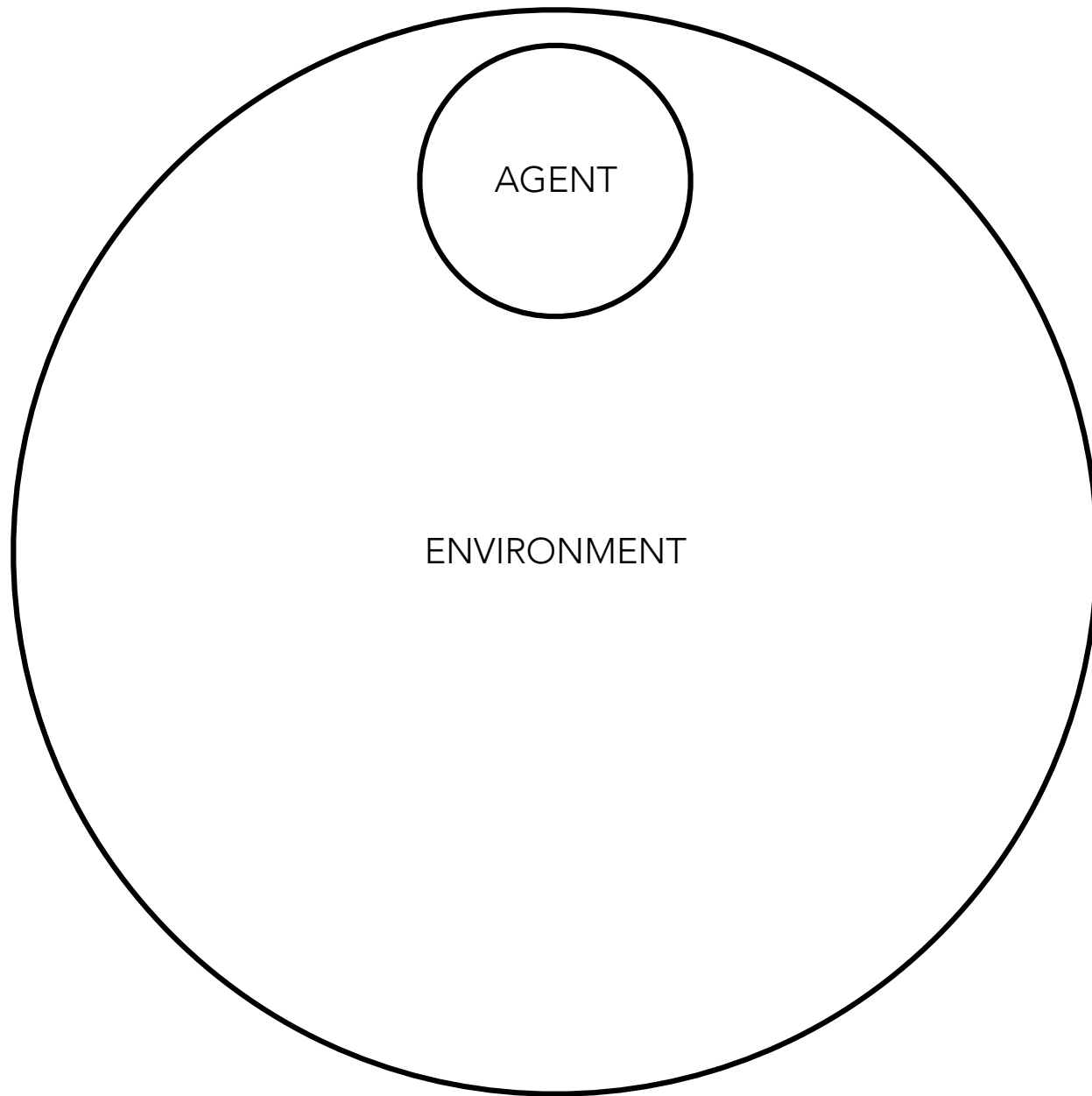
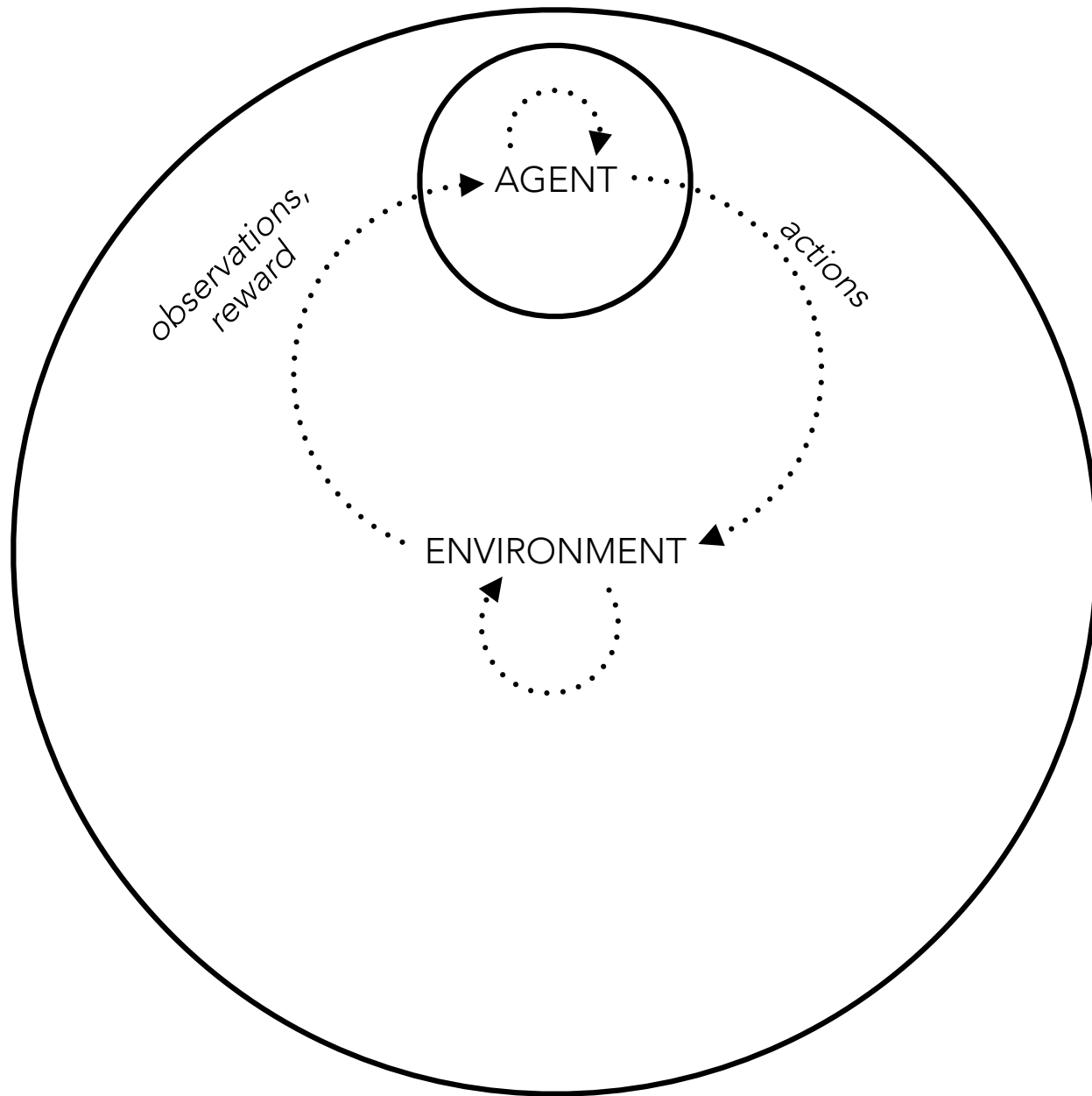
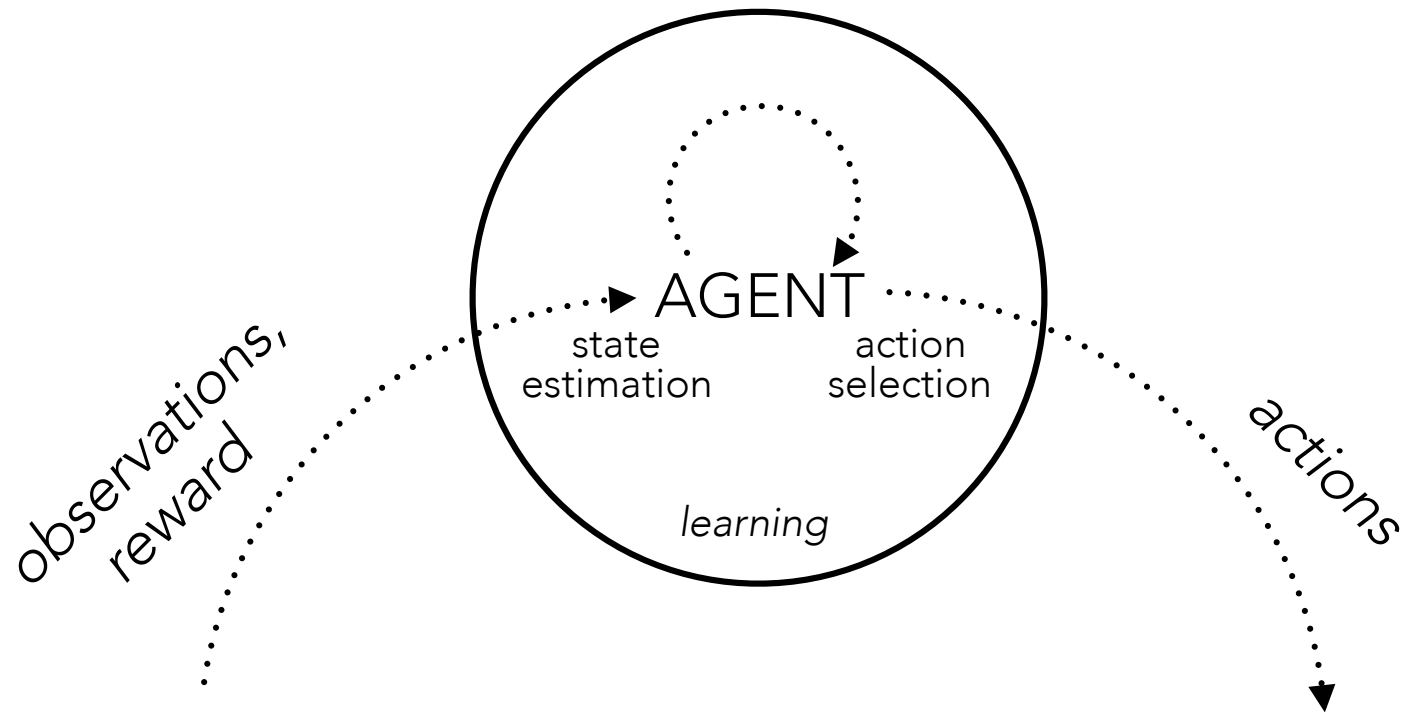
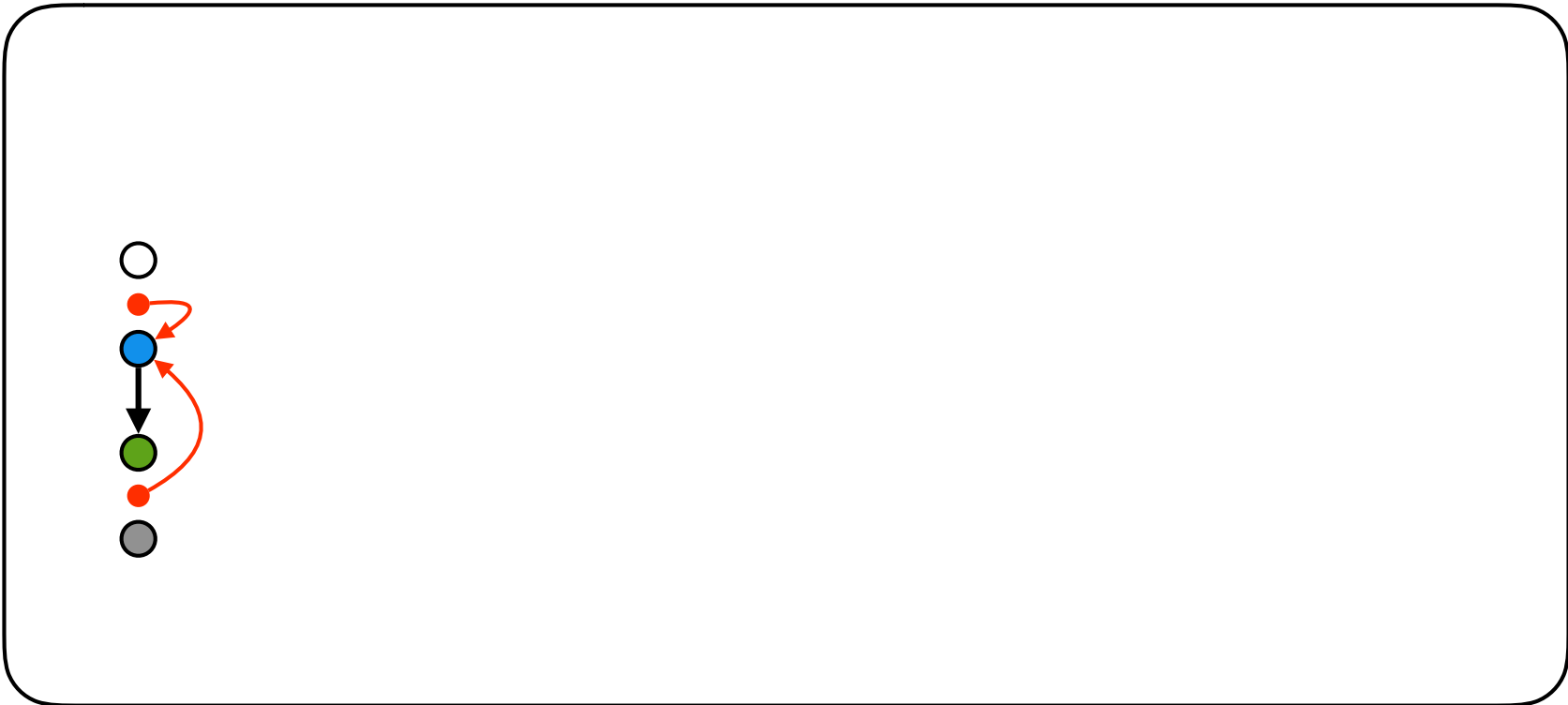

An Inference Perspective on Model-Based Reinforcement Learning

Joseph Marino, Yisong Yue
Caltech

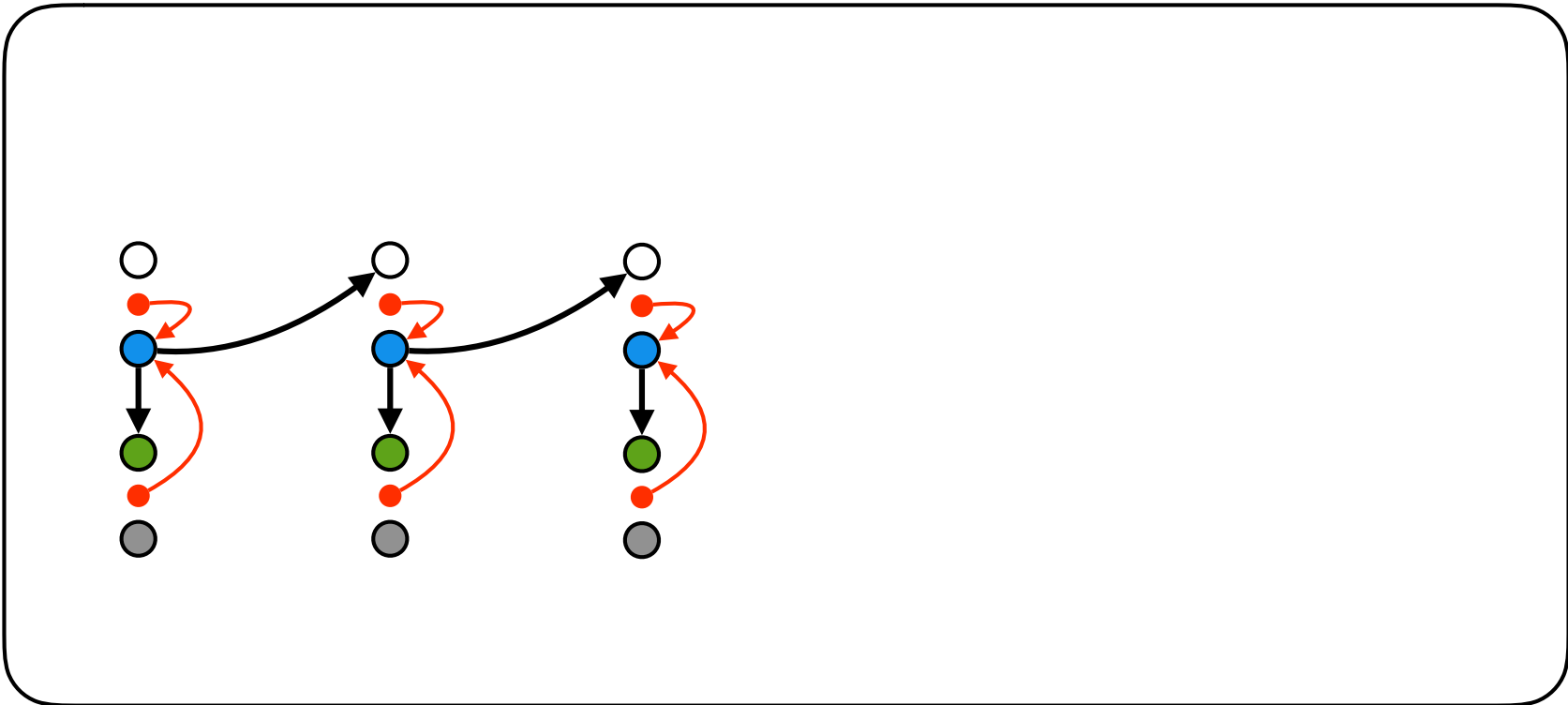






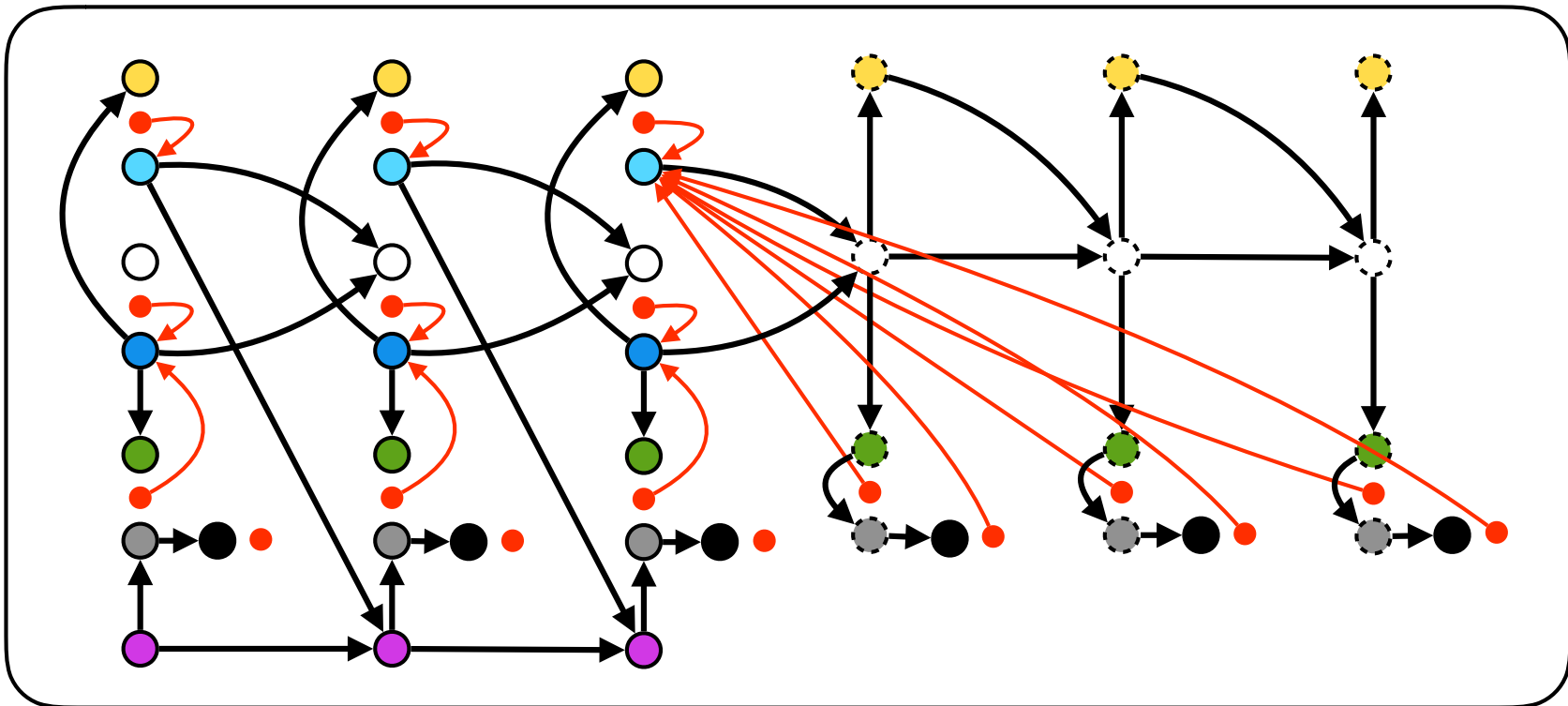


**Iterative
Amortized Inference**
Marino, et al., 2018a



**Iterative
Amortized Inference**
Marino, et al., 2018a

**Amortized
Variational Filtering**
Marino, et al., 2018b



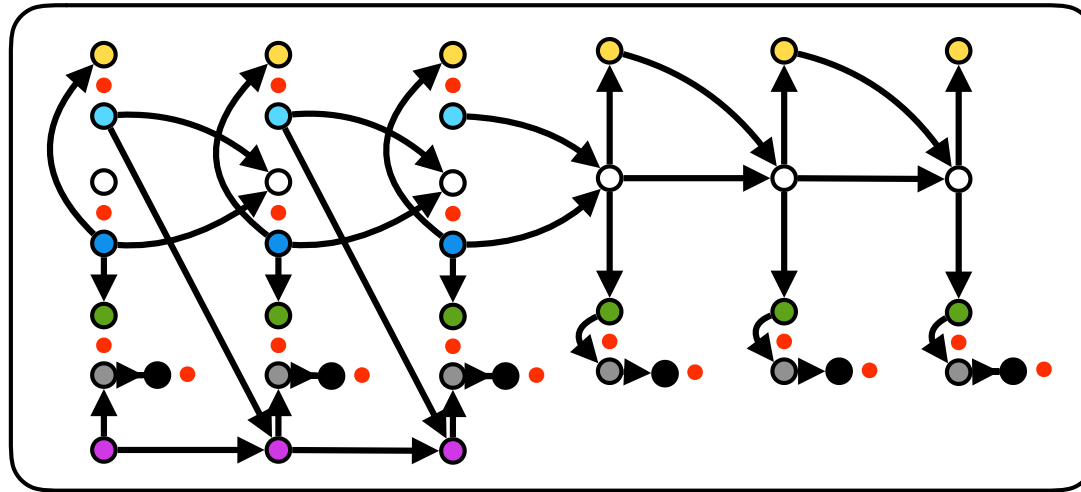
**Iterative
Amortized Inference**
Marino, et al., 2018a

**Amortized
Variational Filtering**
Marino, et al., 2018b

**An Inference Perspective
on Model-Based RL**
Marino & Yue, 2019

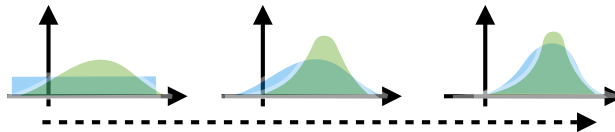
OVERVIEW

frame model-based RL as probabilistic inference & learning



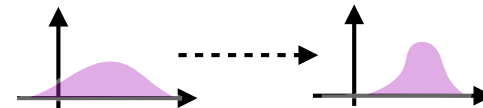
two additional terms:

prior over actions



combine model-based likelihood with a model-free prior

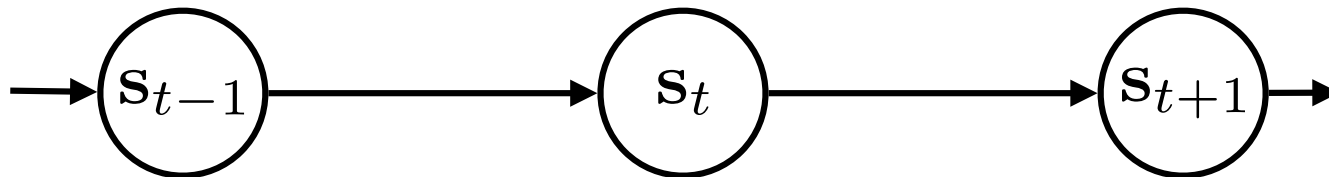
information gain from observations



model task-relevant state information, biases planning

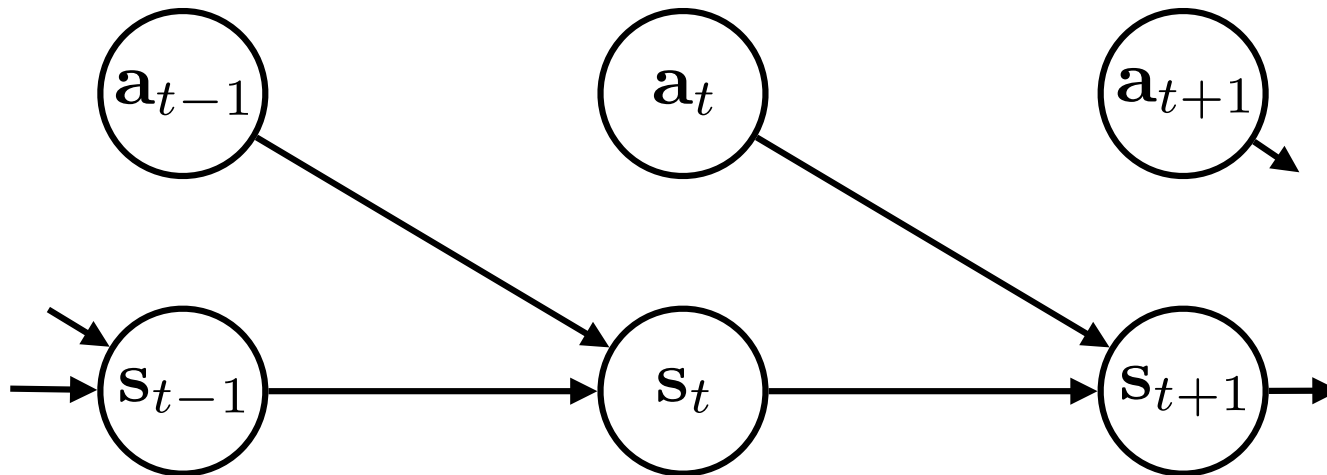
REINFORCEMENT LEARNING AS INFERENCE

reformulate RL as a probabilistic inference problem



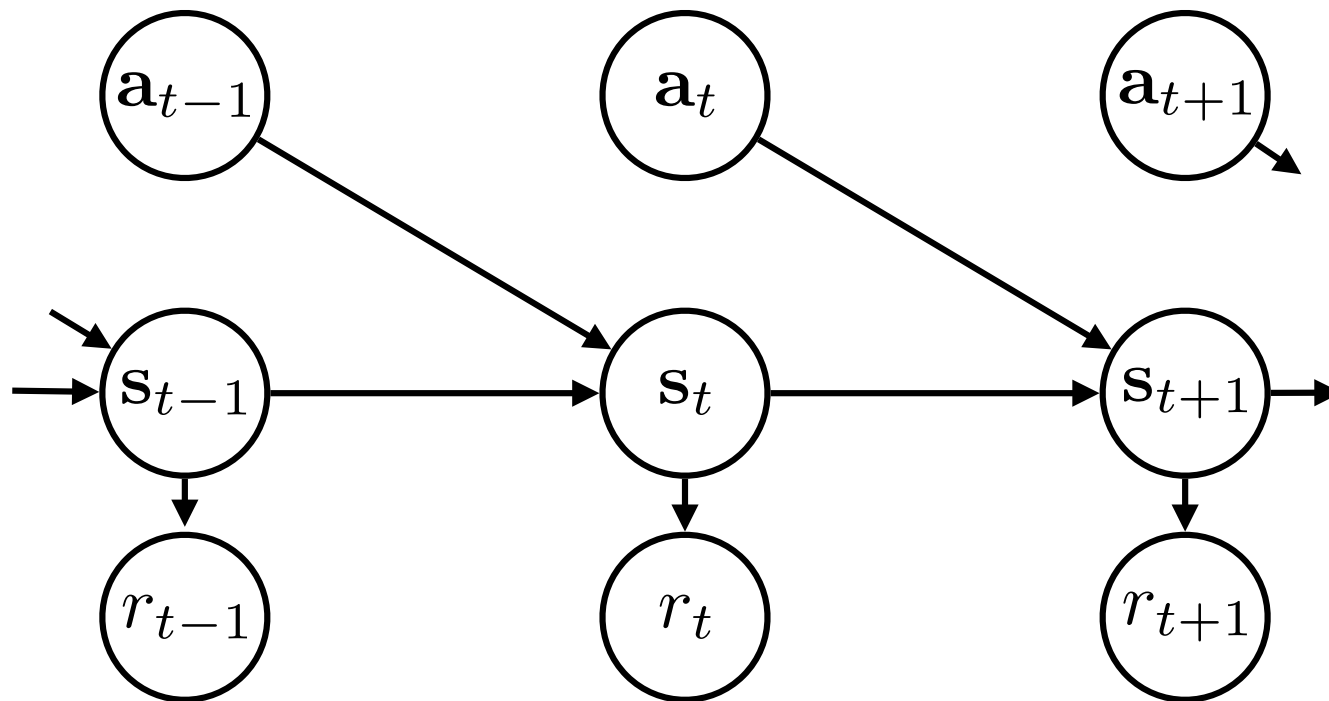
REINFORCEMENT LEARNING AS INFERENCE

reformulate RL as a probabilistic inference problem



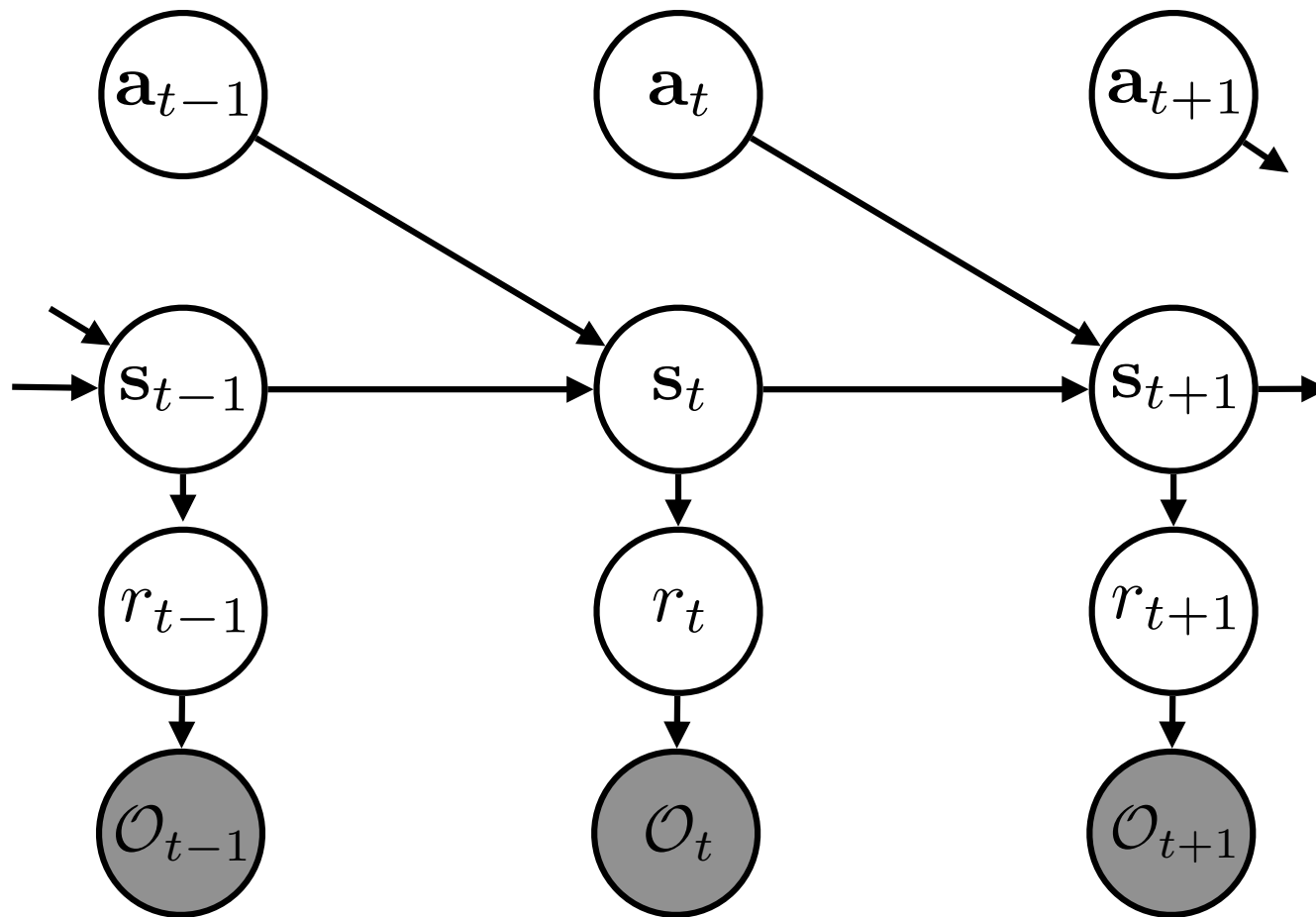
REINFORCEMENT LEARNING AS INFERENCE

reformulate RL as a probabilistic inference problem



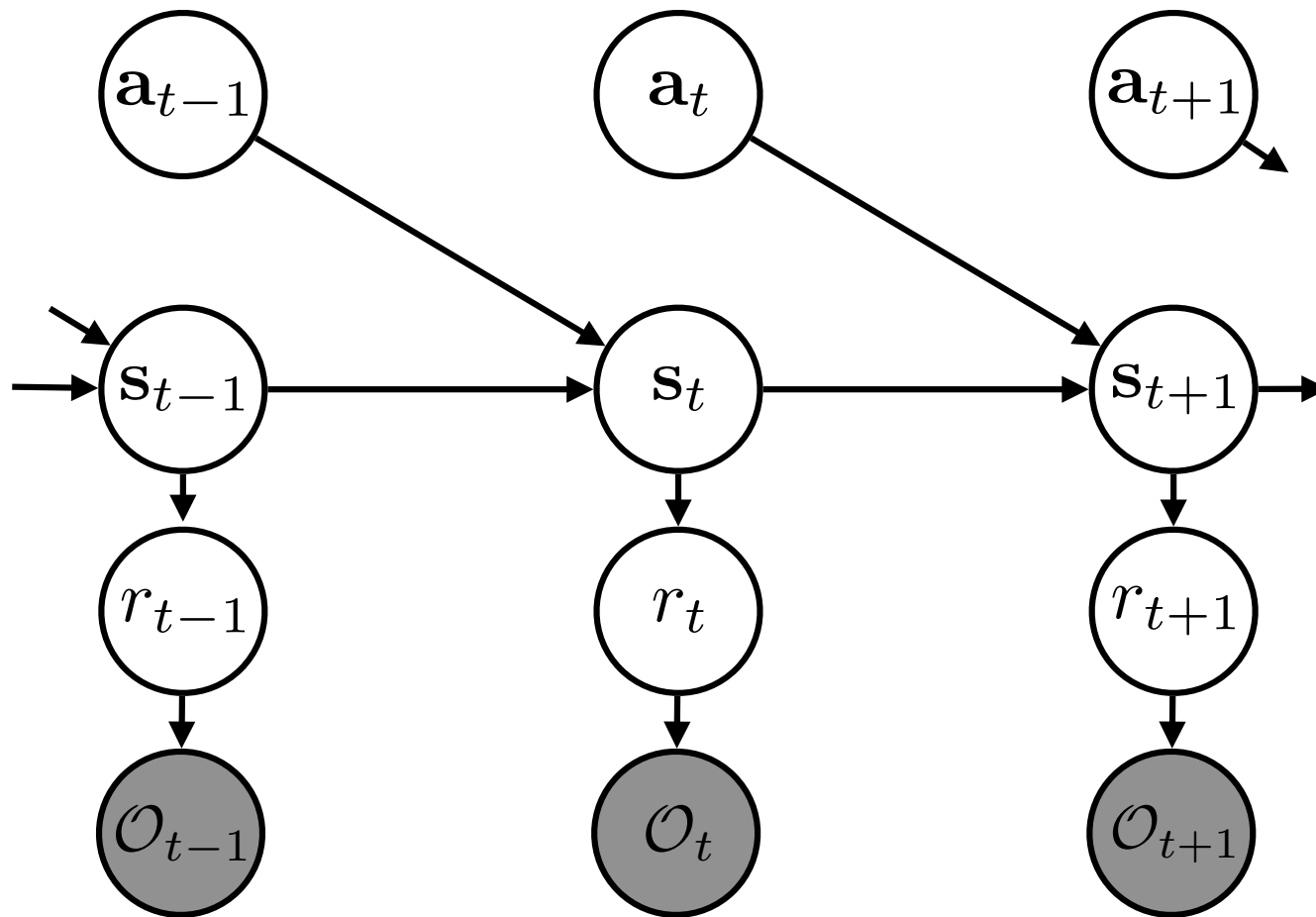
REINFORCEMENT LEARNING AS INFERENCE

reformulate RL as a probabilistic inference problem



PLANNING AS INFERENCE

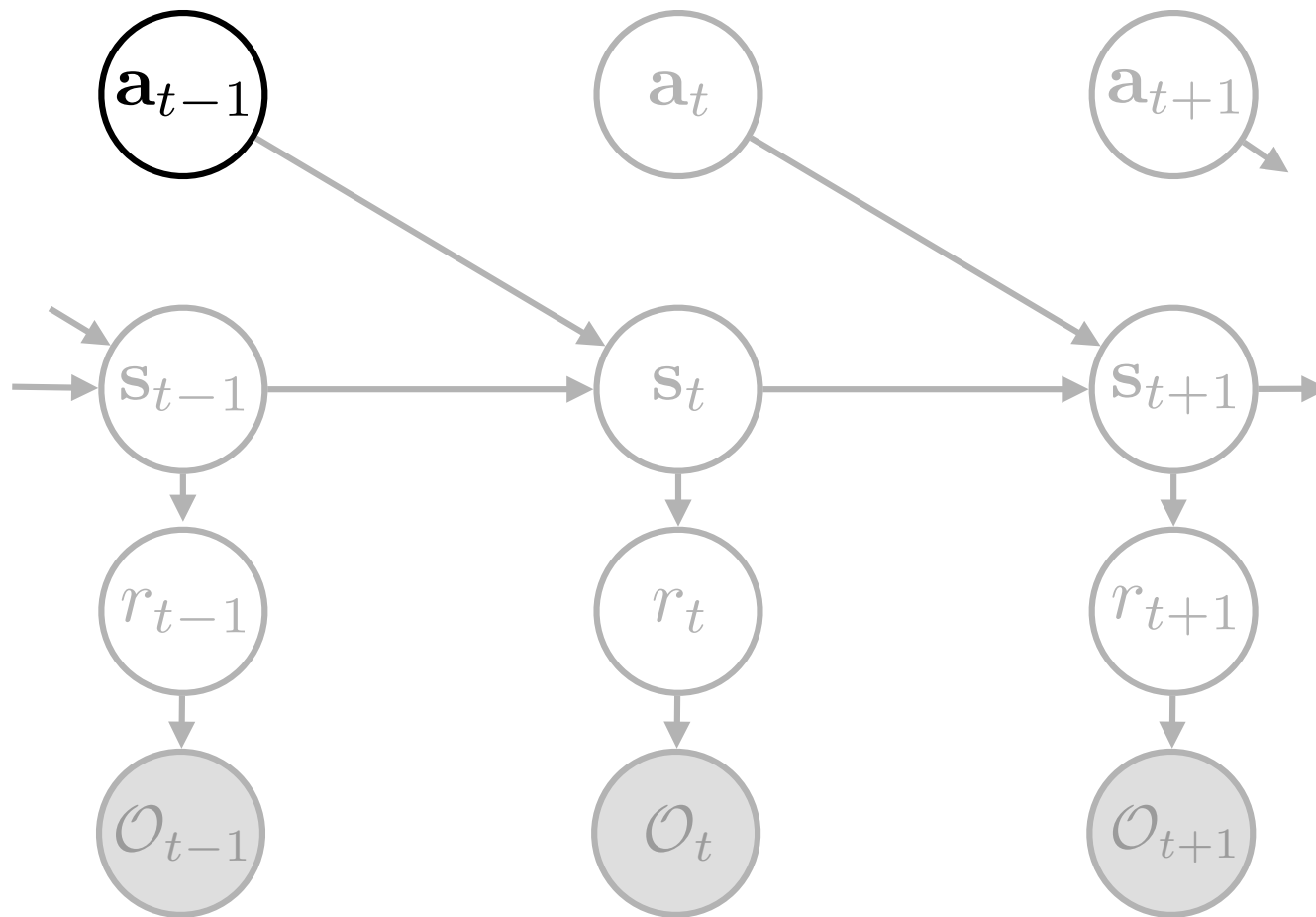
can also reformulate planning as inference



Attias, 2003
Botvinick & Toussaint, 2012

PLANNING AS INFERENCE

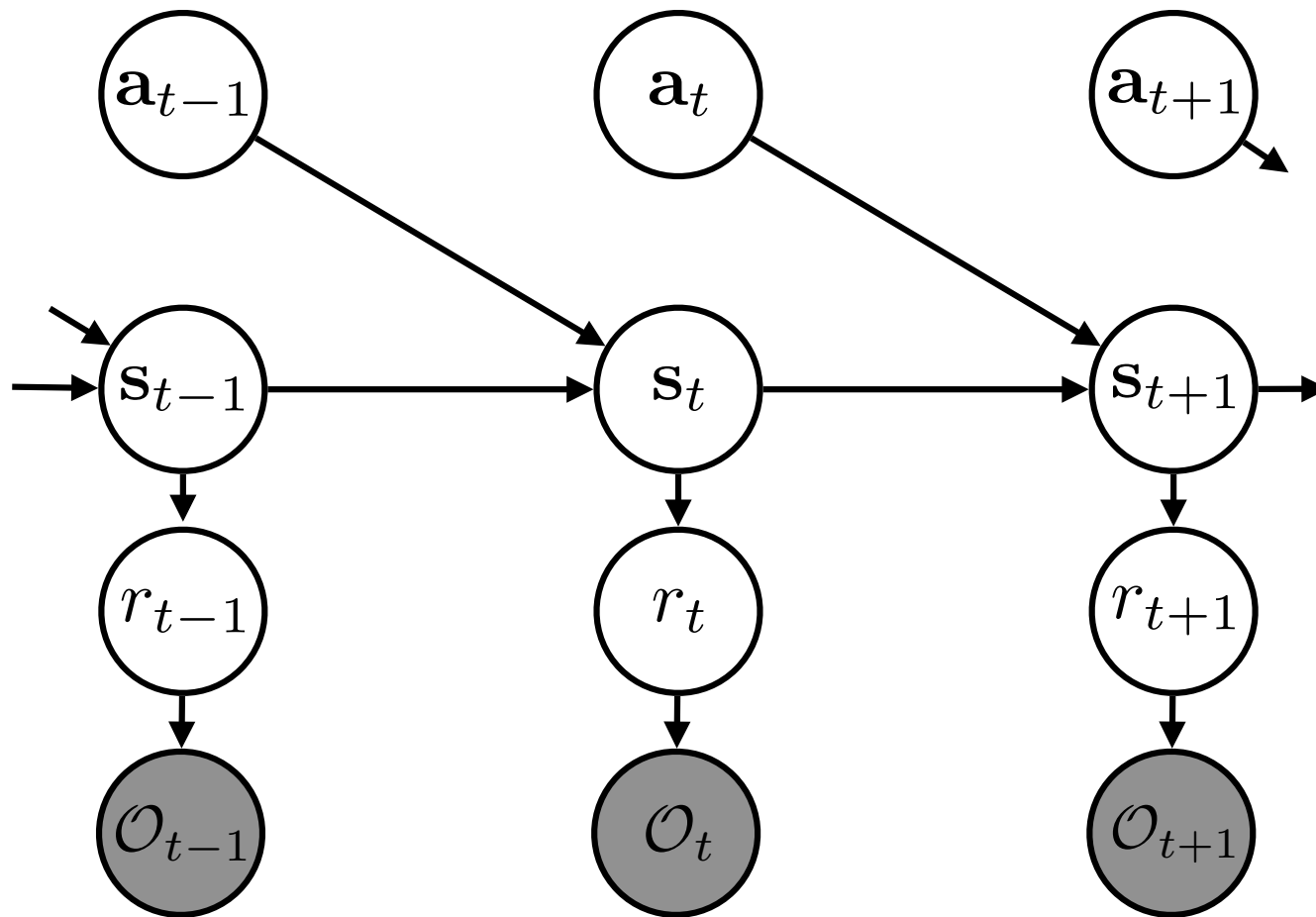
can also reformulate planning as inference



Attias, 2003
Botvinick & Toussaint, 2012

PLANNING AS INFERENCE

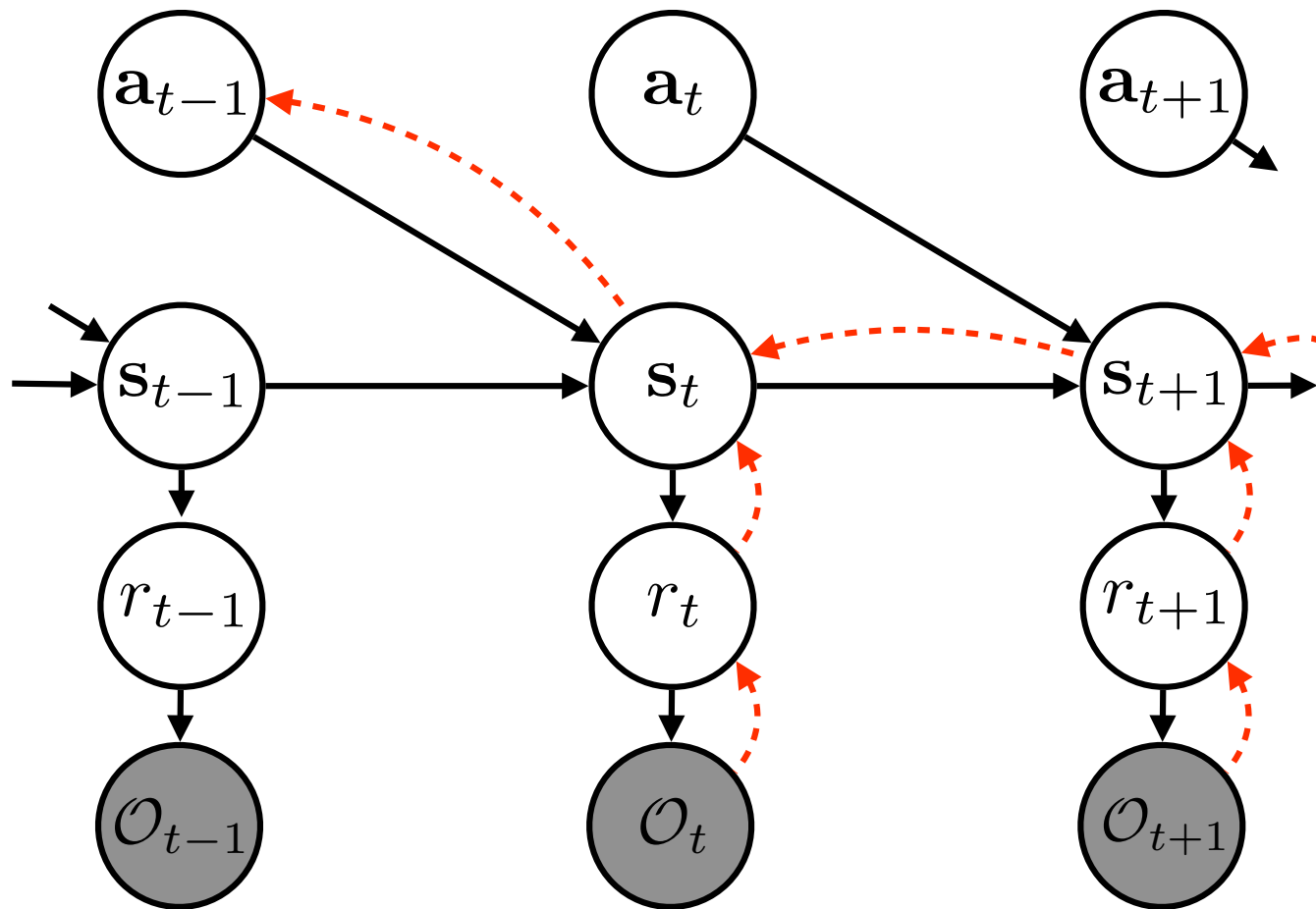
can also reformulate planning as inference



Attias, 2003
Botvinick & Toussaint, 2012

PLANNING AS INFERENCE

can also reformulate planning as inference

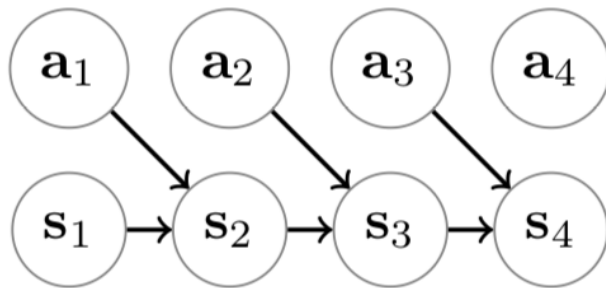


Attias, 2003
Botvinick & Toussaint, 2012

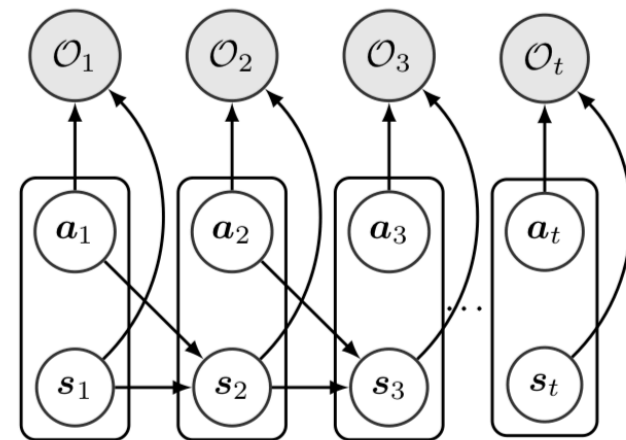
PREVIOUS WORKS

model priors are not learned

focused on fully-observable environments

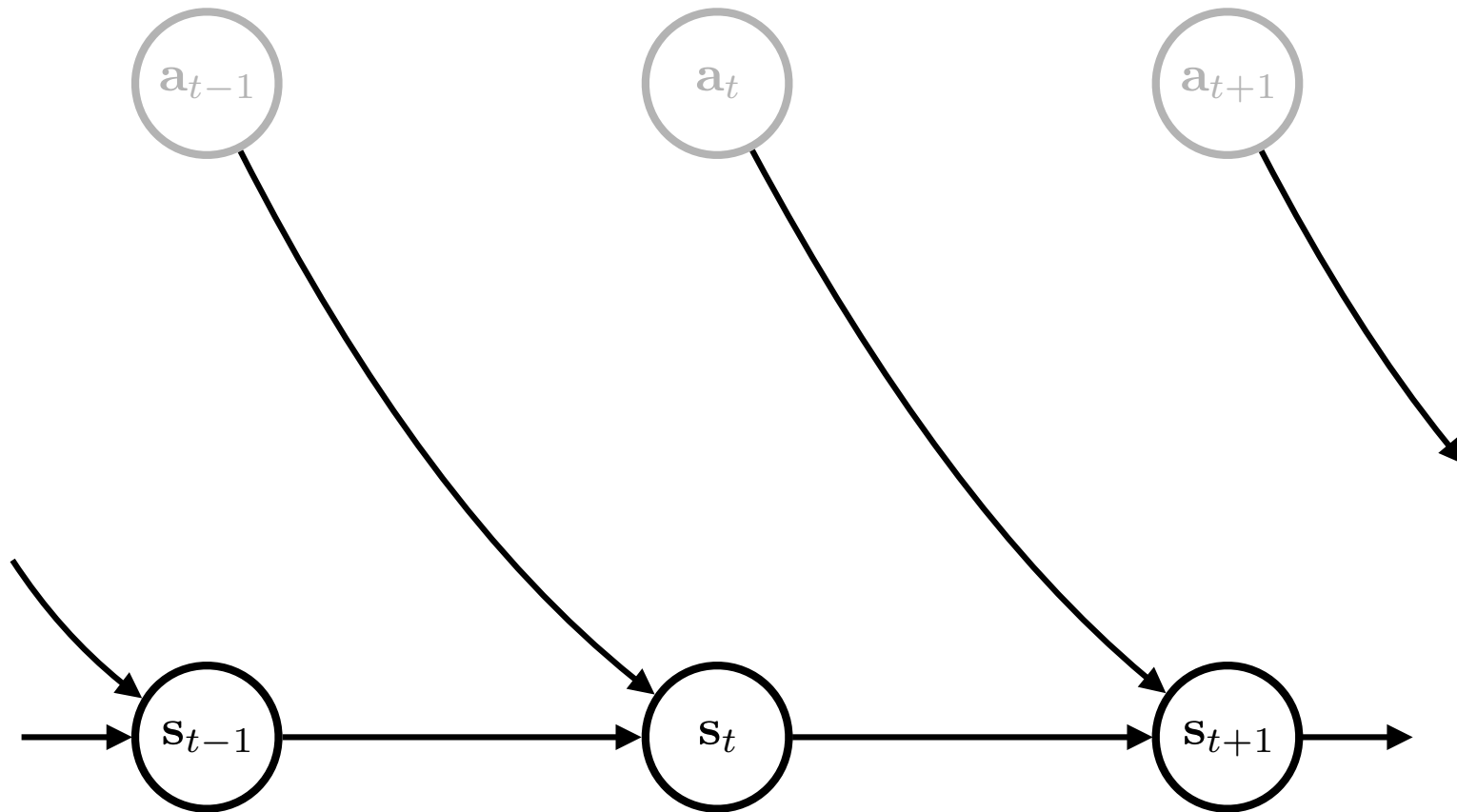


Levine, 2018



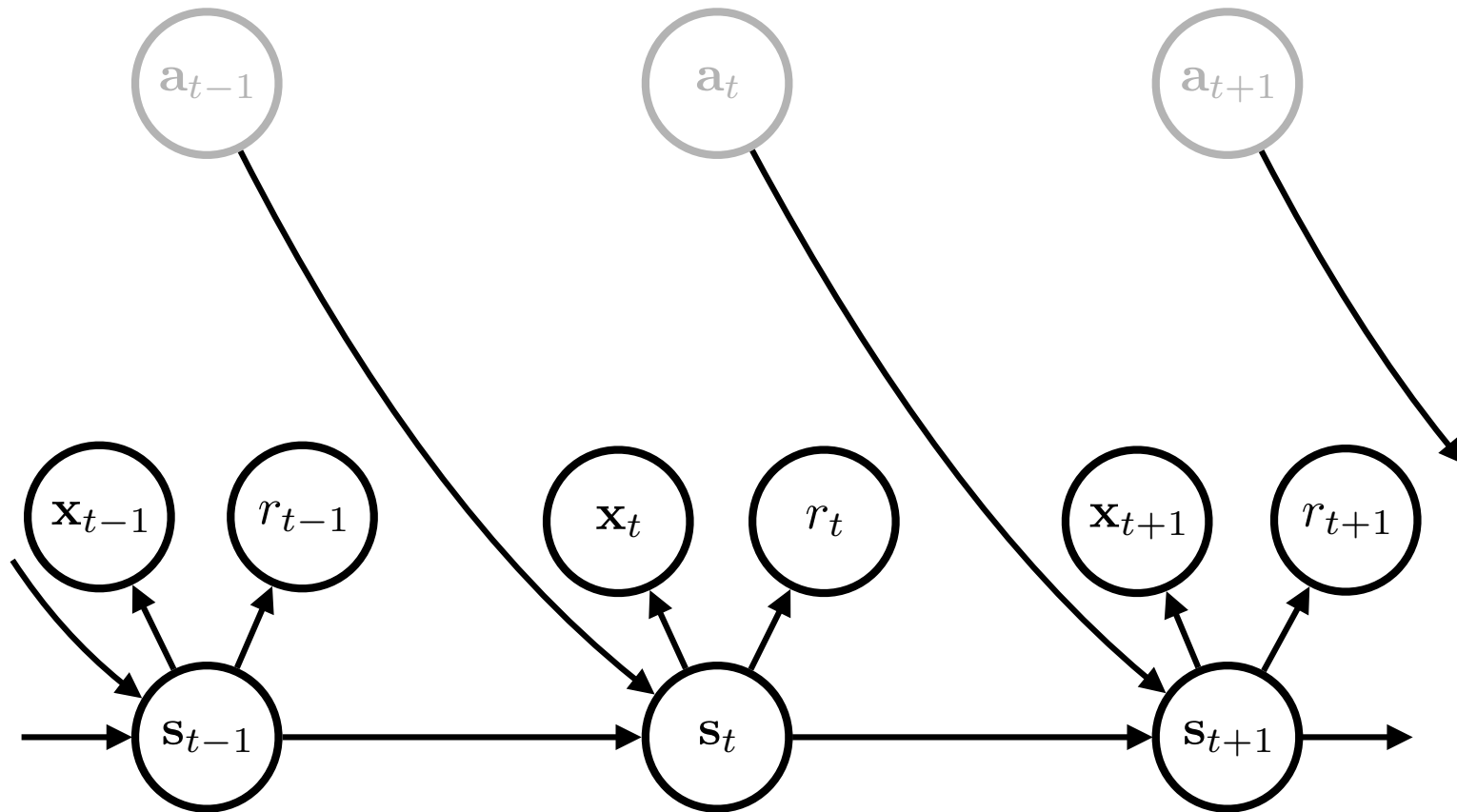
Piché, Thomas, et al., 2019

ENVIRONMENT



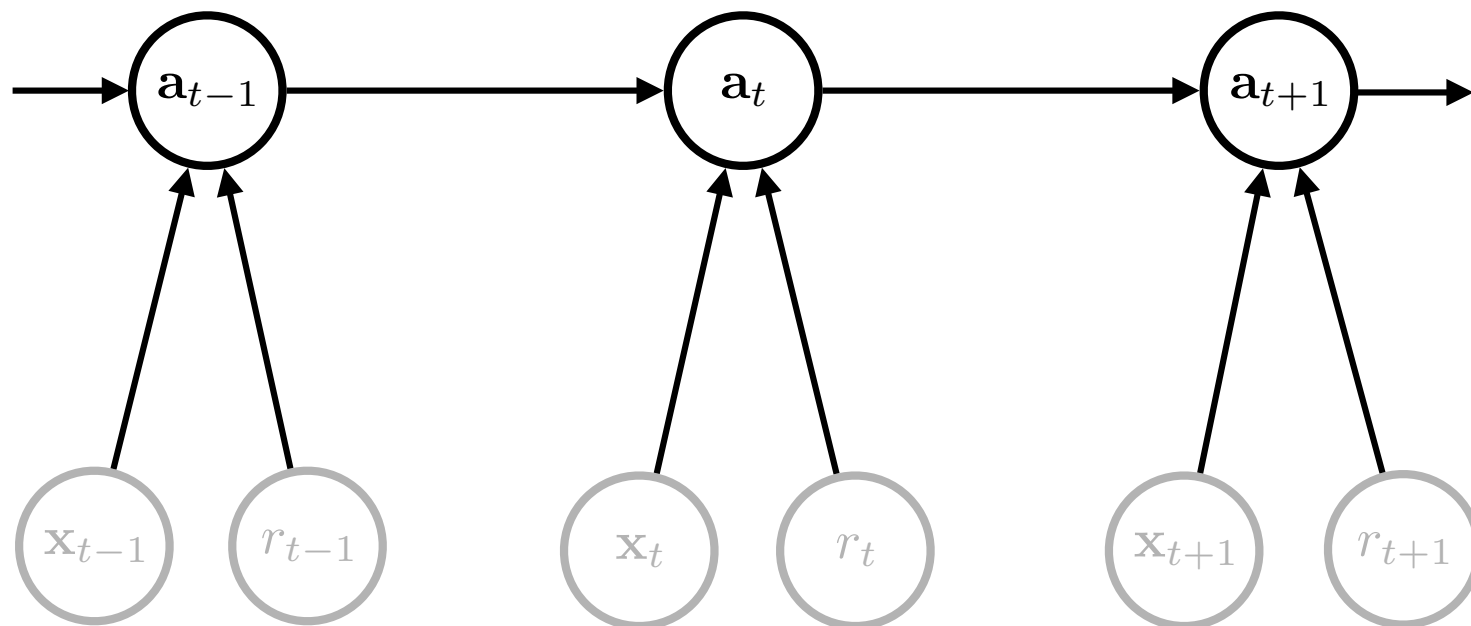
$$p_e(\mathbf{x}_{1:T}, r_{1:T}, \mathbf{s}_{1:T} | \mathbf{a}_{1:T-1}) = \prod_{t=1}^T \underbrace{p_e(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{a}_{t-1})}_{\text{state dynamics}}$$

ENVIRONMENT



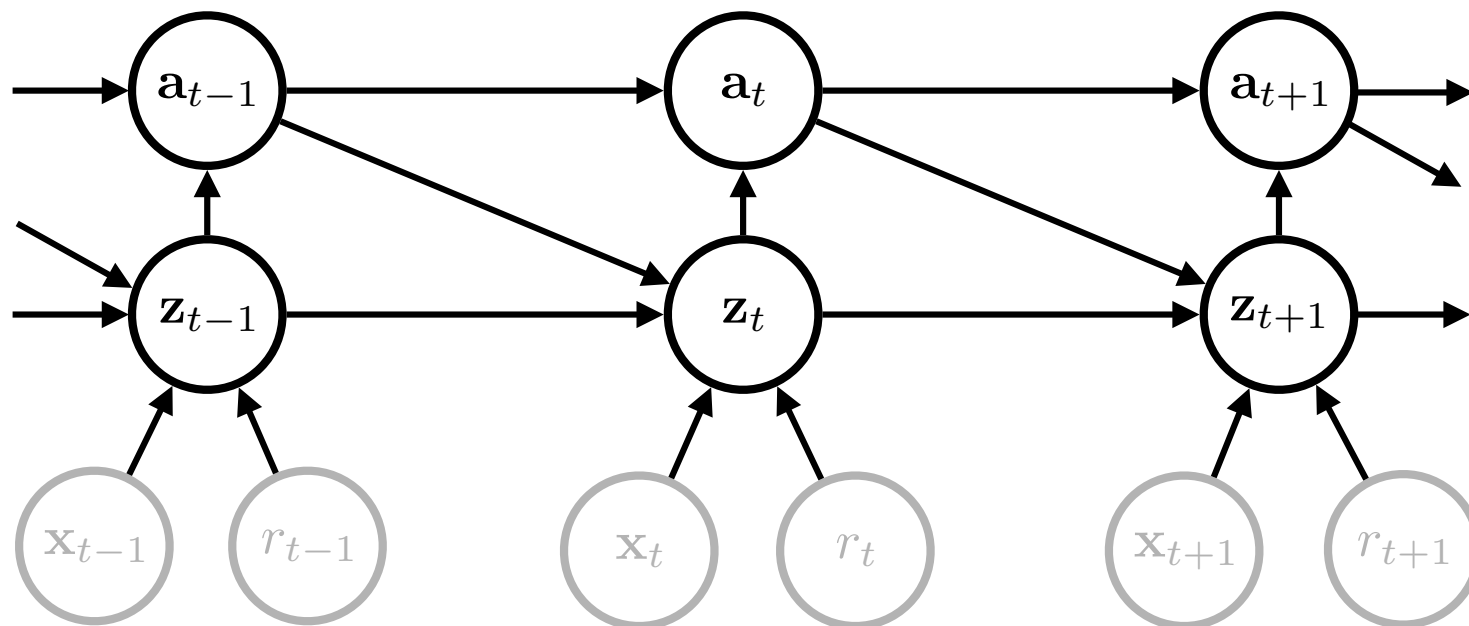
$$p_e(\mathbf{x}_{1:T}, r_{1:T}, \mathbf{s}_{1:T} | \mathbf{a}_{1:T-1}) = \prod_{t=1}^T \underbrace{p_e(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{a}_{t-1})}_{\text{state dynamics}} \underbrace{p_e(\mathbf{x}_t | \mathbf{s}_t) p_e(r_t | \mathbf{s}_t)}_{\text{observations/rewards}}$$

AGENT



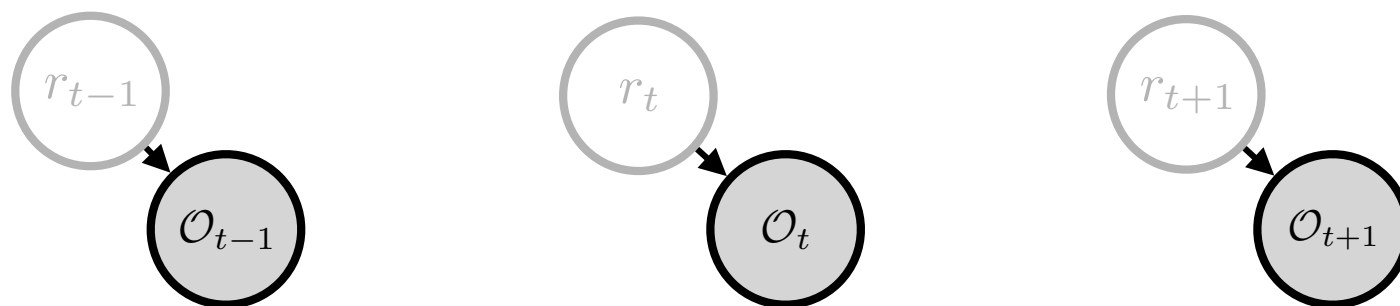
$$p_{\mathbf{a}}(\mathbf{a}_{1:T} | \mathbf{x}_{1:T}, r_{1:t})$$

AGENT



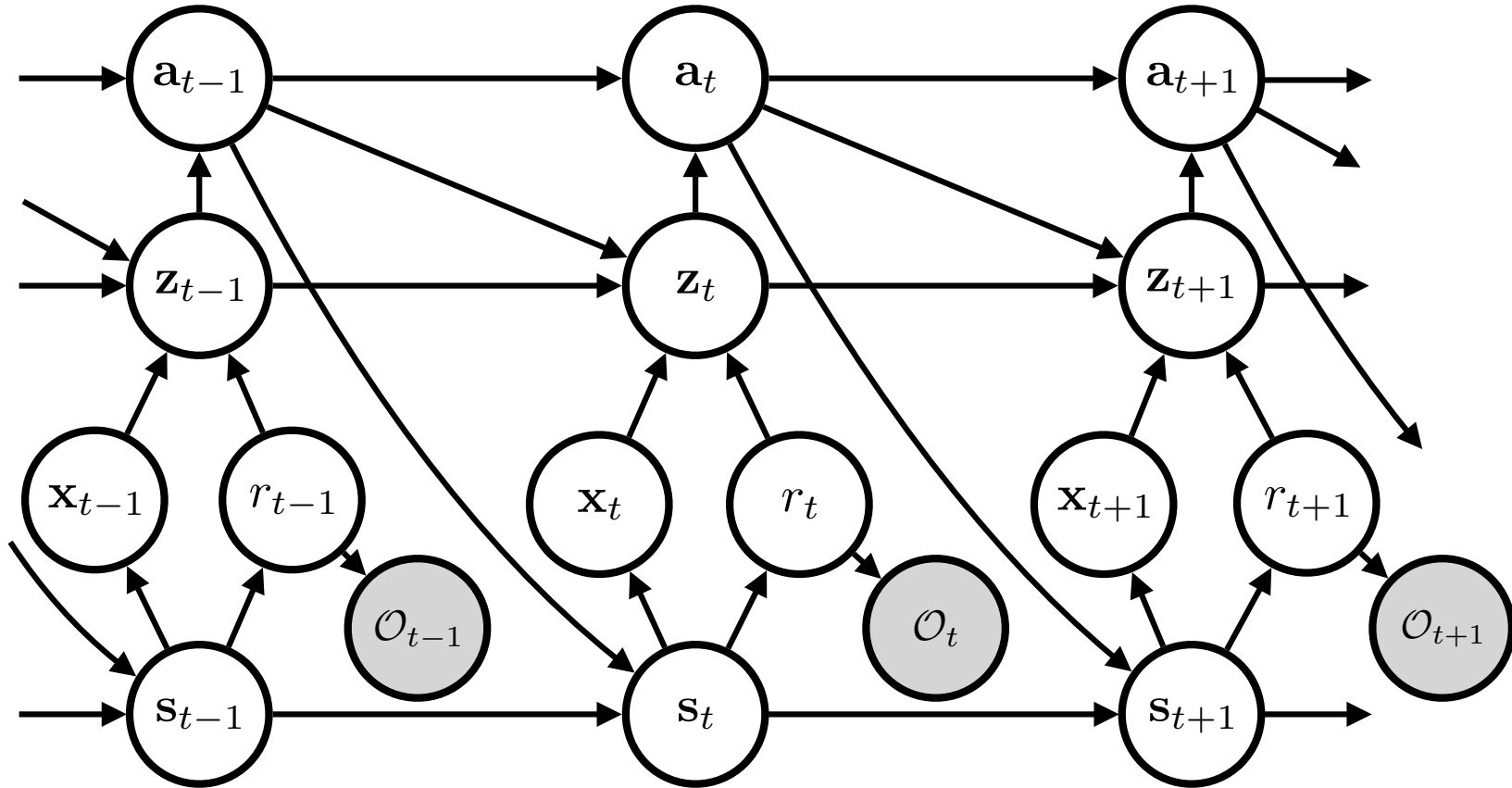
$$p_a(\mathbf{a}_{1:T}, \mathbf{z}_{1:T} | \mathbf{x}_{1:T}, r_{1:T}) = \prod_{t=1}^T \underbrace{p_a(\mathbf{a}_t | \mathbf{a}_{<t}, \mathbf{z}_{\leq t}, \mathbf{x}_{\leq t}, r_{\leq t})}_{\text{action prior}} \underbrace{p_a(\mathbf{z}_t | \mathbf{a}_{<t}, \mathbf{z}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t})}_{\text{internal state prior}}$$

OPTIMALITY



$$p(\mathcal{O}_{1:T} | r_{1:T}) = \prod_{t=1}^T \underbrace{p(\mathcal{O}_t | r_t)}_{\text{cond. likelihood of optimality}}$$

ENVIRONMENT-AGENT-OPTIMALITY



$$p(\mathbf{x}_{1:T}, r_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathbf{z}_{1:T}, \mathcal{O}_{1:T})$$

joint distribution

LEARNING

The agent's distributions are parameterized by θ .

Maximum Log-Likelihood Objective

$$\begin{aligned}\theta^* &= \arg \max_{\theta} \mathbb{E}_{\mathcal{O}_{1:T} \sim \delta(\mathbf{1})} [\log p(\mathcal{O}_{1:T})] \\ &= \arg \max_{\theta} \log p(\mathcal{O}_{1:T} = \mathbf{1}).\end{aligned}$$

where

$$p(\mathcal{O}_{1:T}) = \int p(\mathbf{x}_{1:T}, r_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathbf{z}_{1:T}, \mathcal{O}_{1:T}) d\mathbf{x}_{1:T} dr_{1:T} d\mathbf{s}_{1:T} d\mathbf{a}_{1:T} d\mathbf{z}_{1:T}$$

VARIATIONAL INFERENCE

We cannot evaluate $\log p(\mathcal{O}_{1:T})$ due to the intractable marginalization.

Introduce a structured approximate posterior, q :

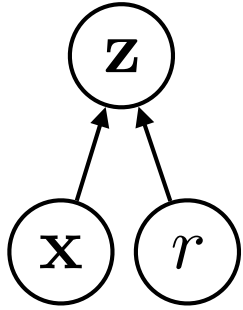
$$q(\mathbf{z}_{1:T}, \mathbf{a}_{1:T} | \mathbf{x}_{1:T}, r_{1:T}, \mathcal{O}_{1:T}) = \prod_{t=1}^T \underbrace{q(\mathbf{z}_t | \mathbf{z}_{<t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T})}_{\text{internal state}} \cdot \underbrace{q(\mathbf{a}_t | \mathbf{z}_{\leq t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T})}_{\text{action}}$$

This results in a lower bound on the objective:

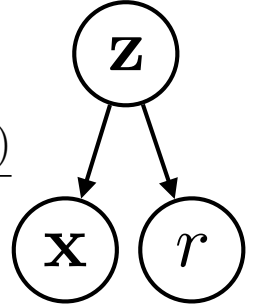
$$\mathcal{L}(q) \leq \log p(\mathcal{O}_{1:T})$$

GENERATIVE AGENT LOWER BOUND

Convert state estimation into a generative mapping using Bayes' Rule:



$$p_a(\mathbf{z}_t | \mathbf{a}_{<t}, \mathbf{z}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}) = \frac{p_a(\mathbf{x}_t, r_t | \mathbf{a}_{<t}, \mathbf{z}_{\leq t}, \mathbf{x}_{<t}, r_{<t}) p_a(\mathbf{z}_t | \mathbf{a}_{<t}, \mathbf{z}_{<t}, \mathbf{x}_{<t}, r_{<t})}{p_a(\mathbf{x}_t, r_t | \mathbf{a}_{<t}, \mathbf{z}_{<t}, \mathbf{x}_{<t}, r_{<t})}$$



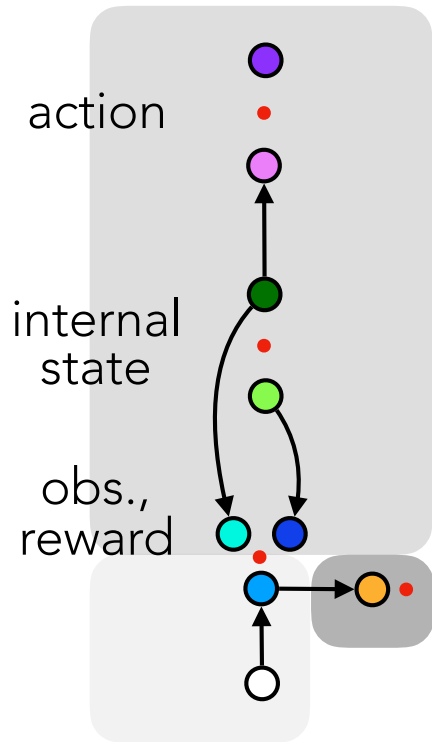
Plugging this into the bound yields:

$$\mathcal{L} = \mathbb{E}_{\mathbf{s}, \mathbf{x}, r \sim p_e, \mathbf{z}, \mathbf{a} \sim q} \left[\underbrace{\sum_{t=1}^T r_t}_{\text{reward}} + \underbrace{\log \frac{p_a(\mathbf{x}_t, r_t | \mathbf{a}_{<t}, \mathbf{z}_{\leq t}, \mathbf{x}_{<t}, r_{<t})}{p_a(\mathbf{x}_t, r_t | \mathbf{a}_{<t}, \mathbf{z}_{<t}, \mathbf{x}_{<t}, r_{<t})}}_{\text{information gain}} - \underbrace{\log \frac{q(\mathbf{z}_t | \mathbf{z}_{<t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T})}{p_a(\mathbf{z}_t | \mathbf{a}_{<t}, \mathbf{z}_{<t}, \mathbf{x}_{<t}, r_{<t})}}_{\text{internal state consistency}} - \underbrace{\log \frac{q(\mathbf{a}_t | \mathbf{z}_{\leq t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T})}{p_a(\mathbf{a}_t | \mathbf{a}_{<t}, \mathbf{z}_{\leq t}, \mathbf{x}_{\leq t}, r_{\leq t})}}_{\text{action consistency}} \right]$$

generative agent bound

GENERATIVE AGENT LOWER BOUND

Computation Graph:



• log likelihood, log ratio

Env.

- s_t dynamics
- x_t, r_t emission

Optimality

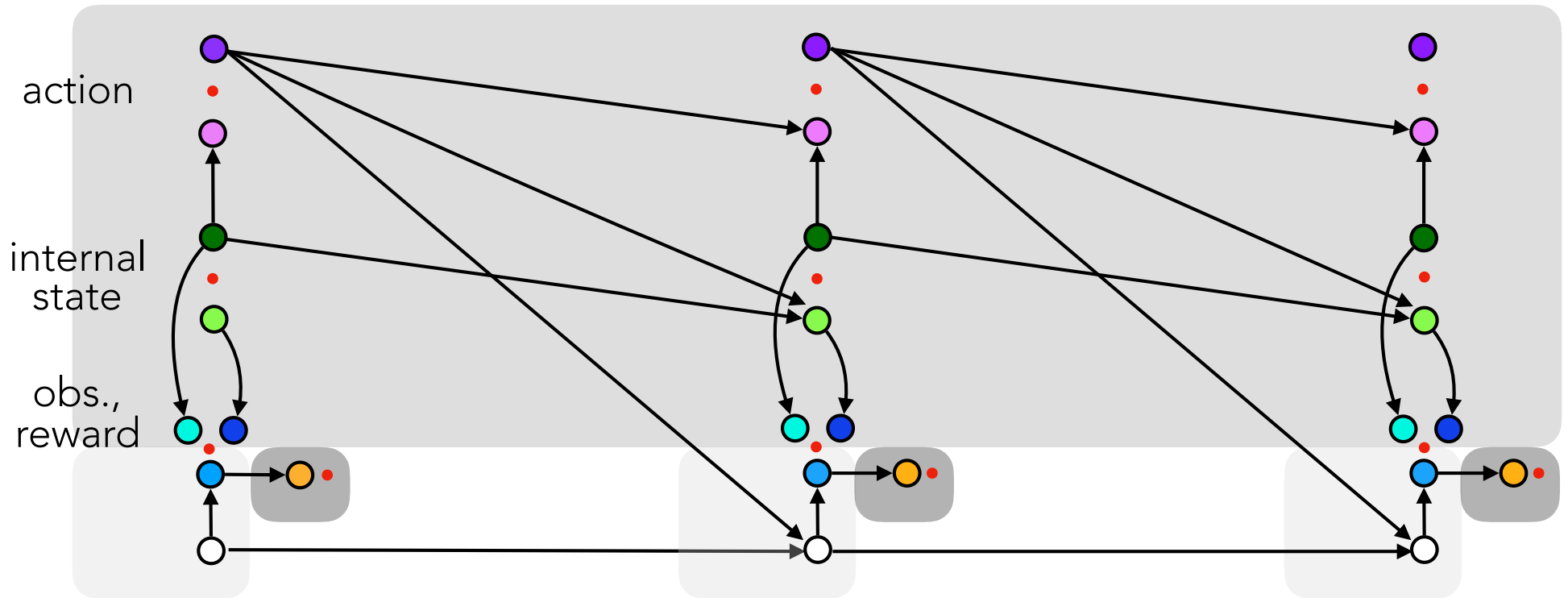
- \mathcal{O}_t cond. likelihood

Agent

- z_t prior
- z_t approx. post.
- a_t prior
- a_t approx. post.
- x_t, r_t cond. likelihood
- x_t, r_t marginal likelihood

GENERATIVE AGENT LOWER BOUND

Computation Graph:



• log likelihood, log ratio

Env.

- s_t dynamics
- x_t, r_t emission

Optimality

- \mathcal{O}_t cond. likelihood

Agent

- z_t prior
- z_t approx. post.
- a_t prior
- a_t approx. post.
- x_t, r_t cond. likelihood
- x_t, r_t marginal likelihood

VARIATIONAL EM

```
while  $\theta$  not converged:
```

```
   $\mathbf{x}_1, r_1, \mathbf{s}_1 \sim p_e(\mathbf{x}_1|\mathbf{s}_1)p_e(r_1|\mathbf{s}_1)p_e(\mathbf{s}_1)$ 
```

```
  for  $t = 1 \dots T$ :
```

```
    # inference (simulated rollouts)
```

```
     $q(\mathbf{z}_t|\mathbf{z}_{<t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T}) \leftarrow \arg \max_q \mathcal{L}_{t:T}$ 
```

```
     $q(\mathbf{a}_t|\mathbf{z}_{\leq t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T}) \leftarrow \arg \max_q \mathcal{L}_{t:T}$ 
```

```
    # interaction
```

```
     $\mathbf{a}_t \sim q(\mathbf{a}_t|\mathbf{z}_{\leq t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T})$ 
```

```
     $\mathbf{x}_{t+1}, r_{t+1}, \mathbf{s}_{t+1} \sim p_e(\mathbf{x}_{t+1}|\mathbf{s}_{t+1})p_e(r_{t+1}|\mathbf{s}_{t+1})p_e(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$ 
```

```
  # learning
```

```
   $\theta \leftarrow \theta + \alpha \nabla_{\theta} \mathcal{L}$ 
```

We typically cannot simulate the environment to evaluate $\mathcal{L}_{t:T}$

ACTION INFERENCE VIA **PLANNING**

Estimate $\mathcal{L}_{t:T}$ using the generative agent's internal model.

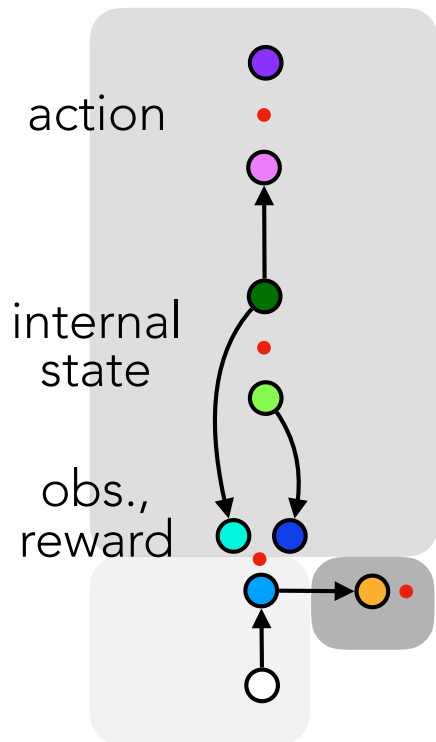
Replace p_e and q in future terms with p_a :

$$\hat{\mathcal{L}}_{t+1:T} = \mathbb{E}_{\mathbf{x}, r, \mathbf{z}, \mathbf{a} \sim p_a} \left[\overset{\text{reward}}{\sum_{\tau=t+1}^T r_{\tau}} + \log \overset{\text{mutual information}}{\frac{p_a(\mathbf{x}_{\tau}, r_{\tau} | \mathbf{a}_{<\tau}, \mathbf{z}_{\leq \tau}, \mathbf{x}_{<\tau}, r_{<\tau})}{p_a(\mathbf{x}_{\tau}, r_{\tau} | \mathbf{a}_{<\tau}, \mathbf{z}_{<\tau}, \mathbf{x}_{<\tau}, r_{<\tau})}} \right]$$

planning bound

PLANNING

Computation Graph:



• log likelihood, log ratio

Env.

- s_t dynamics
- x_t, r_t emission

Optimality

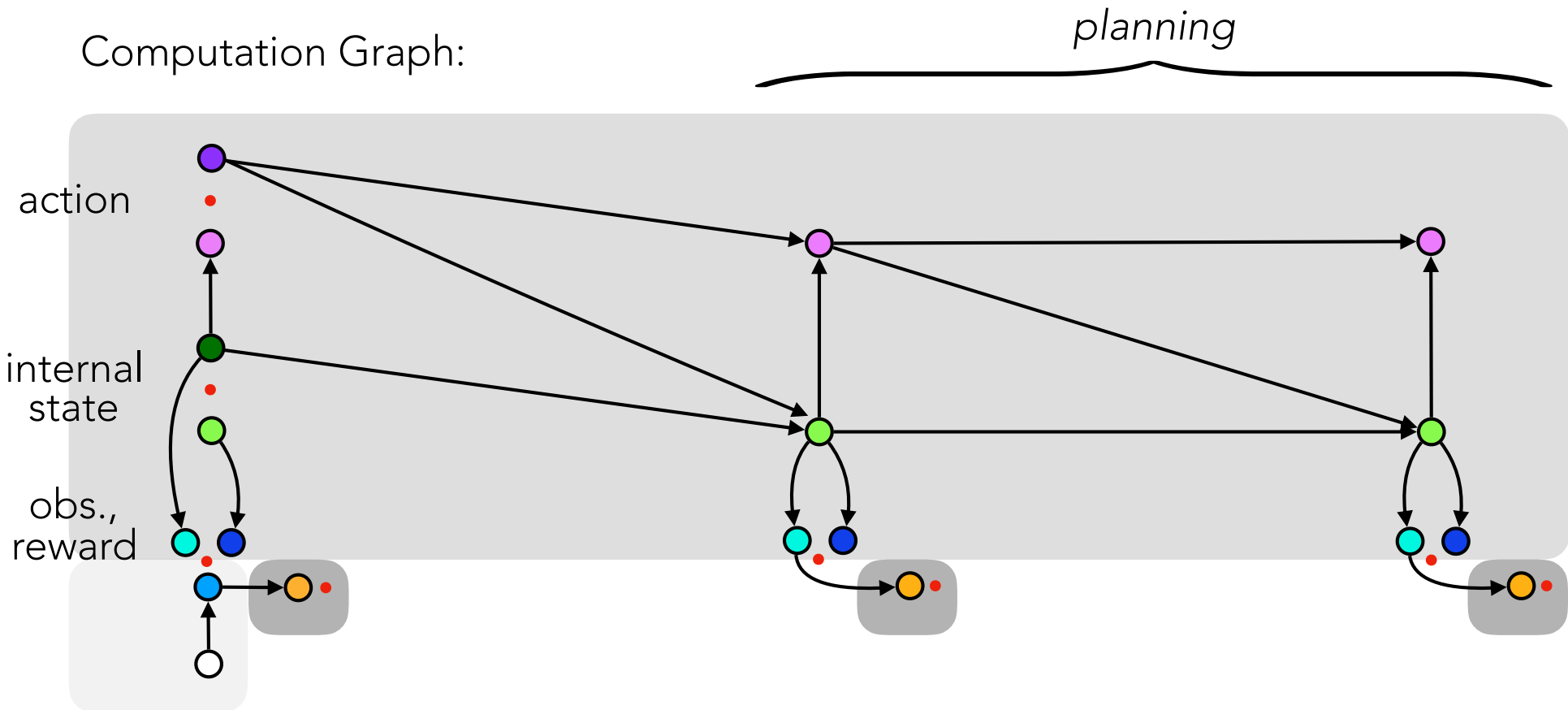
- \mathcal{O}_t cond. likelihood

Agent

- z_t prior
- z_t approx. post.
- a_t prior
- a_t approx. post.
- x_t, r_t cond. likelihood
- x_t, r_t marginal likelihood

PLANNING

Computation Graph:



• log likelihood, log ratio

Env.

- s_t dynamics
- x_t, r_t emission

Optimality

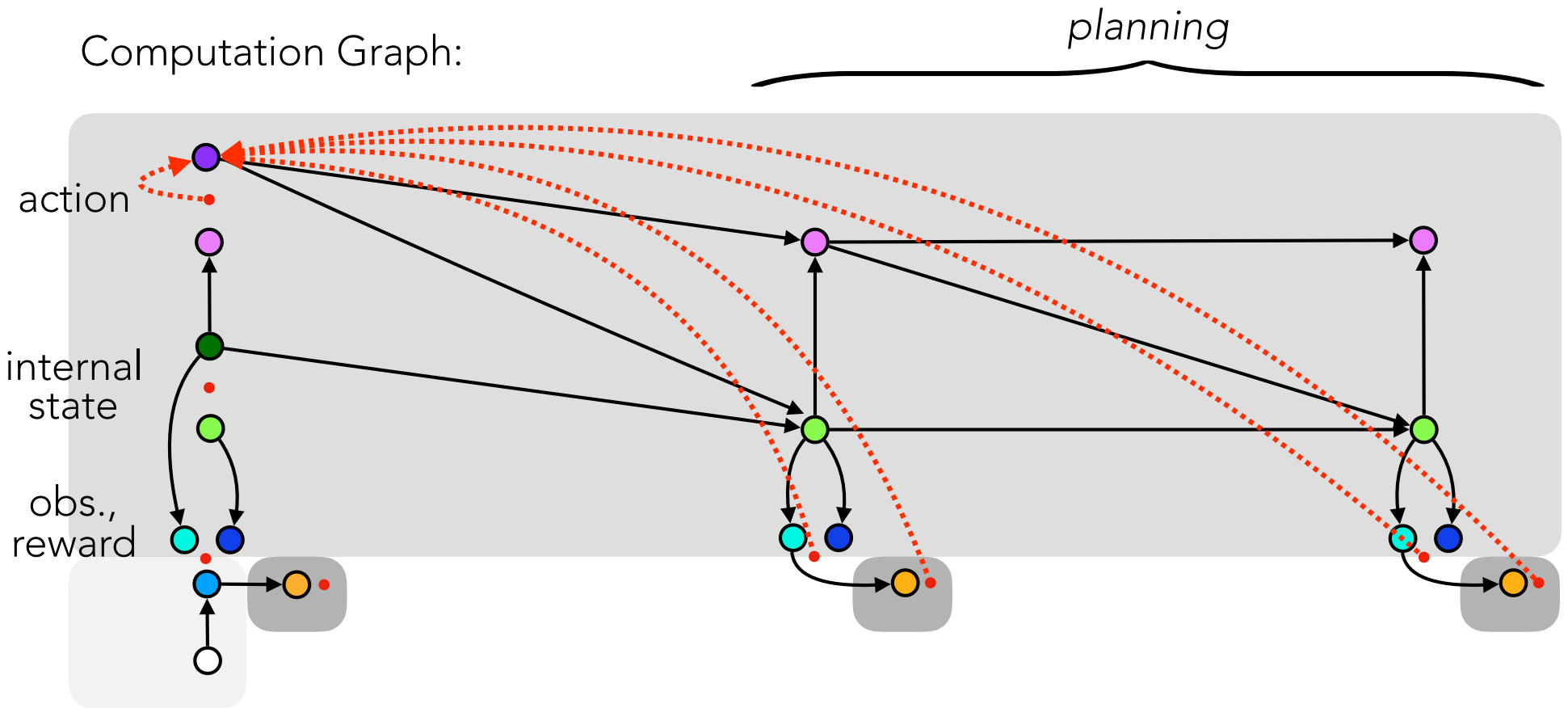
- \mathcal{O}_t cond. likelihood

Agent

- z_t prior
- z_t approx. post.
- a_t prior
- a_t approx. post.
- x_t, r_t cond. likelihood
- x_t, r_t marginal likelihood

PLANNING

Computation Graph:



Action selection depends on:

MODEL-FREE PRIOR

$$p_a(\mathbf{a}_t | \mathbf{a}_{<t}, \mathbf{z}_{\leq t}, \mathbf{x}_{\leq t}, r_{\leq t})$$

MODEL-BASED LIKELIHOOD

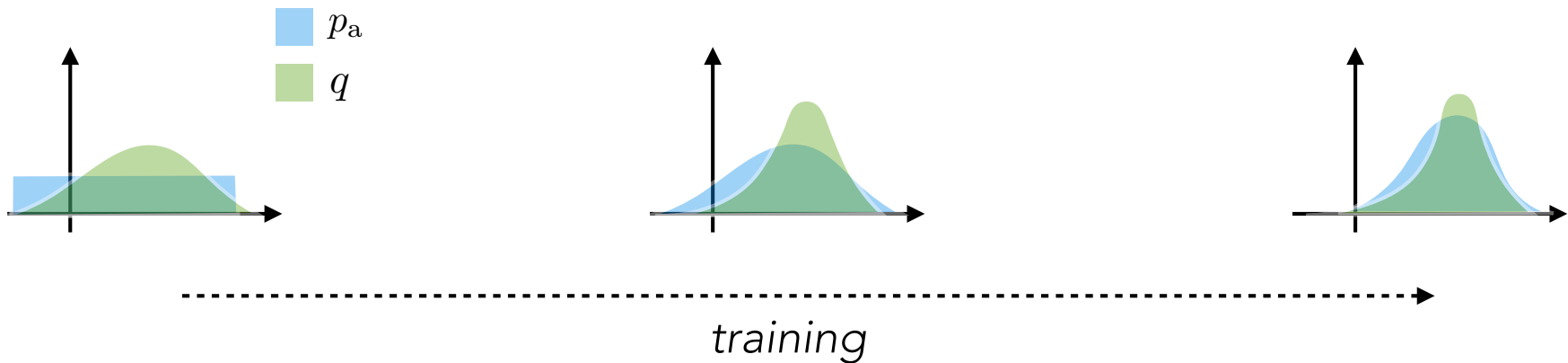
$$\hat{\mathcal{L}}_{t+1:T}$$

PLANNING *DISTILLATION*

“**Habits** are sometimes said to be controlled by *antecedent stimuli*, whereas **goal-directed behavior** is said to be controlled by *its consequences*.”

Actions and habits: the development of behavioural autonomy, Dickinson, 1985
Reinforcement learning: An introduction, Sutton & Barto, 2018

Learning $p_a(\mathbf{a}_t | \mathbf{a}_{<t}, \mathbf{z}_{\leq t}, \mathbf{x}_{\leq t}, r_{\leq t})$ will shift it toward $q(\mathbf{a}_t | \mathbf{z}_{\leq t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T})$



Model-based planning will be “distilled” into a model-free policy, forming a *habit*.

MODELING THE ENVIRONMENT

GENERATIVE AGENT:

$$\mathcal{L} = \mathbb{E}_{\mathbf{s}, \mathbf{x}, r \sim p_e, \mathbf{z}, \mathbf{a} \sim q} \left[\sum_{t=1}^T r_t + \log \frac{p_a(\mathbf{x}_t, r_t | \mathbf{a}_{<t}, \mathbf{z}_{\leq t}, \mathbf{x}_{<t}, r_{<t})}{p_a(\mathbf{x}_t, r_t | \mathbf{a}_{<t}, \mathbf{z}_{<t}, \mathbf{x}_{<t}, r_{<t})} - \log \frac{q(\mathbf{z}_t | \mathbf{z}_{<t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T})}{p_a(\mathbf{z}_t | \mathbf{a}_{<t}, \mathbf{z}_{<t}, \mathbf{x}_{<t}, r_{<t})} - \log \frac{q(\mathbf{a}_t | \mathbf{z}_{\leq t}, \mathbf{a}_{<t}, \mathbf{x}_{\leq t}, r_{\leq t}, \mathcal{O}_{t+1:T})}{p_a(\mathbf{a}_t | \mathbf{a}_{<t}, \mathbf{z}_{\leq t}, \mathbf{x}_{\leq t}, r_{\leq t})} \right]$$

learning bound

$$\hat{\mathcal{L}}_{t+1:T} = \mathbb{E}_{\mathbf{x}, r, \mathbf{z}, \mathbf{a} \sim p_a} \left[\sum_{\tau=t+1}^T r_\tau + \log \frac{p_a(\mathbf{x}_\tau, r_\tau | \mathbf{a}_{<\tau}, \mathbf{z}_{\leq \tau}, \mathbf{x}_{<\tau}, r_{<\tau})}{p_a(\mathbf{x}_\tau, r_\tau | \mathbf{a}_{<\tau}, \mathbf{z}_{<\tau}, \mathbf{x}_{<\tau}, r_{<\tau})} \right]$$

inference (planning) bound

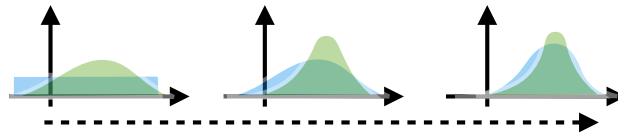
The likelihood ratio

learning: encourages the agent to learn a task-relevant model

inference: biases planning toward less stochastic/uncertain outcomes

RECAP: BENEFITS

PLANNING DISTILLATION



convert model-based planning into a model-free policy

- + fewer interactions during training
- + after training, fast to act

MODELING THE ENVIRONMENT



estimate information gain / mutual information

- + learn a more task-oriented internal state
- + improve robustness during planning



`joelouismarino.github.io`



`jmarino@caltech.edu`