r_t^{\text{int}}$

 - 12: $\rho \leftarrow \rho'$
 - 13: **end for**
-

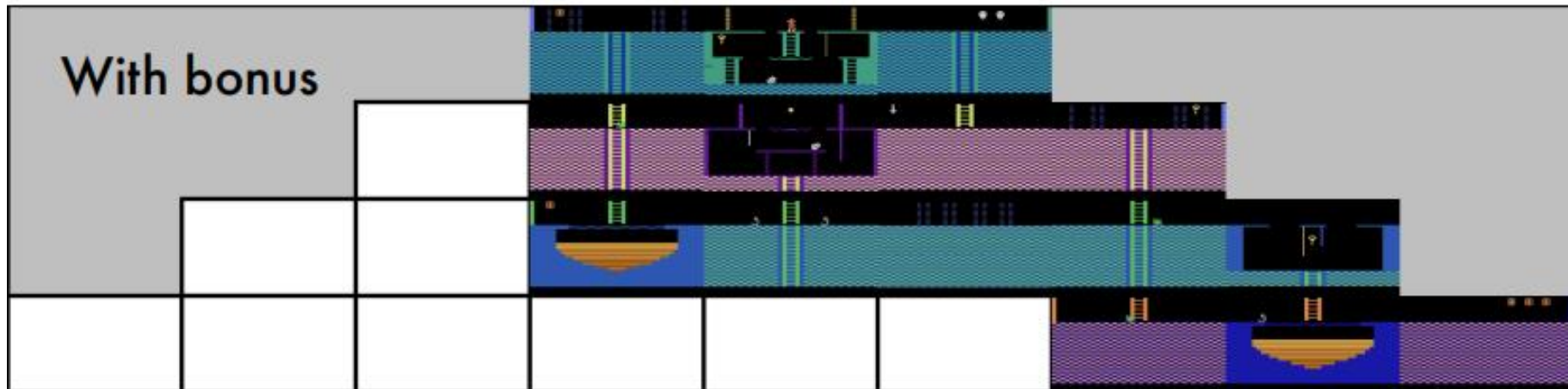
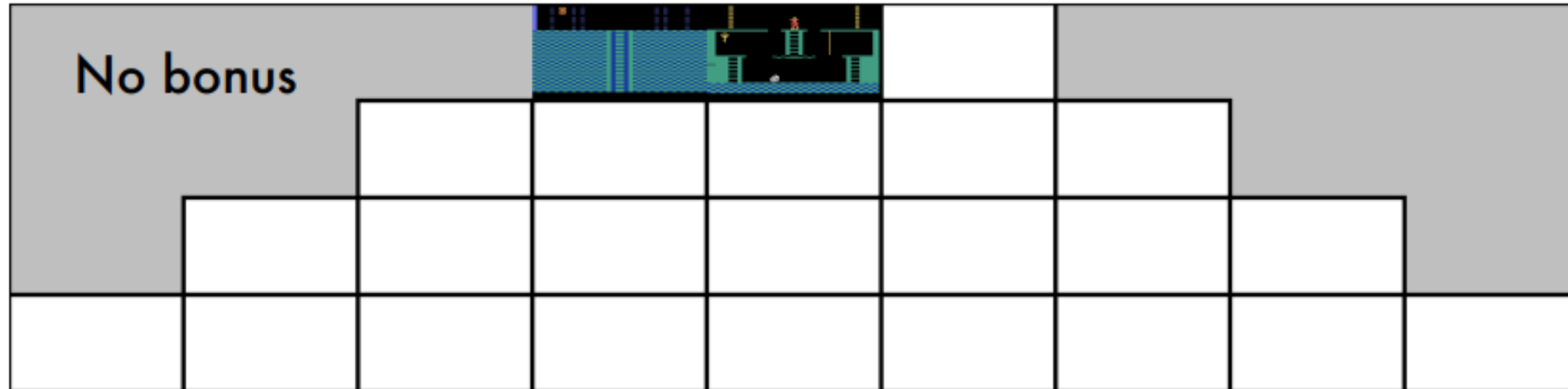
$$r_t^{\text{int}} = \frac{1}{\sqrt{0.01 + N(\mathbf{s}_{t+1})}}$$

Exploration Pseudo-Count




Unifying Count-Based Exploration and Intrinsic Motivation
[Bellemare et al. 2016]

Exploration Pseudo-Count



Unifying Count-Based Exploration and Intrinsic Motivation
[Bellemare et al. 2016]

Count-Based Reward

$$r_t^{\text{int}} = \frac{1}{1 + N(\mathbf{s}_{t+1})}$$


$\rho(\mathbf{s})$

Entropy-Based Reward

$$r_t^{\text{int}} = -\log(\rho(\mathbf{s}_{t+1}))$$

maximize state entropy

$$\mathcal{H}(\mathbf{s})$$

No need to estimate likelihood before and after state visitation

$$\rho(\mathbf{s}) \approx p(\mathbf{s}) \quad \rho'(\mathbf{s}) \approx p'(\mathbf{s})$$

Exploration Pseudo-Count

ALGORITHM: Exploration with Pseudo Counts

- 1: $\mathcal{D} \leftarrow$ initialize dataset
 - 2: Fit density model $\rho(\mathbf{s})$ to \mathcal{D}

 - 3: **for** every timestep t **do**
 - 4: Observe state \mathbf{s}_t
 - 5: Sample action from policy $\mathbf{a}_t \sim \pi(\mathbf{a}_t|\mathbf{s}_t)$
 - 6: Apply \mathbf{a}_t and observe new state \mathbf{s}_{t+1} and extrinsic reward r_t^{ext}
 - 7: Store \mathbf{s}_{t+1} in \mathcal{D}
 - 8: Fit density model $\rho'(\mathbf{s})$ to \mathcal{D}
 - 9: Use $\rho(\mathbf{s}_{t+1})$ and $\rho'(\mathbf{s}_{t+1})$ to estimate pseudo-count $\hat{N}(\mathbf{s}_{t+1})$
 - 10: Calculate intrinsic reward r_t^{int} with $\hat{N}(\mathbf{s}_{t+1})$
 - 11: Calculate reward $r_t = r_t^{\text{ext}} + \beta r_t^{\text{int}}$

 - 12: $\rho \leftarrow \rho'$
 - 13: **end for**
-

Density estimation
is hard!

Surprise Maximization



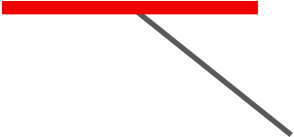
[[Sikana](#)]



[[Science Channel](#)]

Surprise Maximization

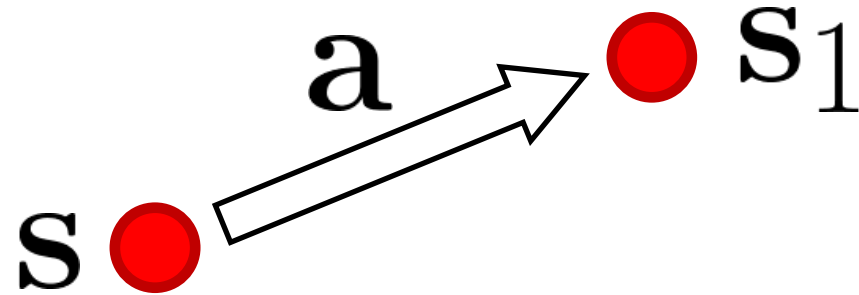
- Use surprise as a proxy for novelty



How do we measure
surprise?

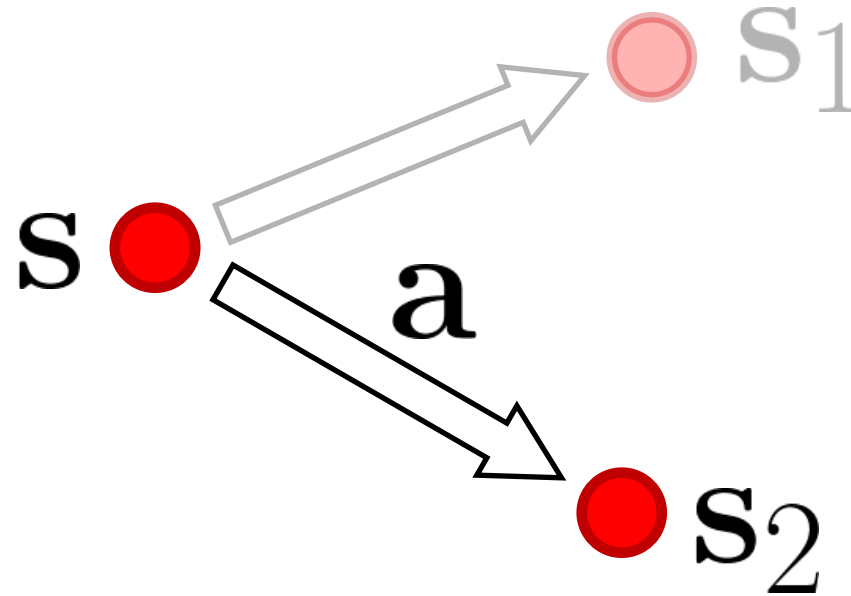
Surprise Maximization

- Use surprise as a proxy for novelty
- Detect surprise via prediction error



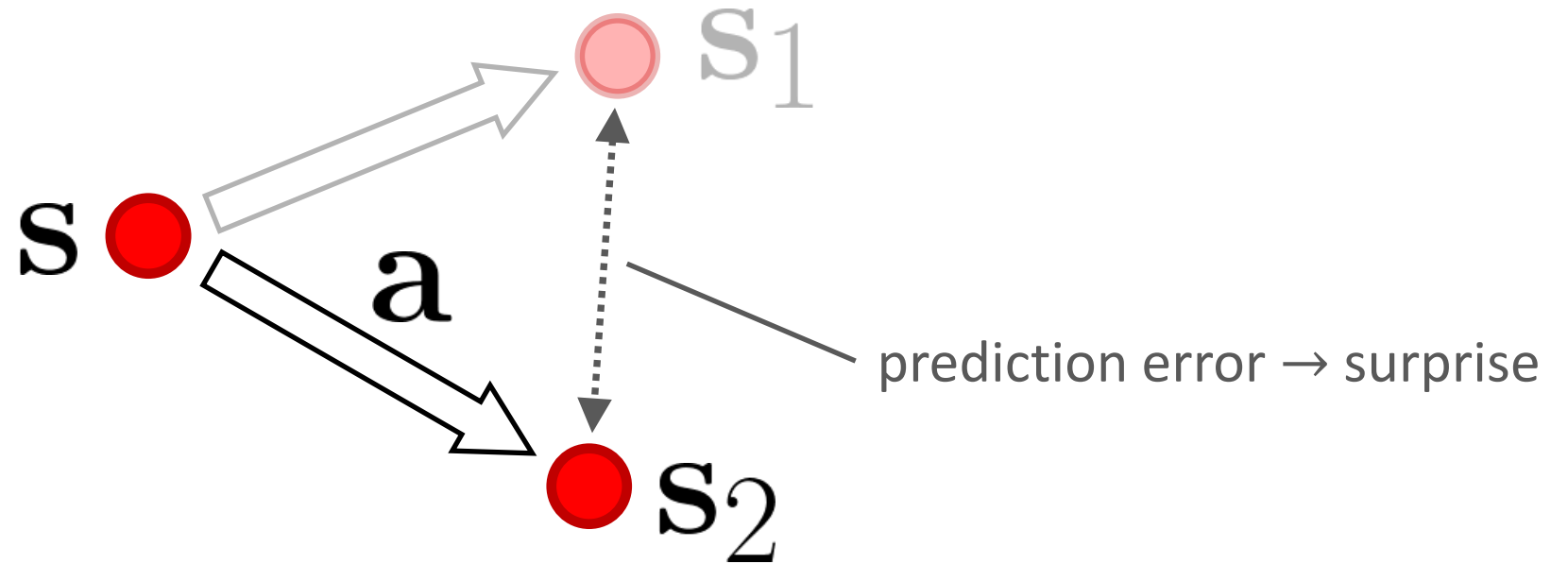
Surprise Maximization

- Use surprise as a proxy for novelty
- Detect surprise via prediction error



Surprise Maximization

- Use surprise as a proxy for novelty
- Detect surprise via prediction error



Dynamics Model

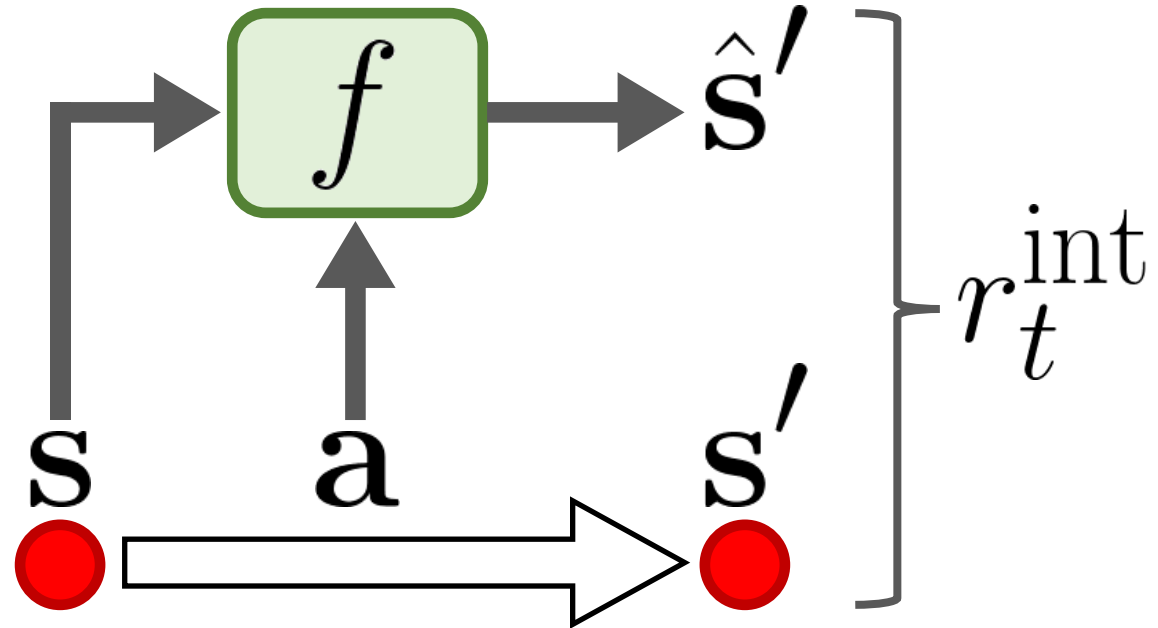
Detect surprise using a dynamics model

$$f(\mathbf{s}' | \mathbf{s}, \mathbf{a})$$

Intrinsic reward maximizes prediction error

$$r_t^{\text{int}} = -\log f(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

Dynamics Model

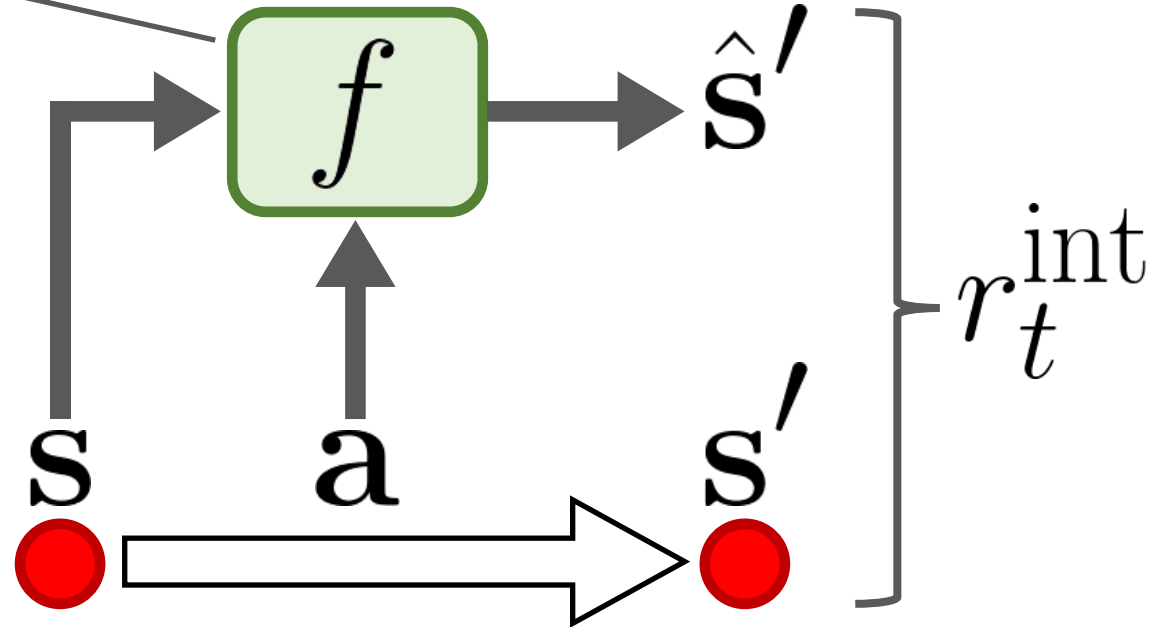


Surprise-Based Intrinsic Motivation for Deep Reinforcement Learning
[Achiam and Sastry 2017]

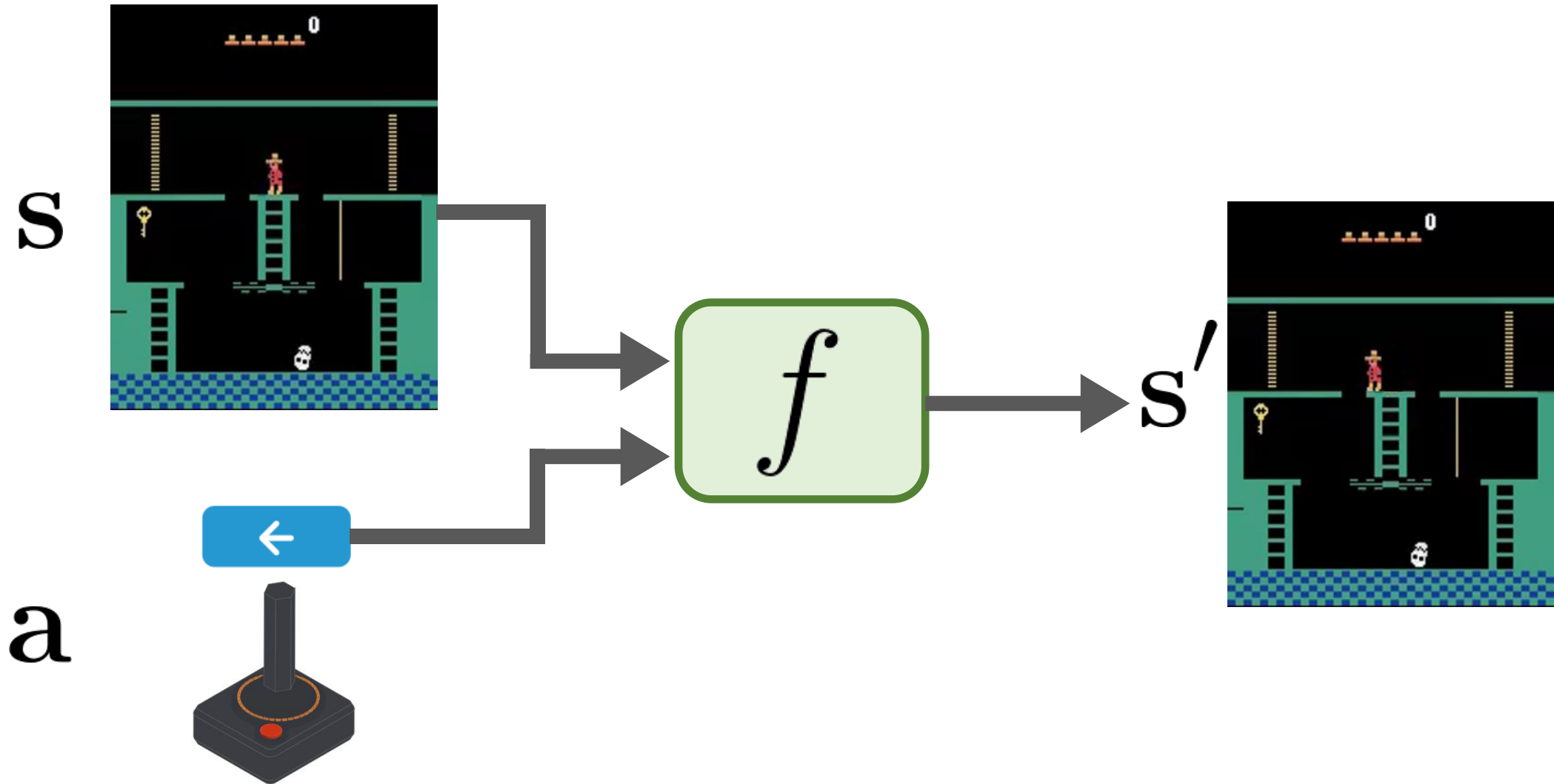
Dynamics Model

Trained with data
from policy

$\{(s_i, a_i, s'_i)\}$



Dynamics Model



Dynamics Models

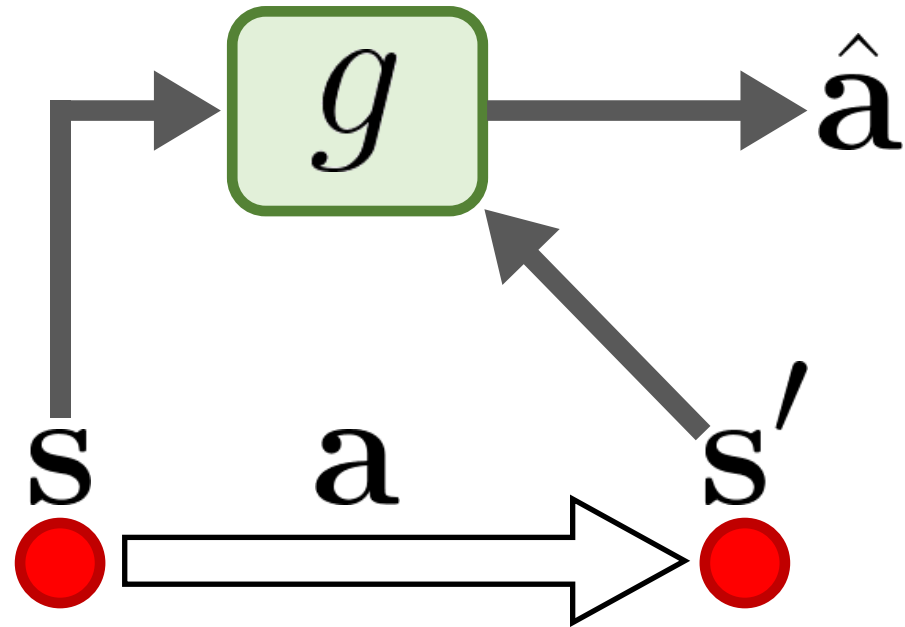
Forward dynamics model:

$$f(\mathbf{s}' | \mathbf{s}, \mathbf{a})$$

Inverse dynamics model:

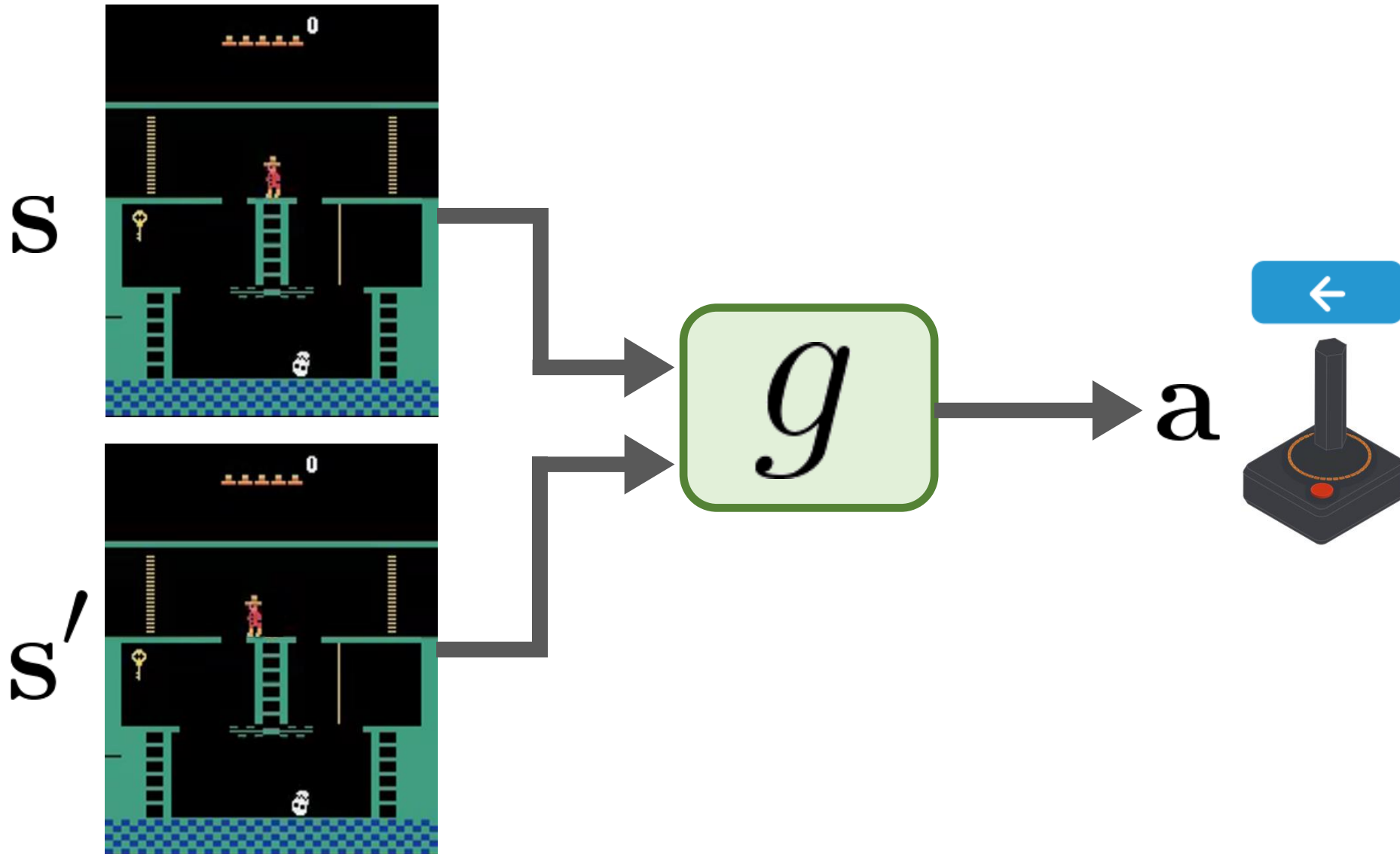
$$g(\mathbf{a} | \mathbf{s}, \mathbf{s}')$$

Inverse Dynamics Model

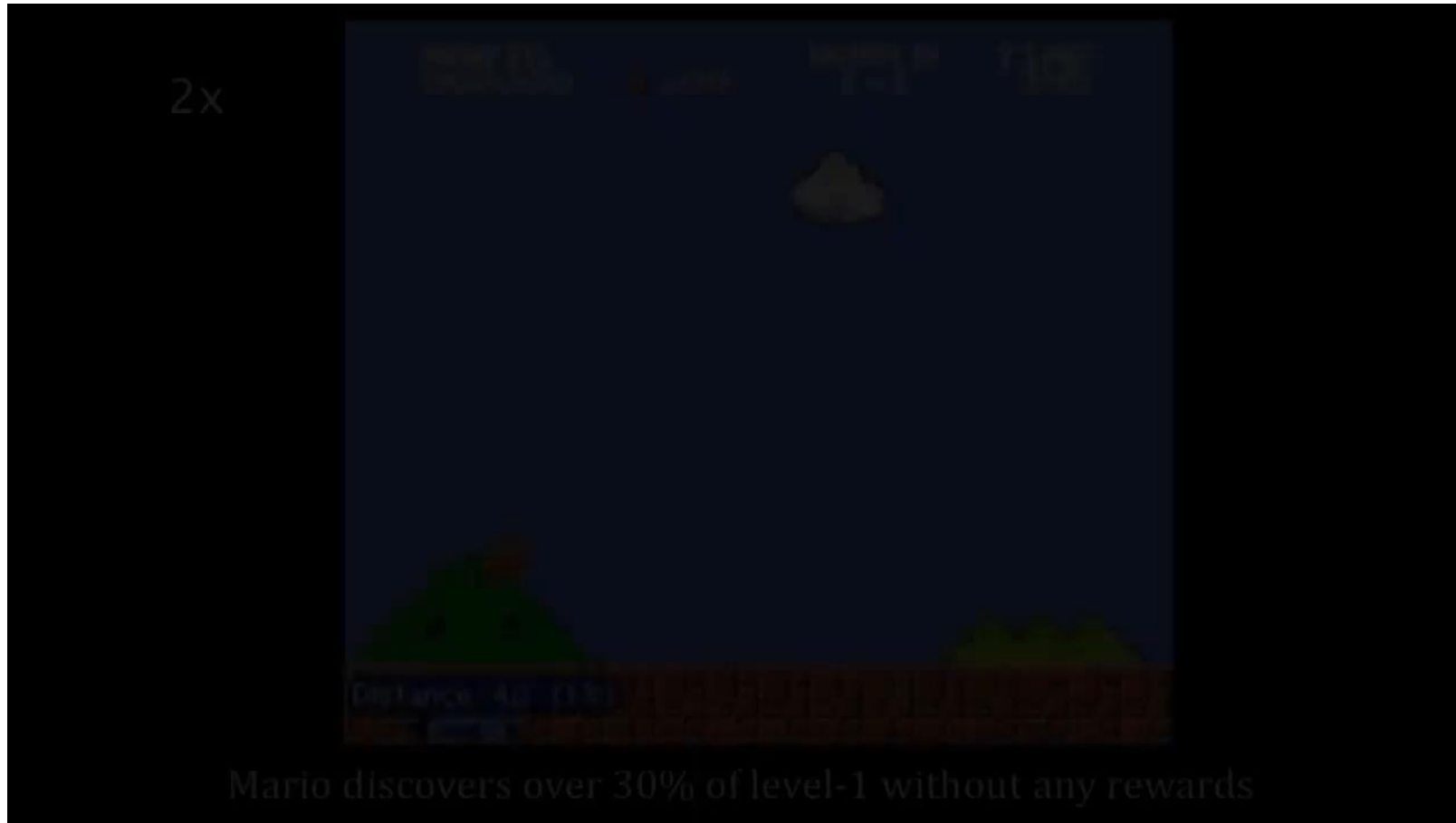


$$r_t^{\text{int}} = -\log g(\mathbf{a}_t | \mathbf{s}_t, \mathbf{s}_{t+1})$$

Inverse Dynamics Model



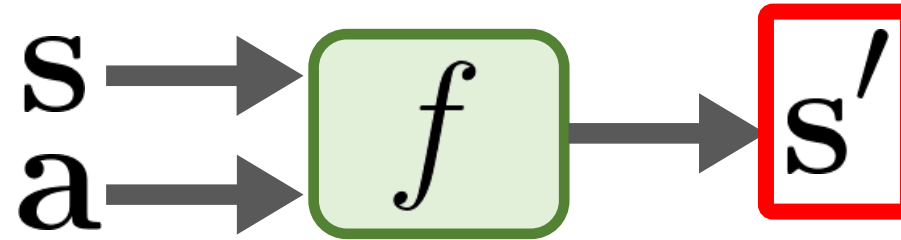
Intrinsic Curiosity Module (ICM)



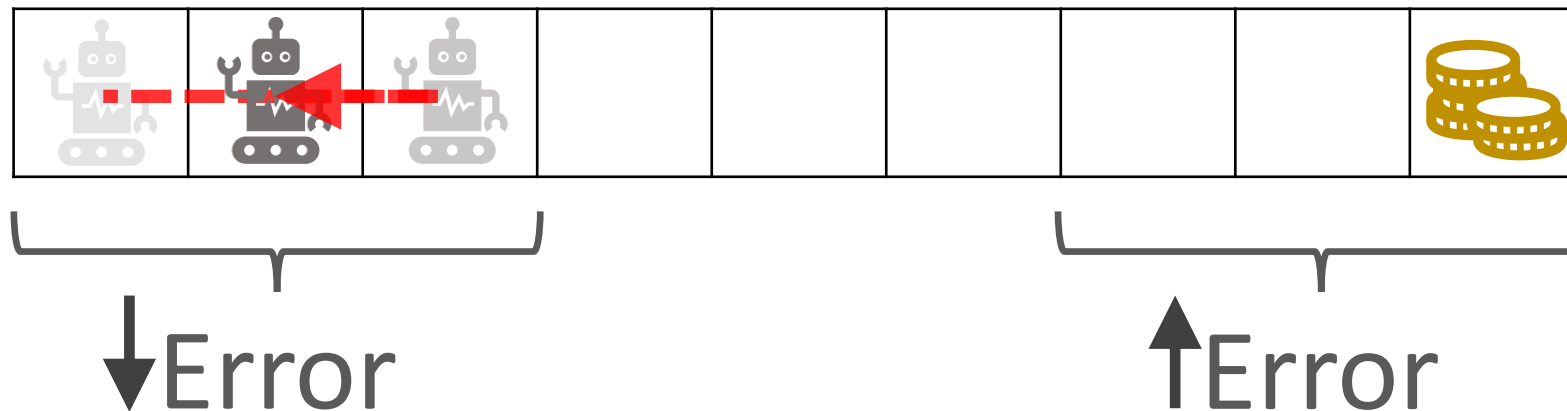
No extrinsic reward!

Curiosity-driven Exploration by Self-supervised Prediction
[Pathak et al. 2017]

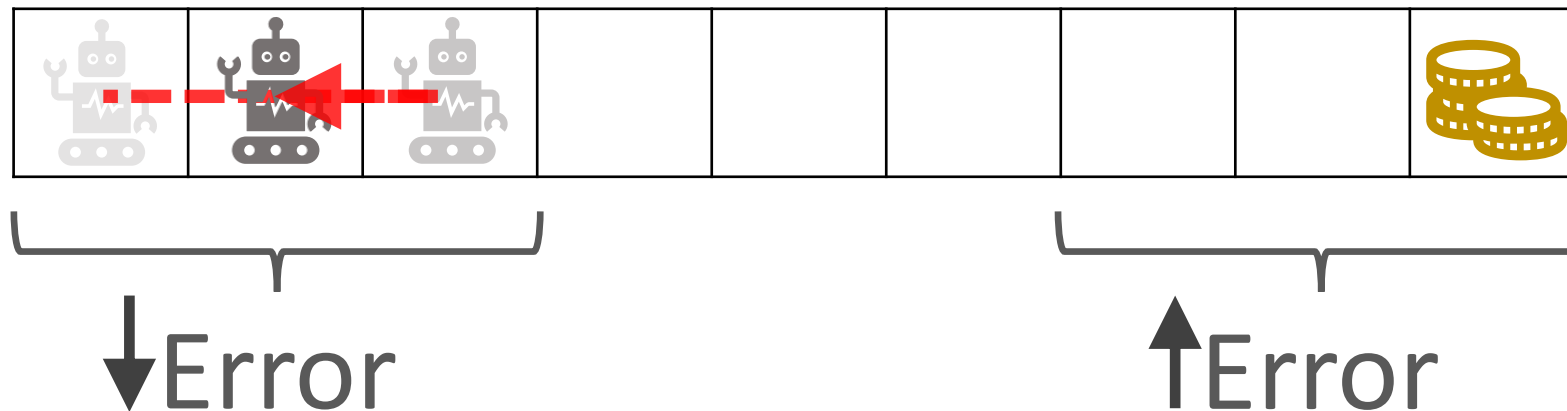
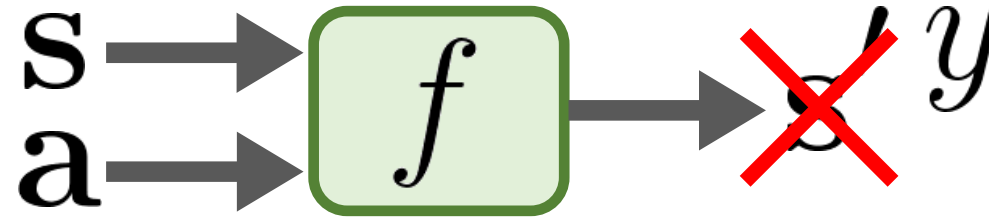
Prediction Error



Why dynamics?

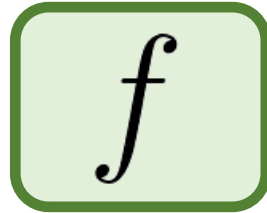


Prediction Error

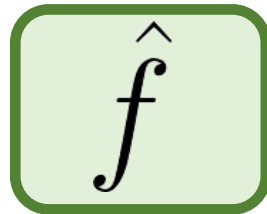


Random Network Distillation

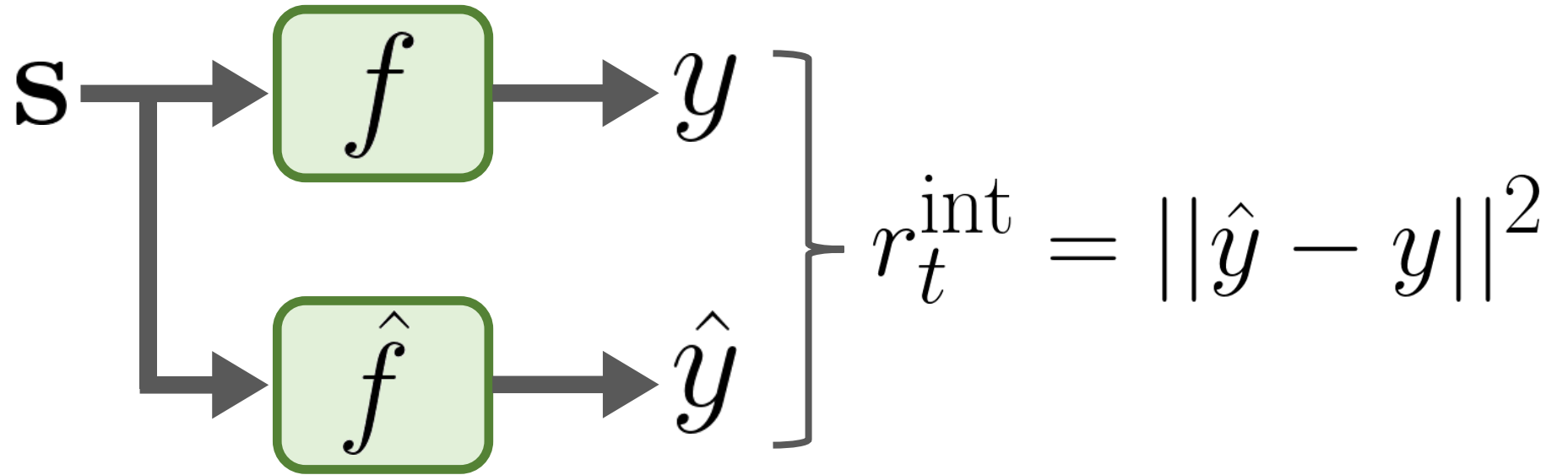
Predictor
Network



Target
Network

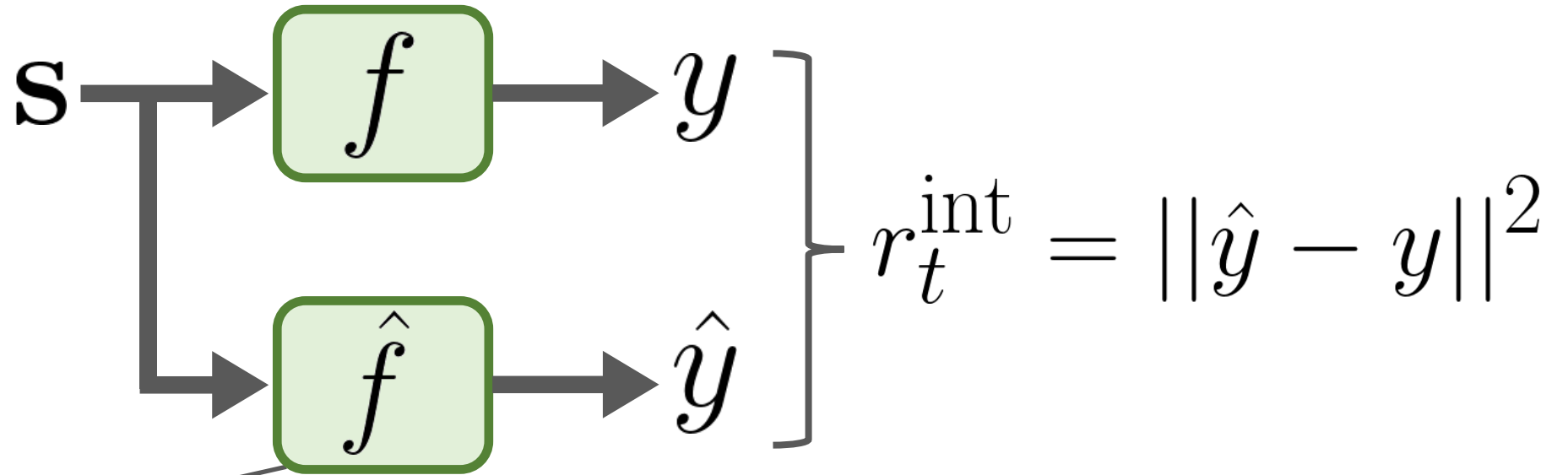


Random Network Distillation



$$\arg \max_f \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \left[\|\hat{f}(\mathbf{s}) - f(\mathbf{s})\|^2 \right]$$

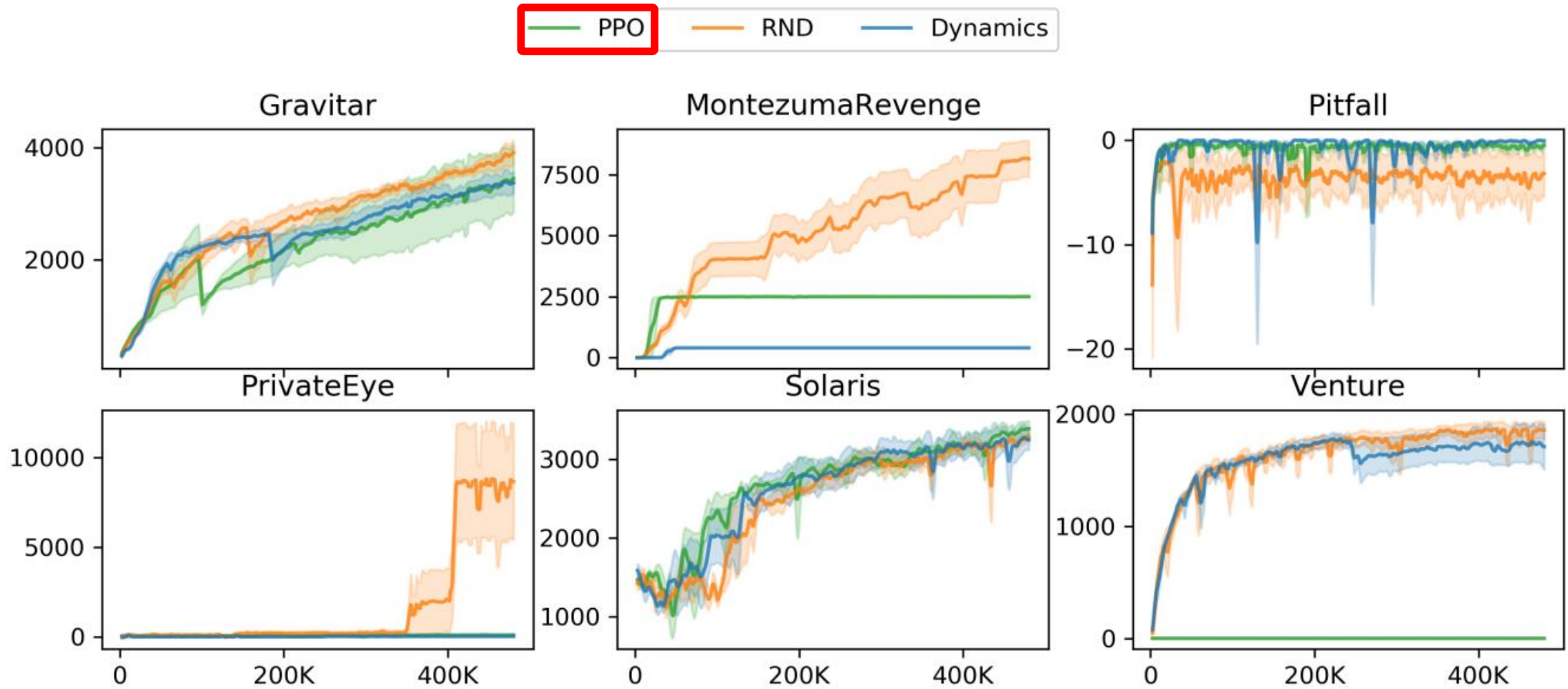
Random Network Distillation



randomly
initialized

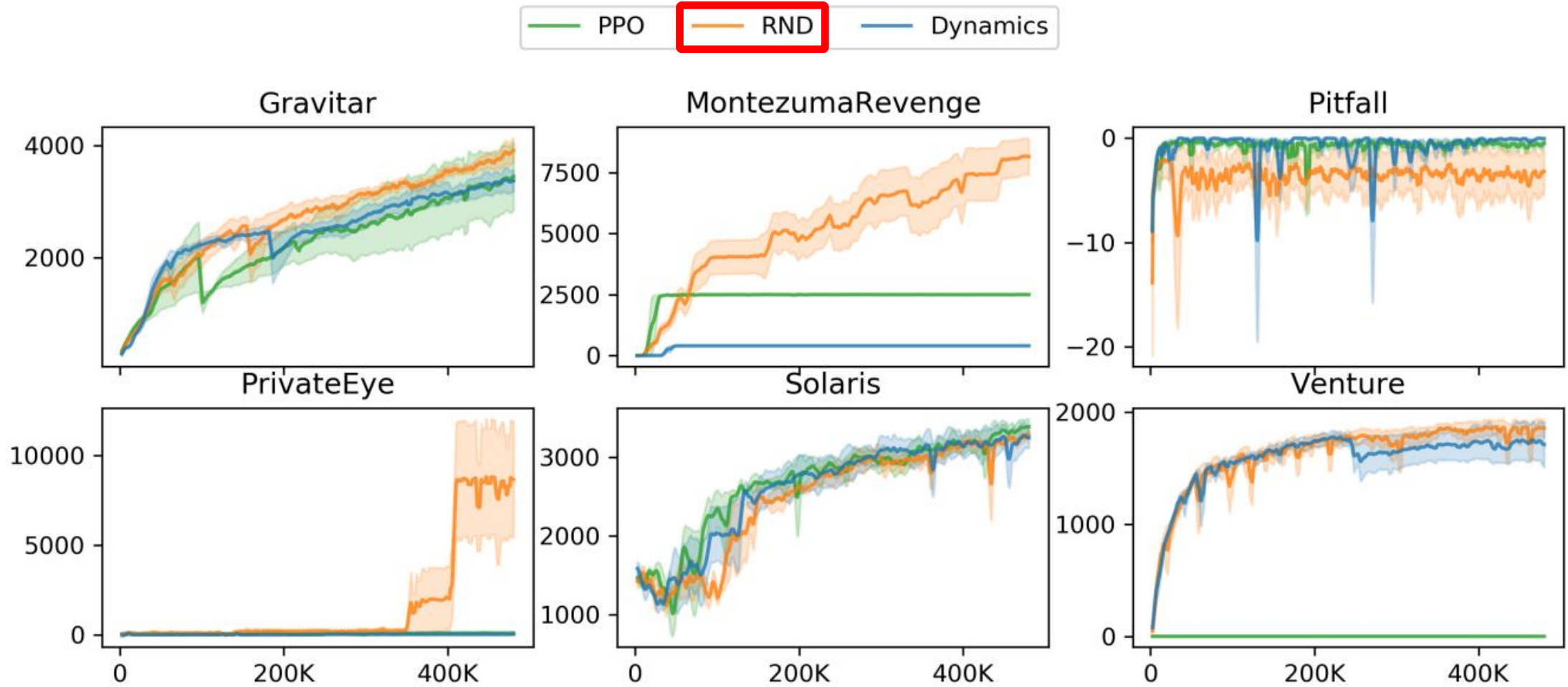
$$\arg \max_f \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} \left[\|\hat{f}(\mathbf{s}) - f(\mathbf{s})\|^2 \right]$$

Random Network Distillation



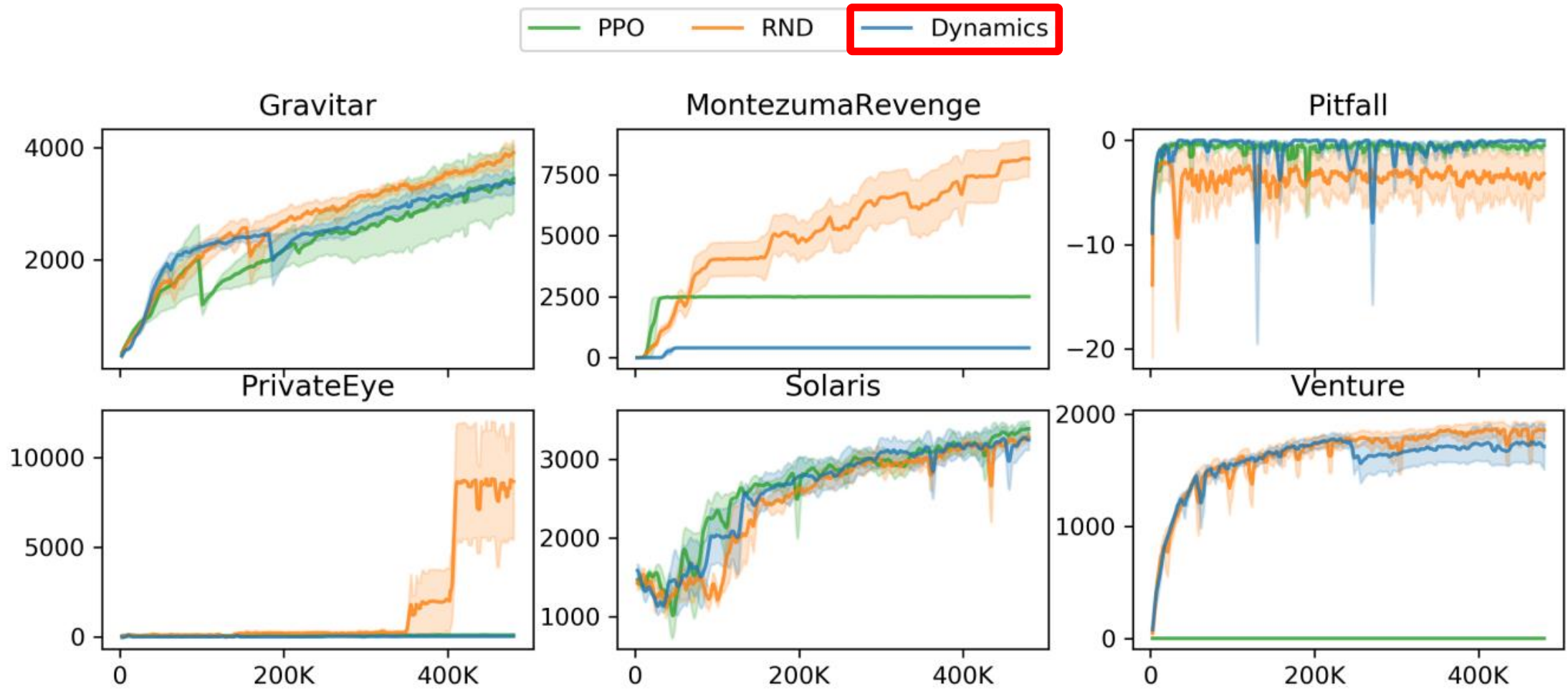
Exploration by Random Network Distillation
[Burda et al. 2019]

Random Network Distillation



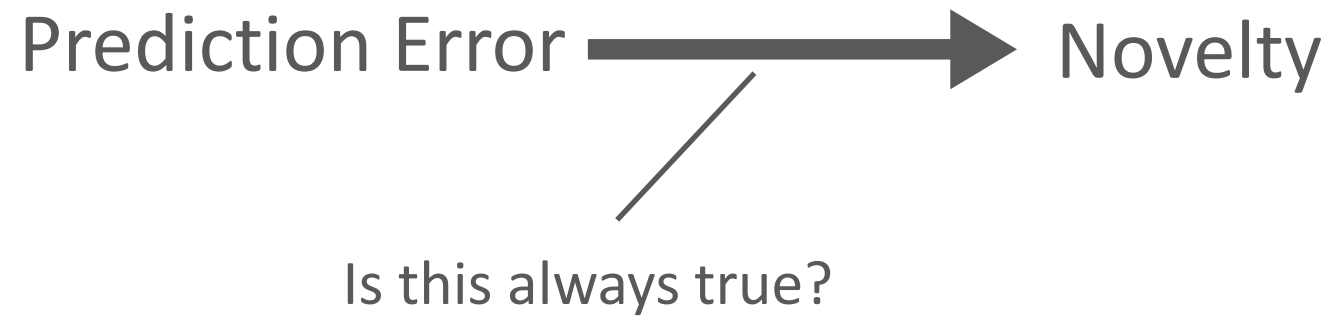
Exploration by Random Network Distillation
[Burda et al. 2019]

Random Network Distillation

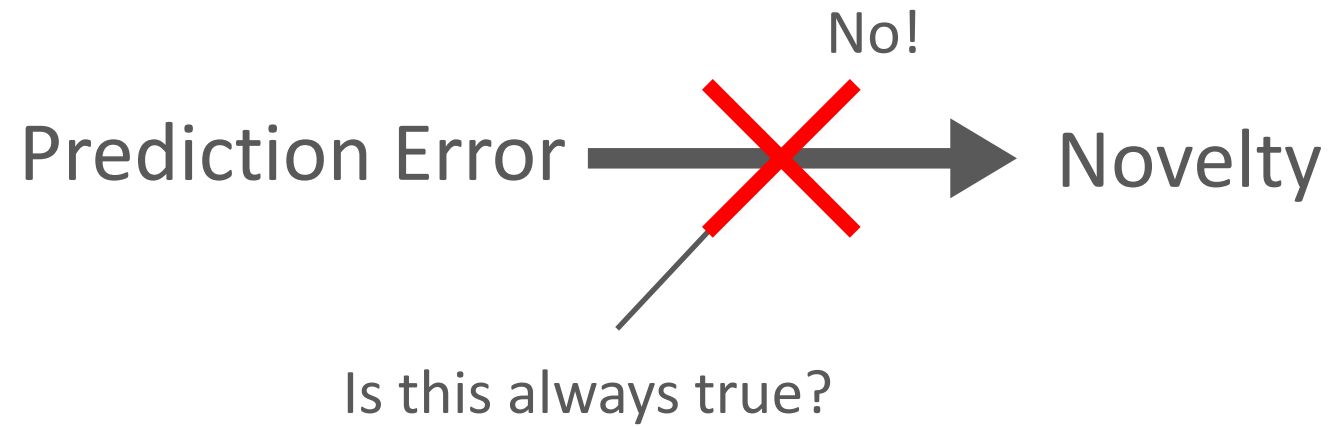


Exploration by Random Network Distillation
[Burda et al. 2019]

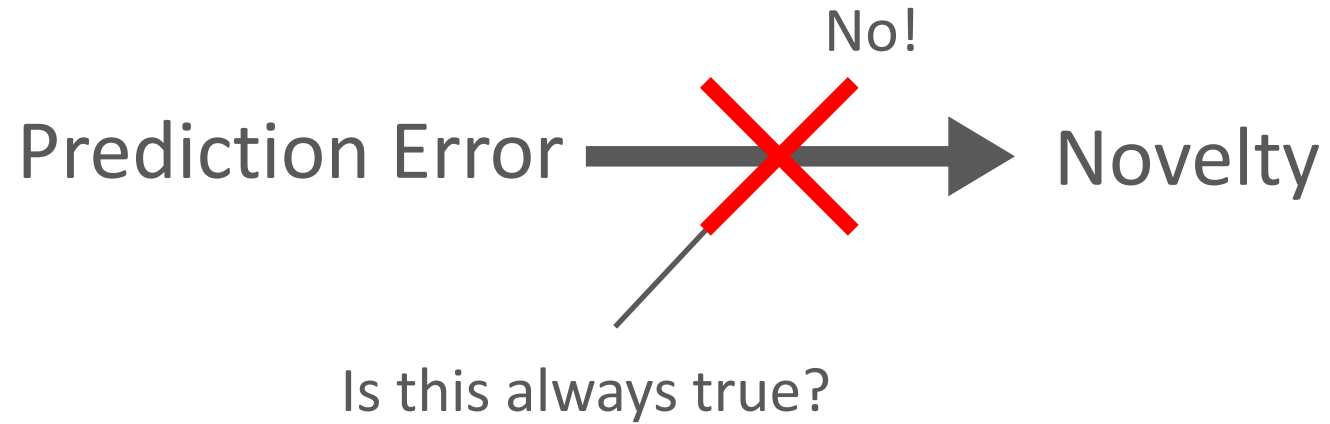
Novelty



Novelty



Novelty



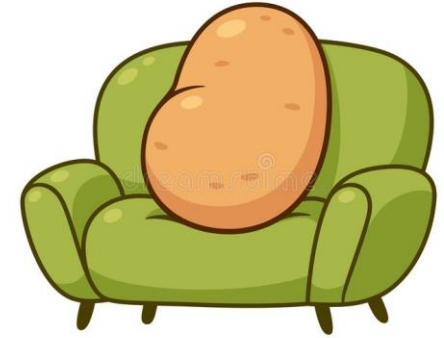
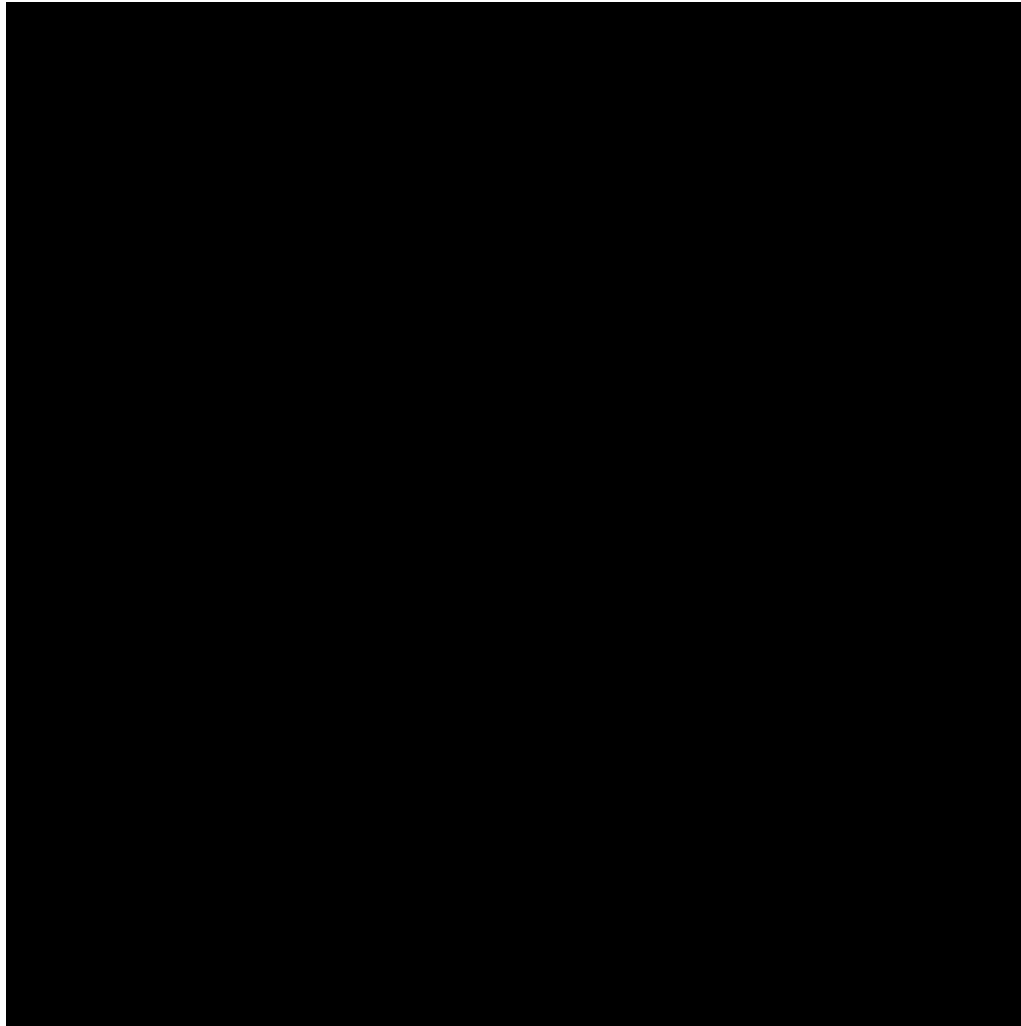
Unpredictable \neq Novel

Noisy-TV Problem



Large-Scale Study of Curiosity-Driven Learning
[Burda et al. 2018]

Noisy-TV Problem

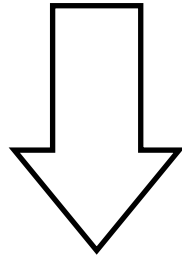


- Prediction error works well in static envs
- But can get distracted by variability in dynamic envs

Surprise Minimization

Maximize Surprise

$$r_t^{\text{int}} = -\log(\rho(\mathbf{s}_{t+1}))$$



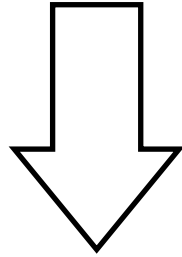
Minimize Surprise

$$r_t^{\text{int}} = \log(\rho(\mathbf{s}_{t+1}))$$

Surprise Minimization

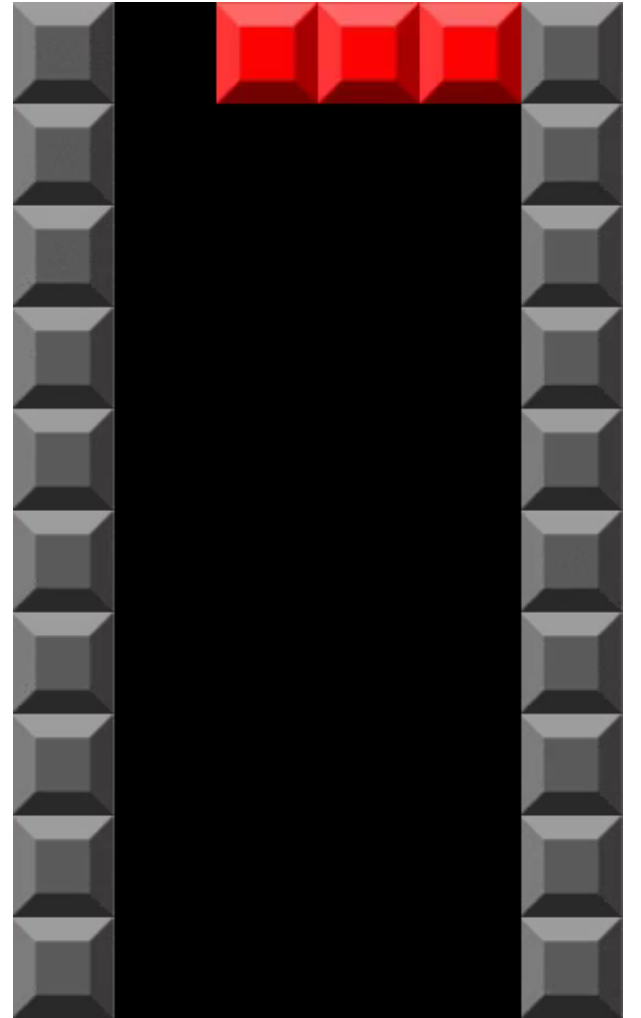
Maximize Surprise

$$r_t^{\text{int}} = -\log(\rho(\mathbf{s}_{t+1}))$$



Minimize Surprise

$$r_t^{\text{int}} = \log(\rho(\mathbf{s}_{t+1}))$$



Summary

- Exploration Exploitation Tradeoff
- Dense vs Sparse Rewards
- Intrinsic Motivation
- Count-Based Exploration
- Surprise Maximization