

# Formal Verification of Efficiently Distilled RL Policies with Many-sided Guarantees

*Florent Delgrange, Ann Nowé, Guillermo A. Pérez*



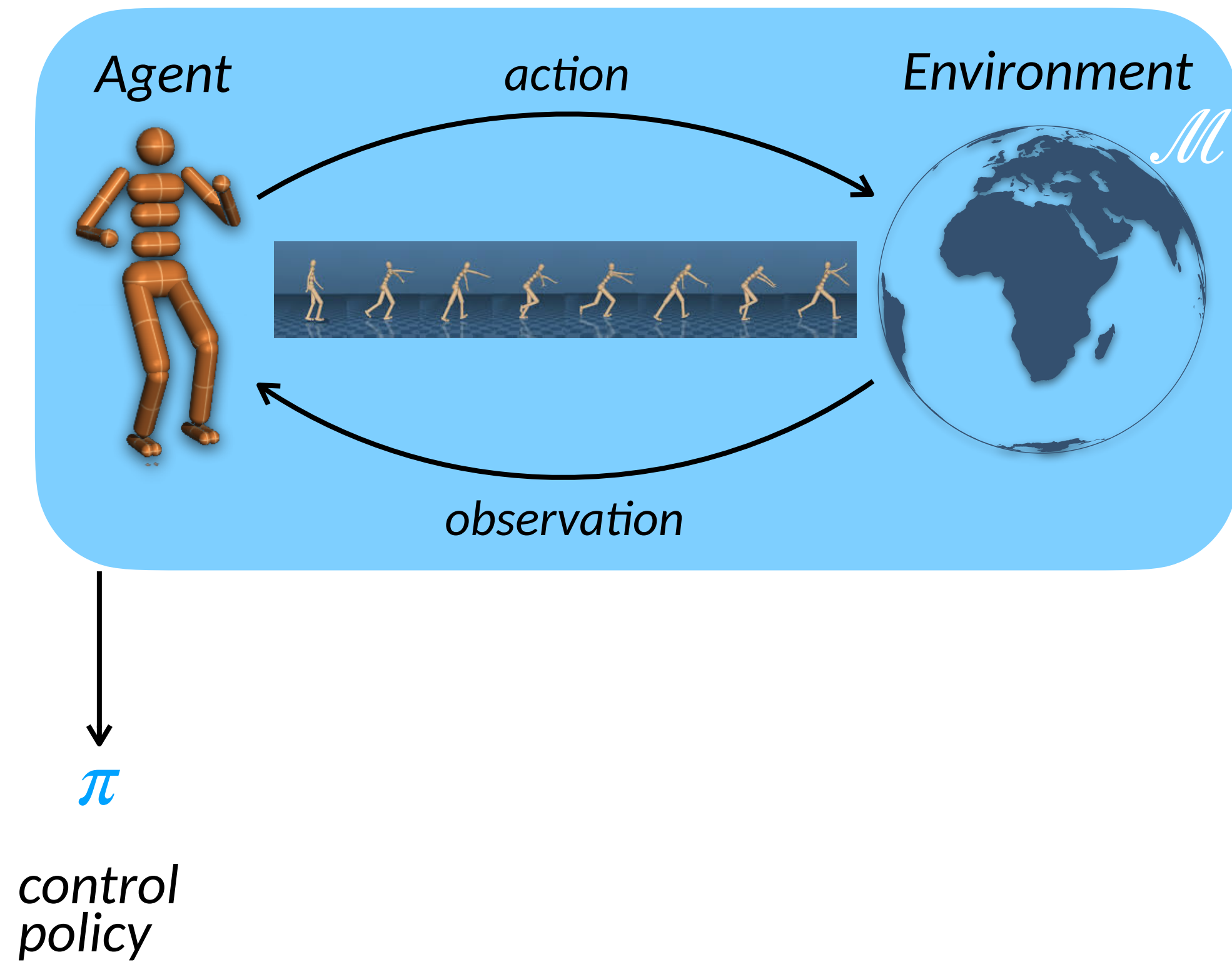
ARTIFICIAL  
INTELLIGENCE  
RESEARCH GROUP



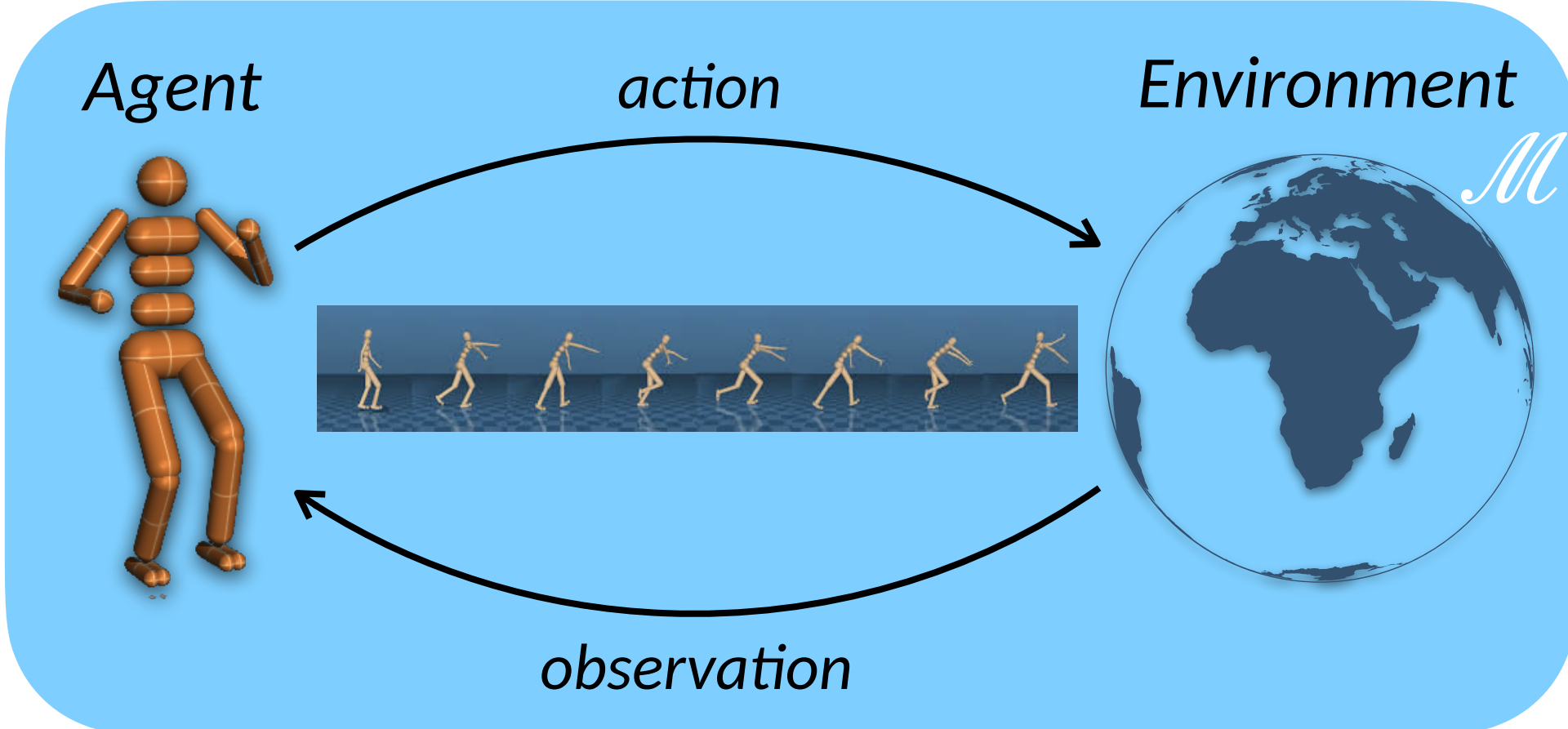
Universiteit  
Antwerpen



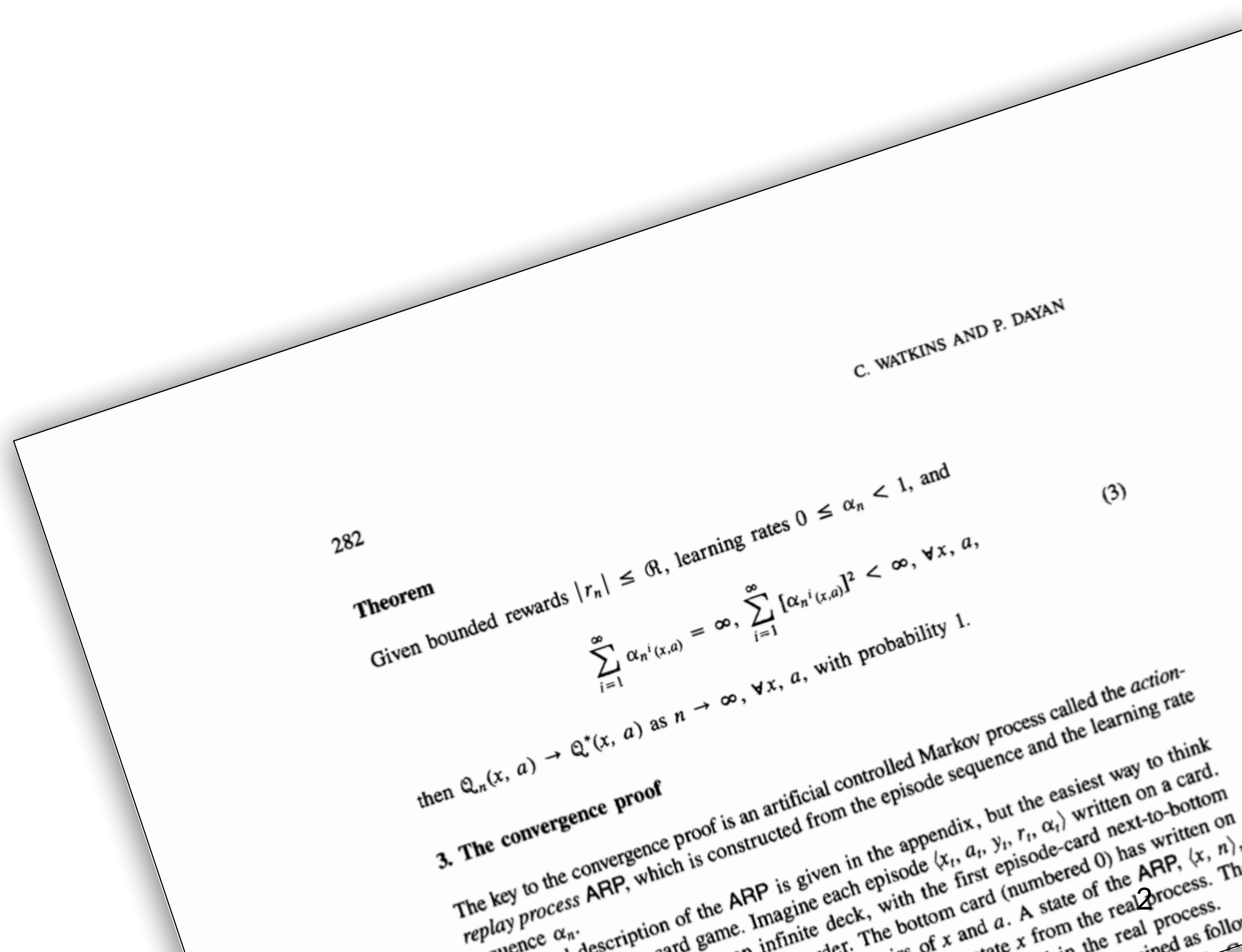
## Reinforcement Learning



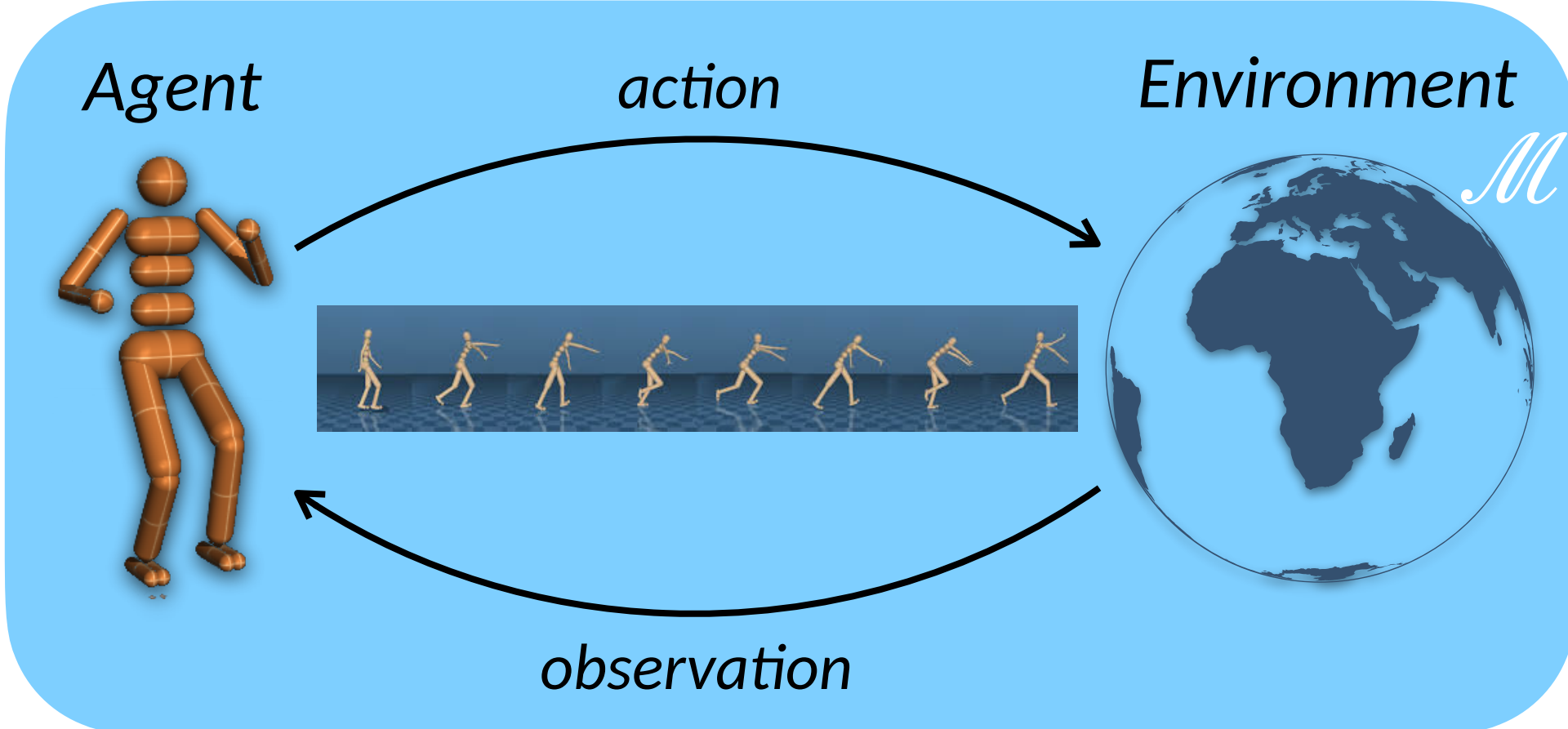
## Reinforcement Learning



control policy

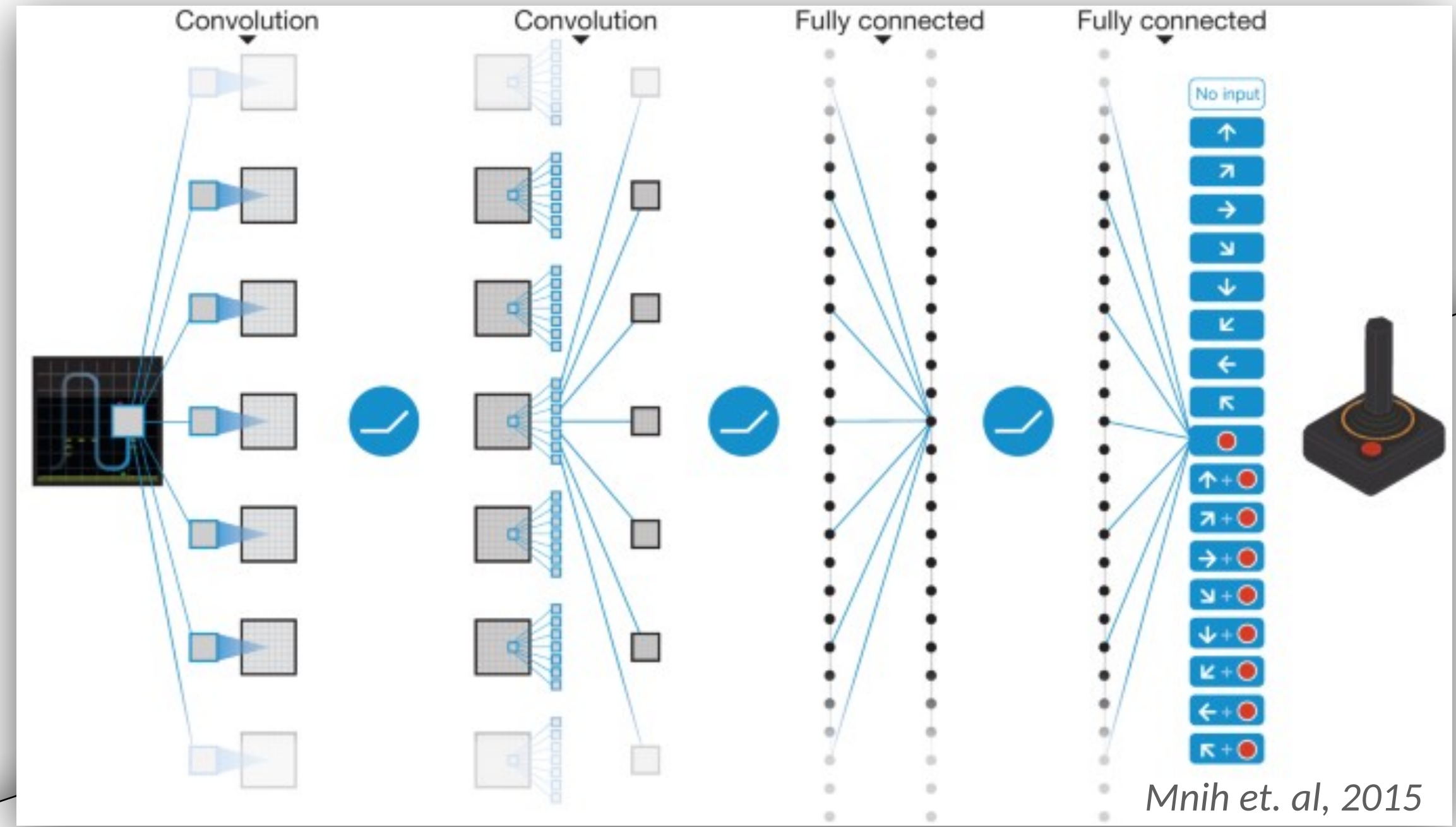


## Reinforcement Learning



$\pi$

control policy



282

**Theorem**

Given bounded rewards  $|r_n| \leq R$ , learning rates  $0 \leq \alpha_n < 1$ , and

$$\sum_{i=1}^{\infty} \alpha_n^{i(x,a)} = \infty, \sum_{i=1}^{\infty} [\alpha_n^{i(x,a)}]^2 < \infty, \forall x, a,$$

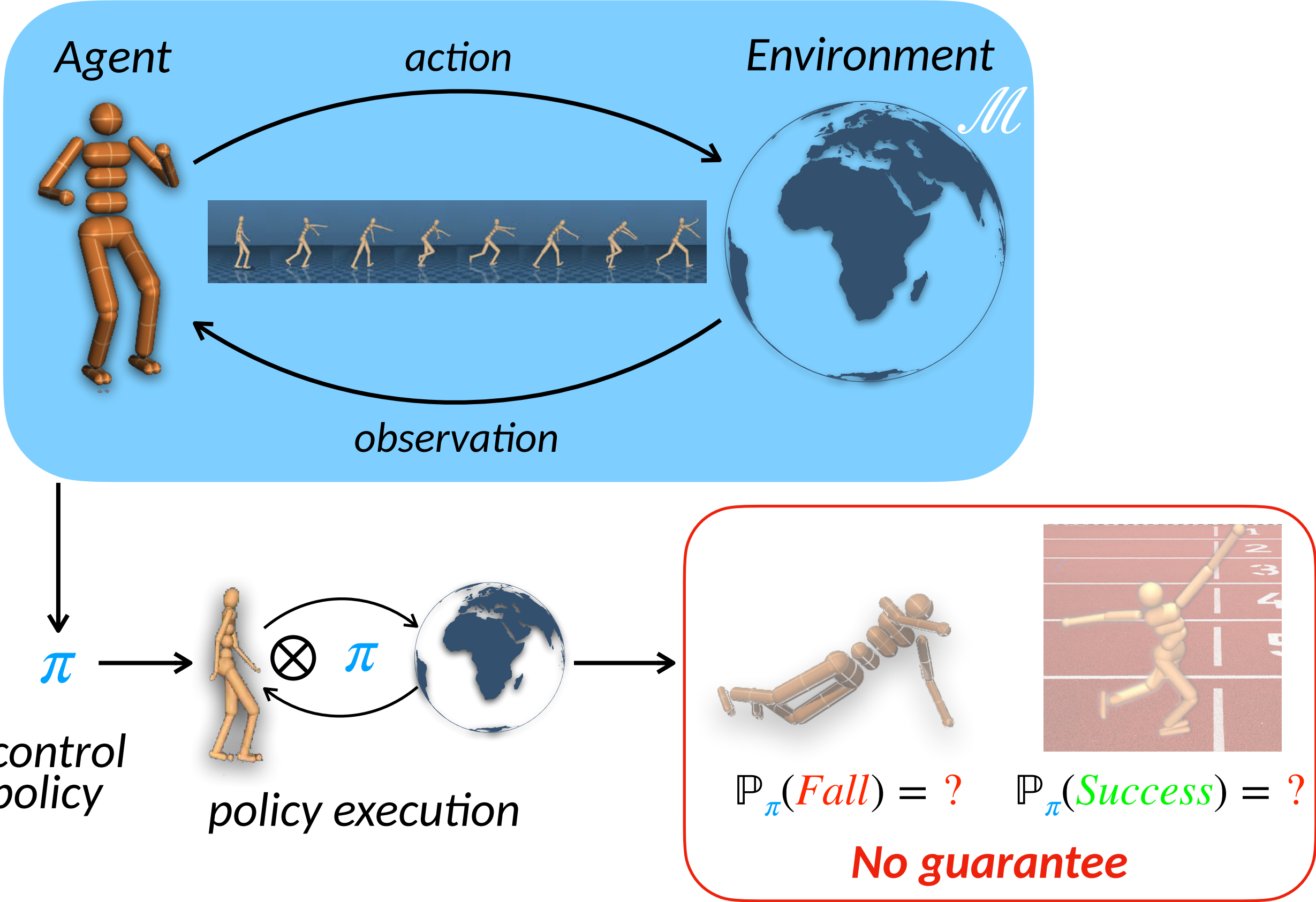
then  $Q_n(x, a) \rightarrow Q^*(x, a)$  as  $n \rightarrow \infty, \forall x, a$ , with probability 1.

**3. The convergence proof**

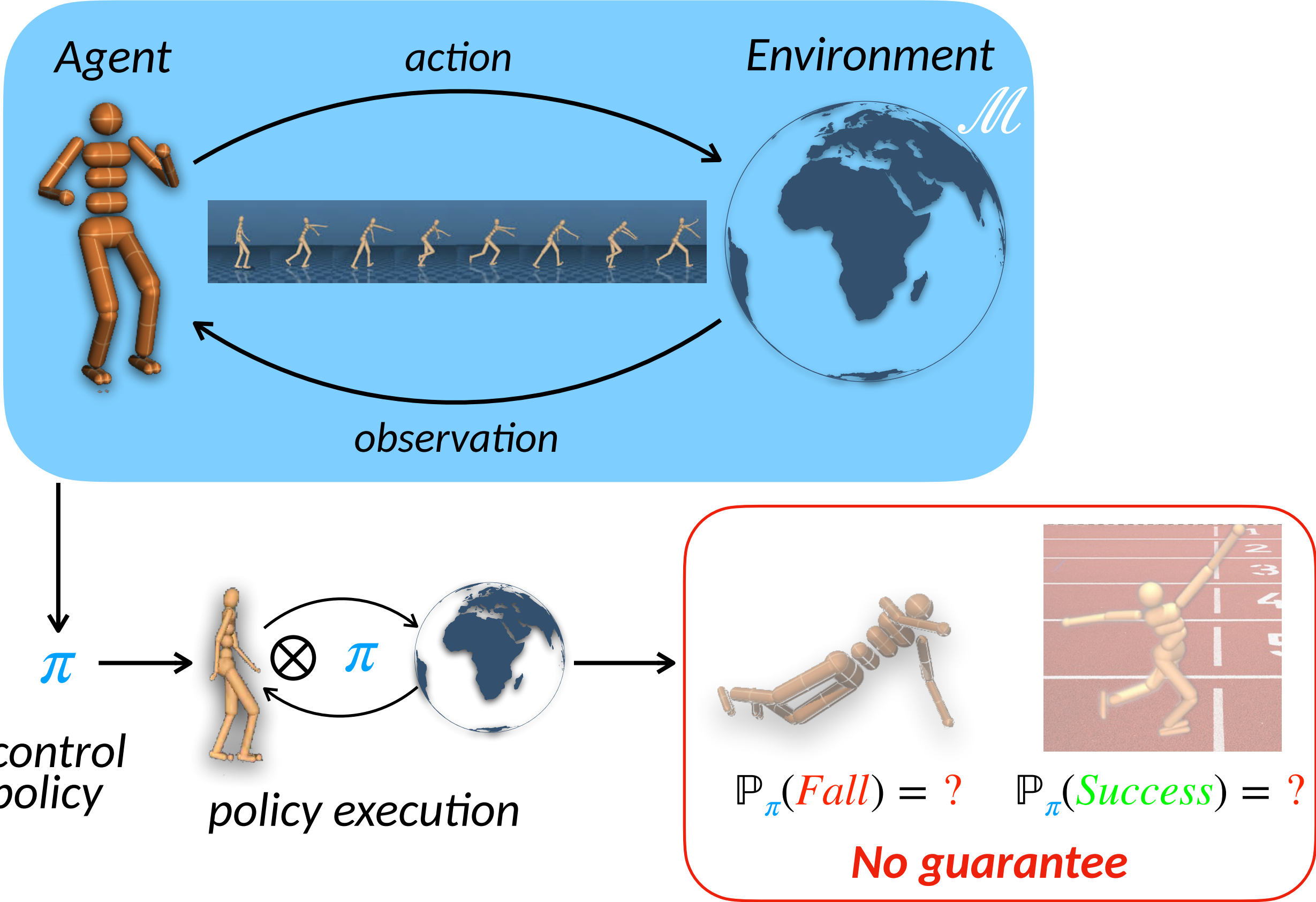
The key to the convergence proof is an artificial controlled Markov process called the *action-replay process ARP*, which is constructed from the episode sequence and the learning rate  $\alpha_n$ .

A description of the ARP is given in the appendix, but the easiest way to think of the ARP is as a card game. Imagine each episode  $(x_i, a_i, y_i, r_i, \alpha_i)$  written on a card. The bottom card (numbered 0) has written on it the state  $x$  and  $a$ . A state of the ARP,  $(x, n)$ , is a card from the real process. The ARP is defined as follows:

## Reinforcement Learning

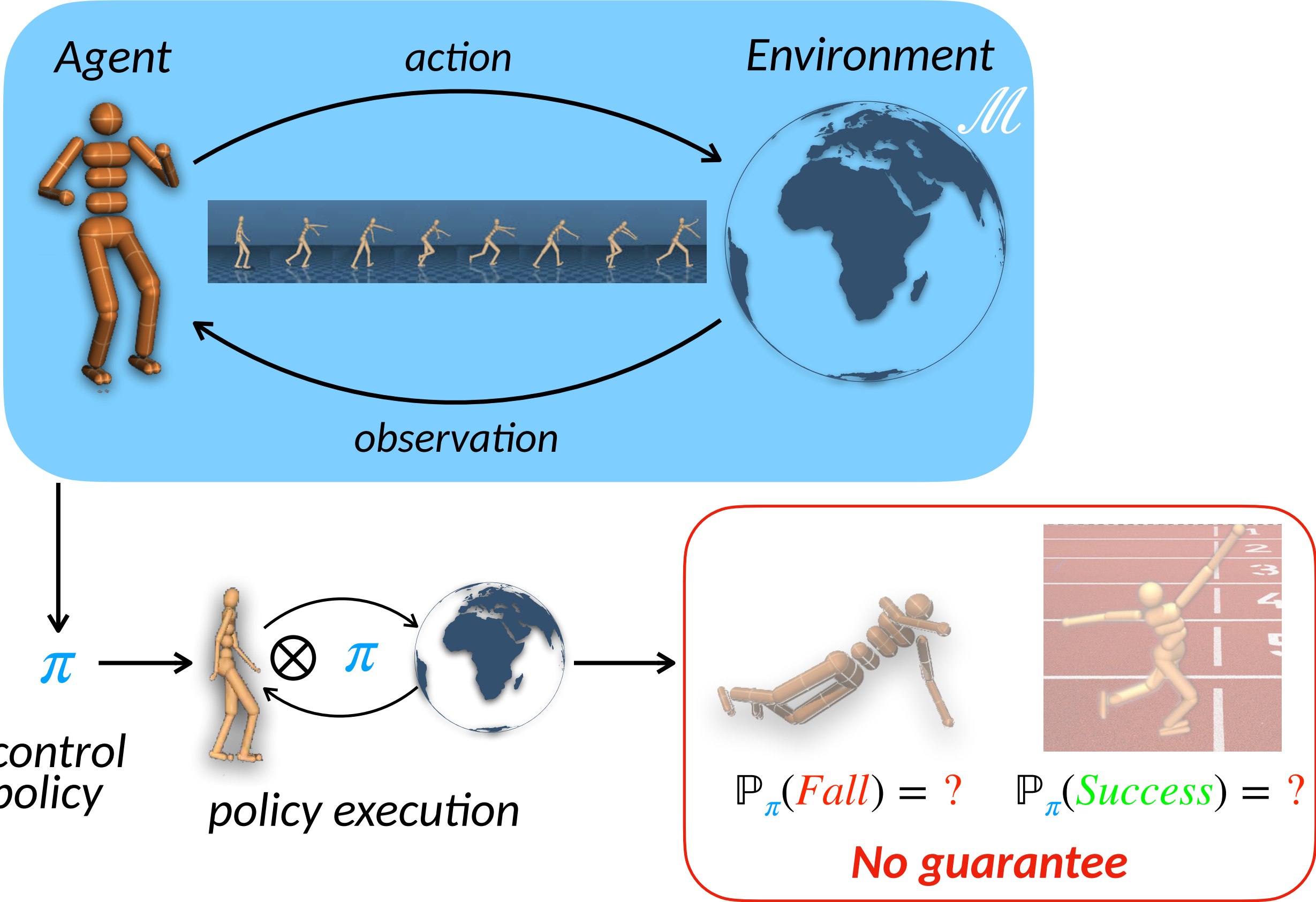


## Reinforcement Learning



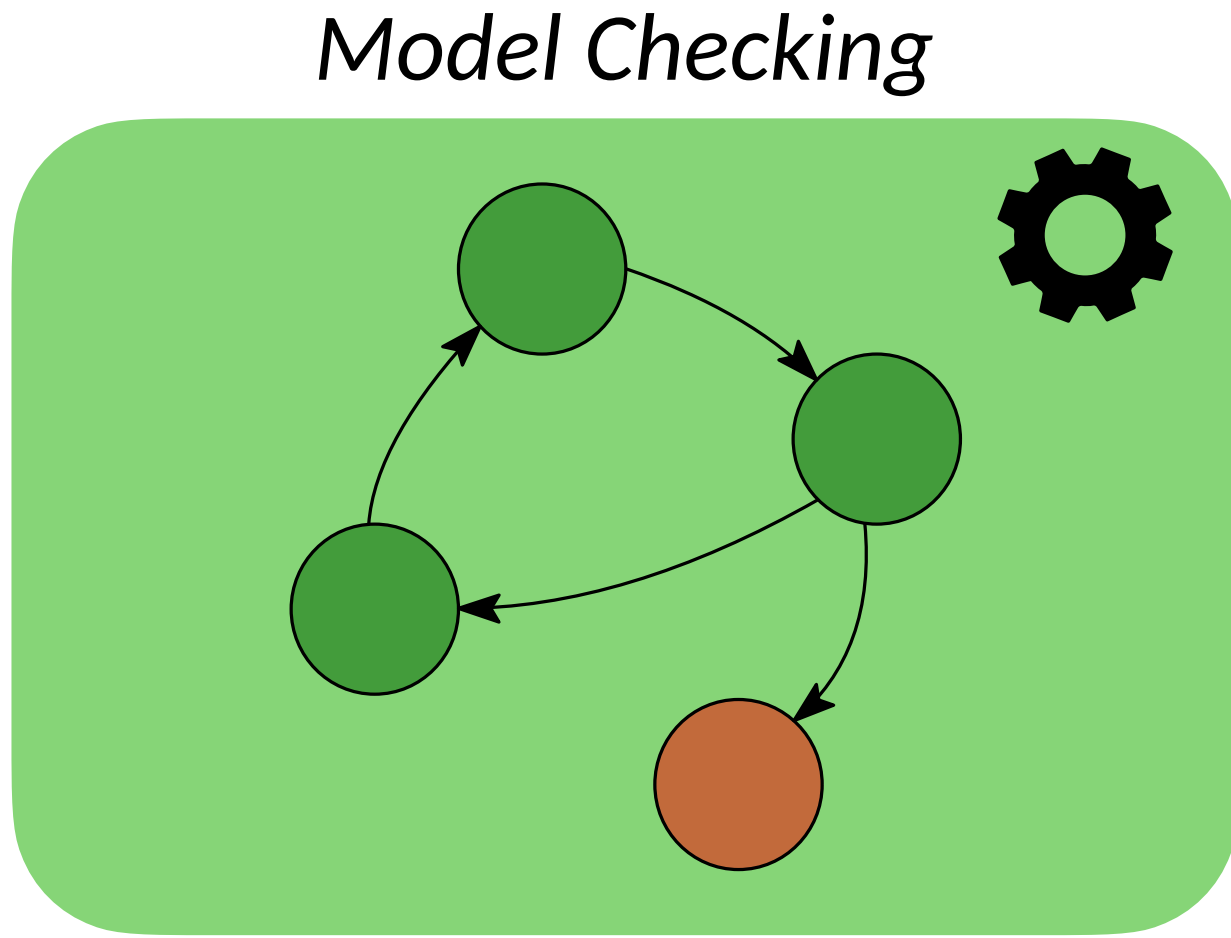
- Continuous state/action spaces
- Unknown environment

## Reinforcement Learning



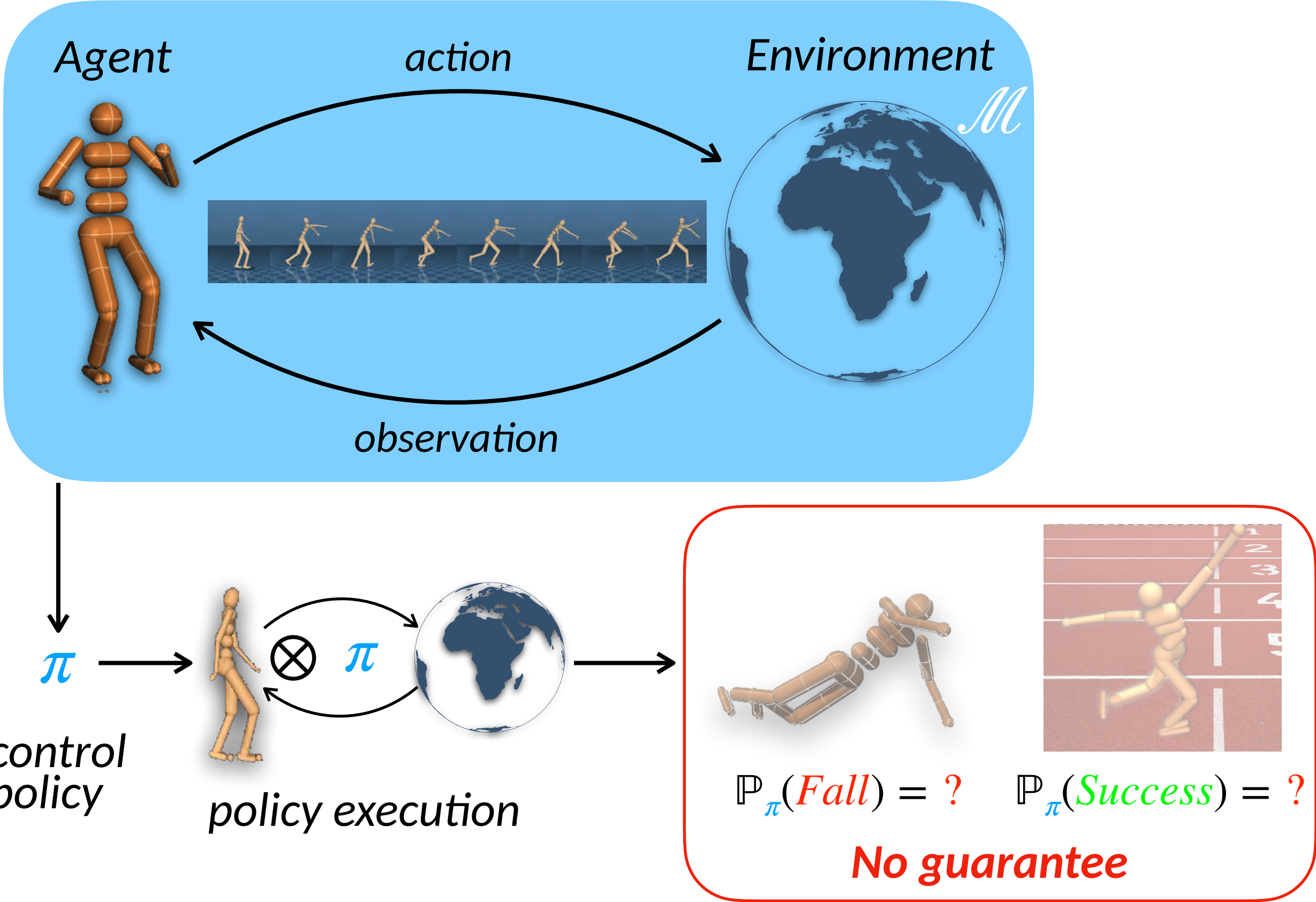
- Continuous state/action spaces
- Unknown environment

## Formal Guarantees



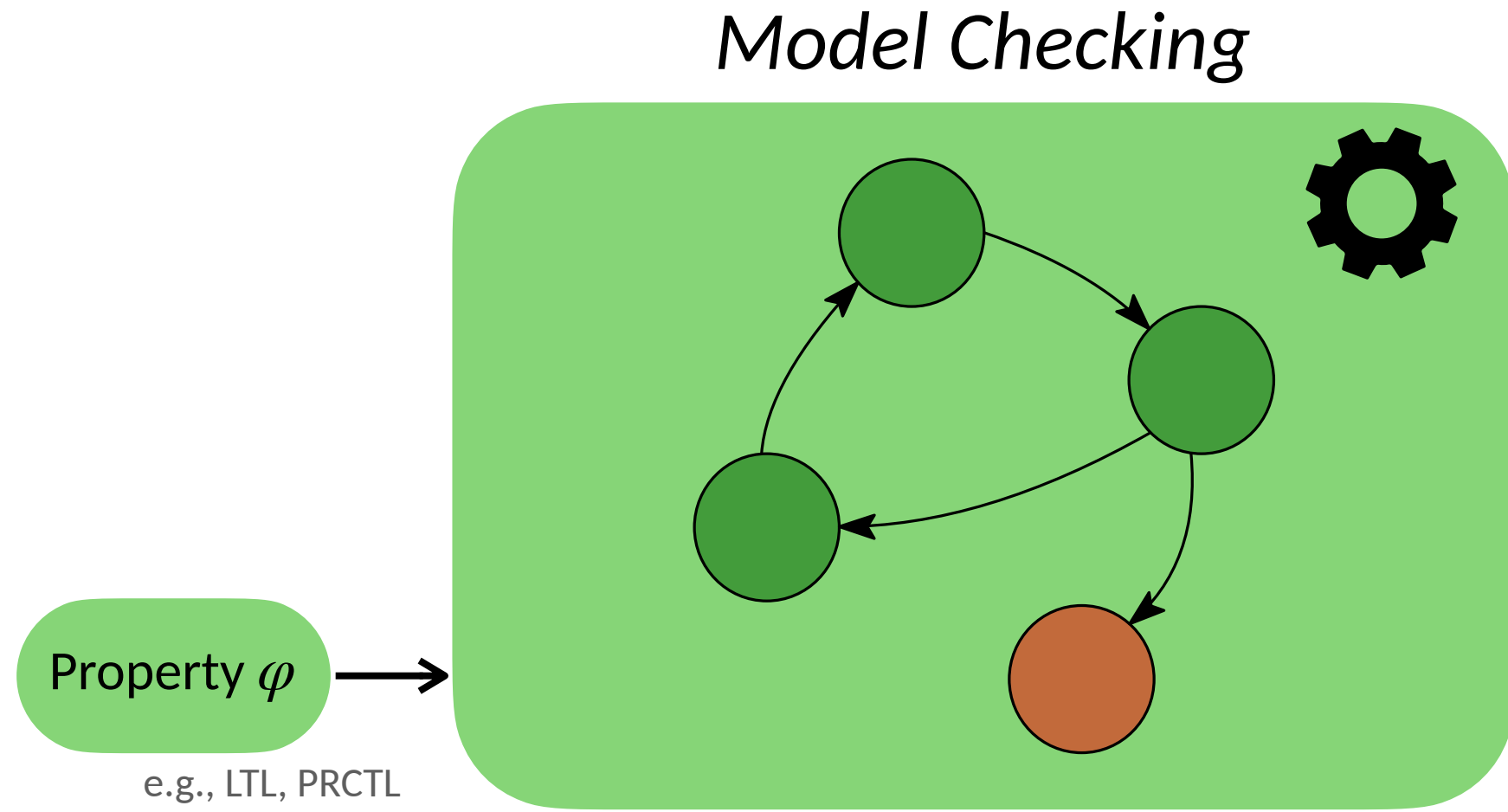
- Full knowledge of the model of the interaction

## Reinforcement Learning



- Continuous state/action spaces
- Unknown environment

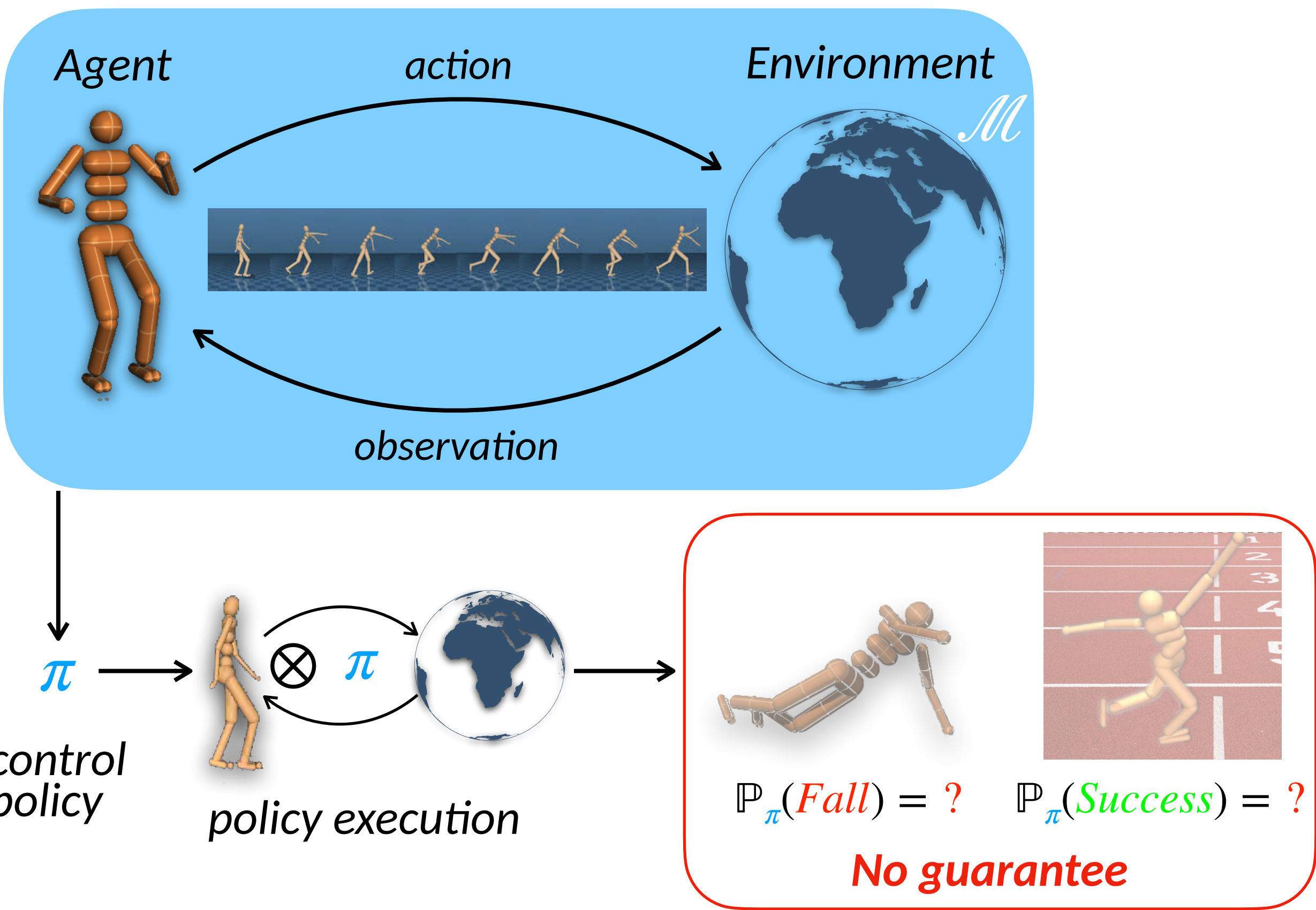
## Formal Guarantees



- Full knowledge of the model of the interaction

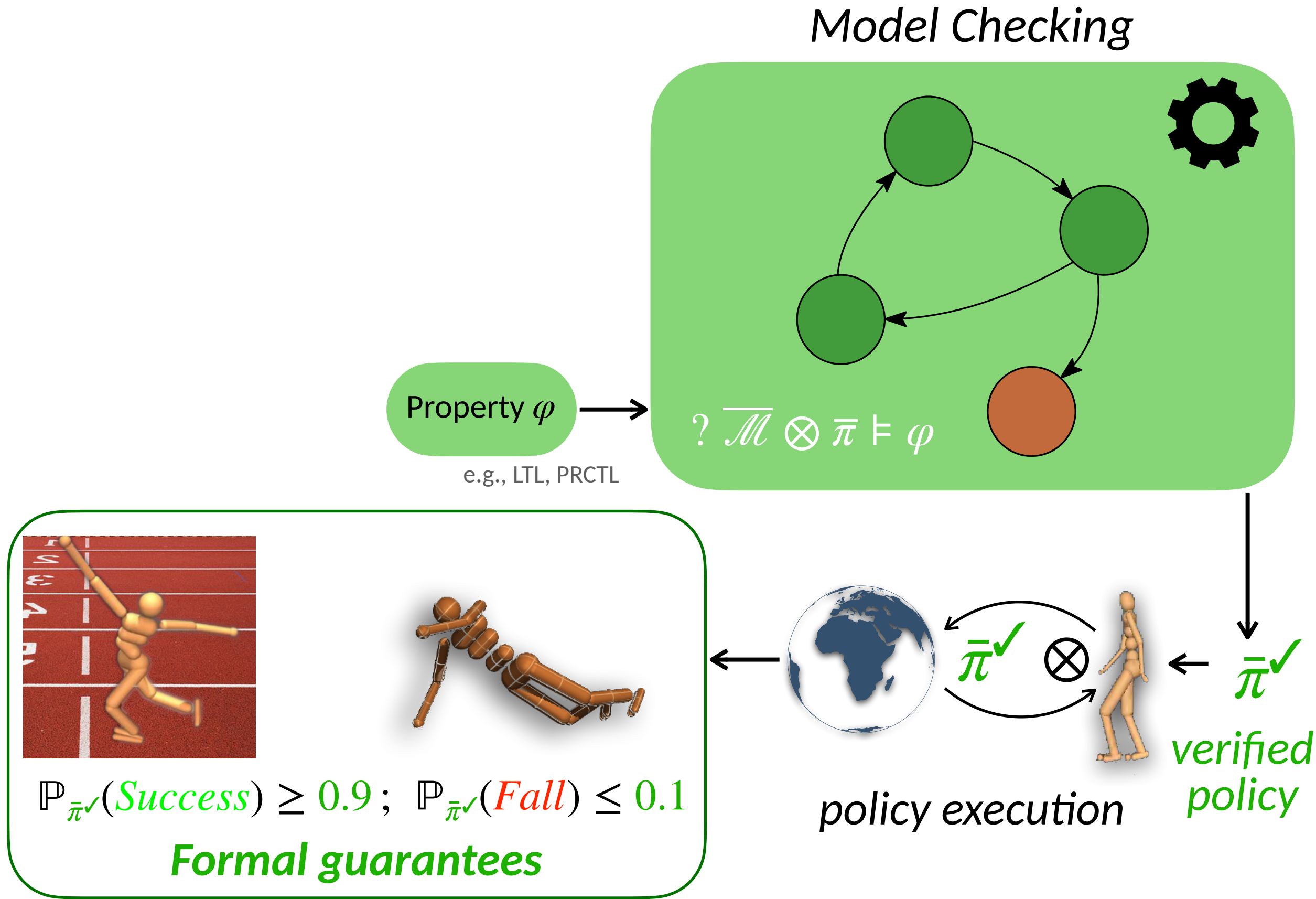


## Reinforcement Learning



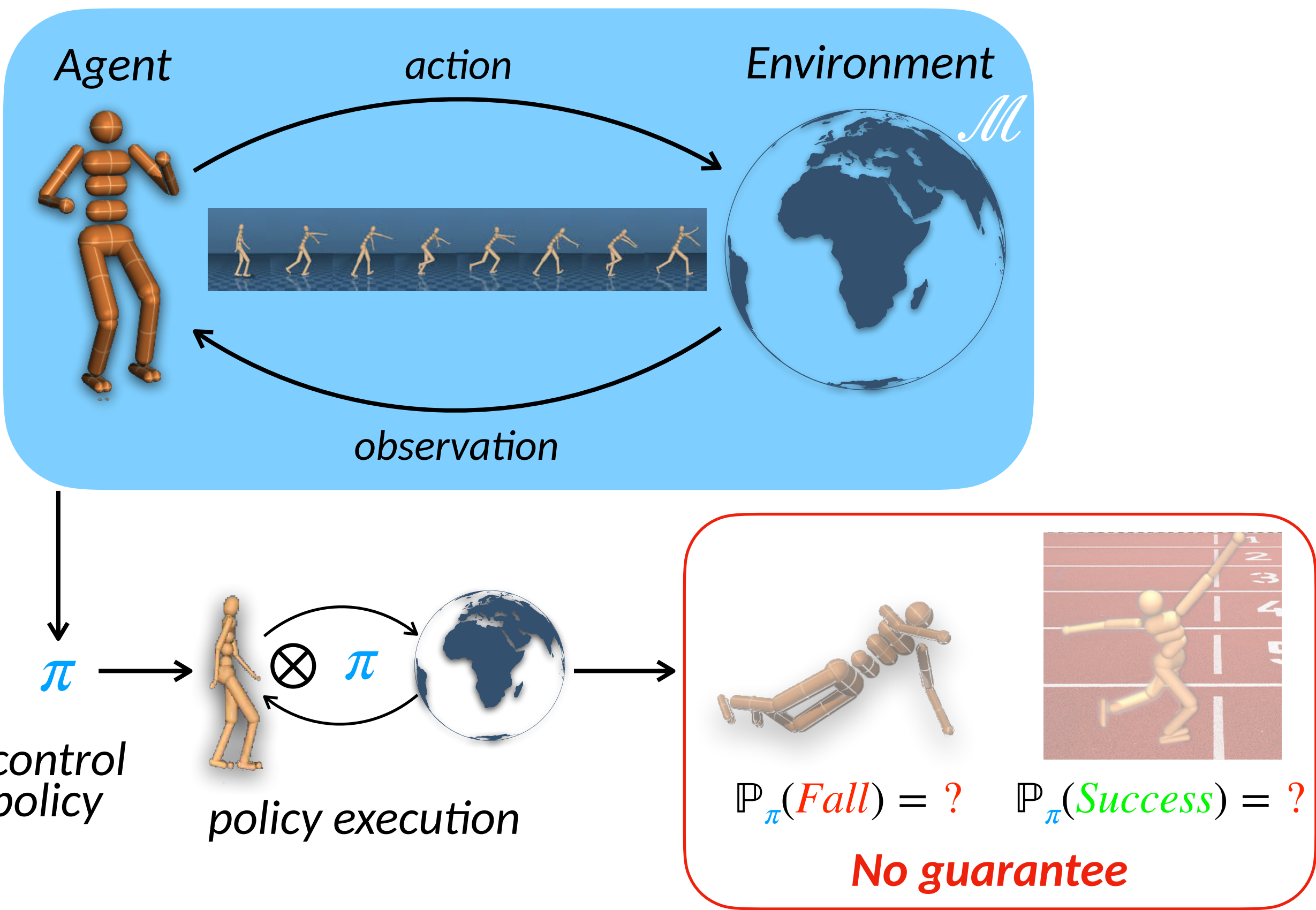
- Continuous state/action spaces
- Unknown environment

## Formal Guarantees



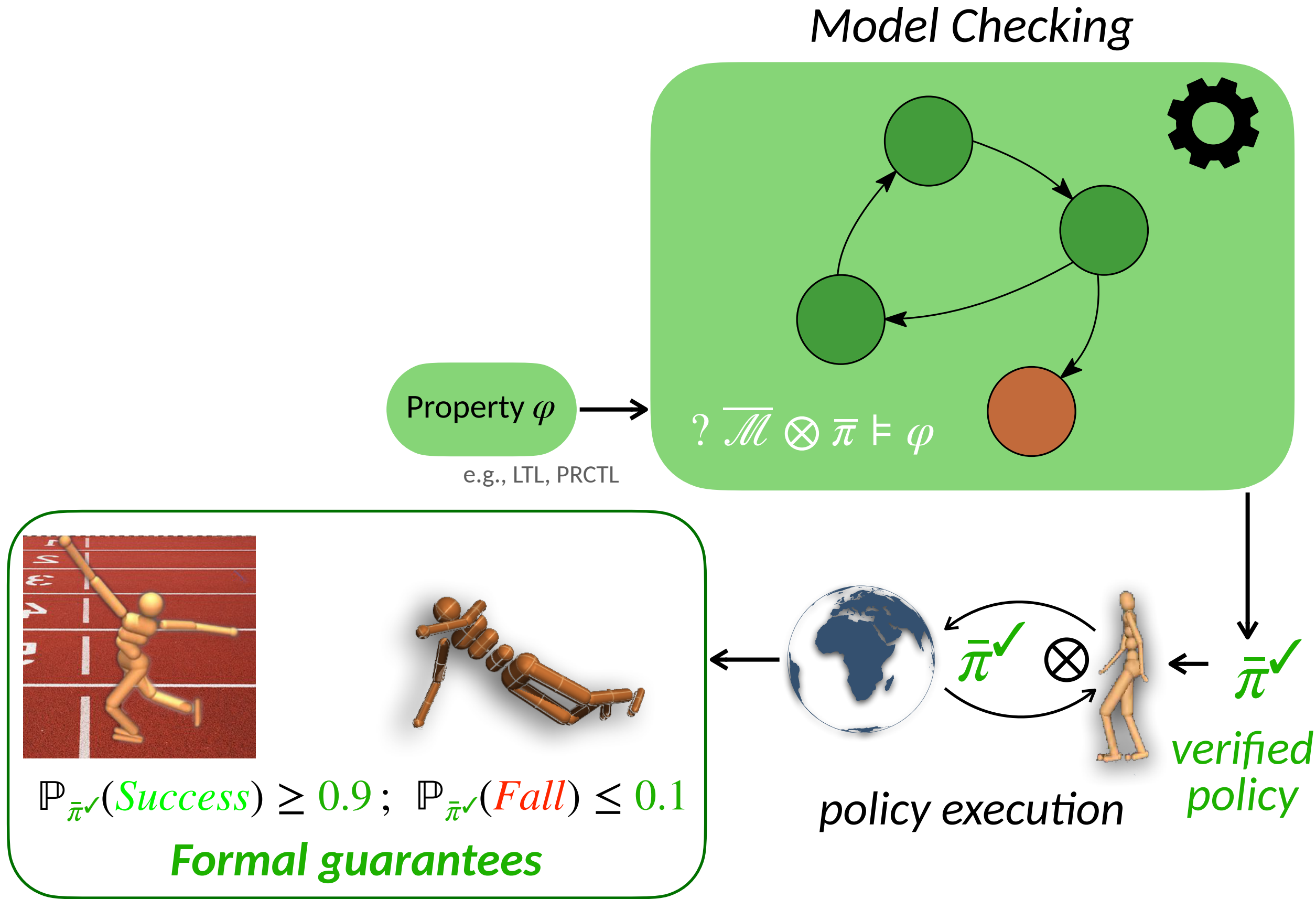
- Full knowledge of the model of the interaction

## Reinforcement Learning



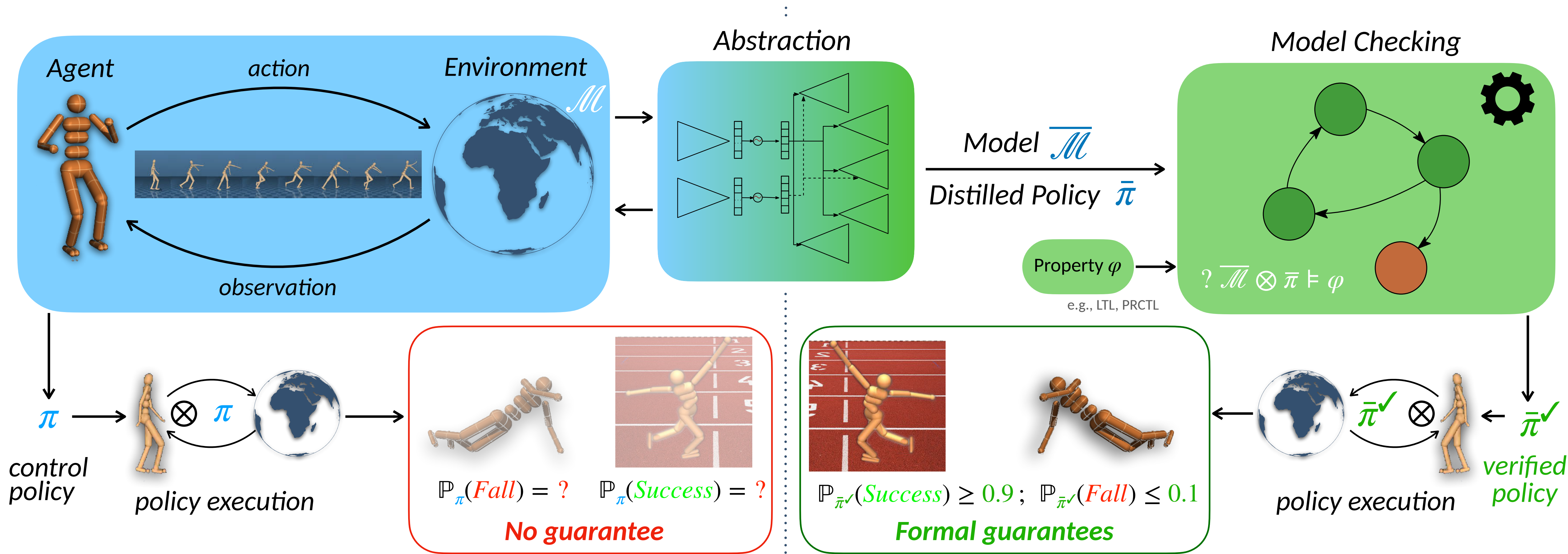
- Continuous state/action spaces
- Unknown environment

## Formal Guarantees



- Full knowledge of the model of the interaction
- Exhaustive exploration of the model
- Sensitive to the state space explosion problem

## Reinforcement Learning Policies with Formal Guarantees



- Continuous state/action spaces
- Unknown environment

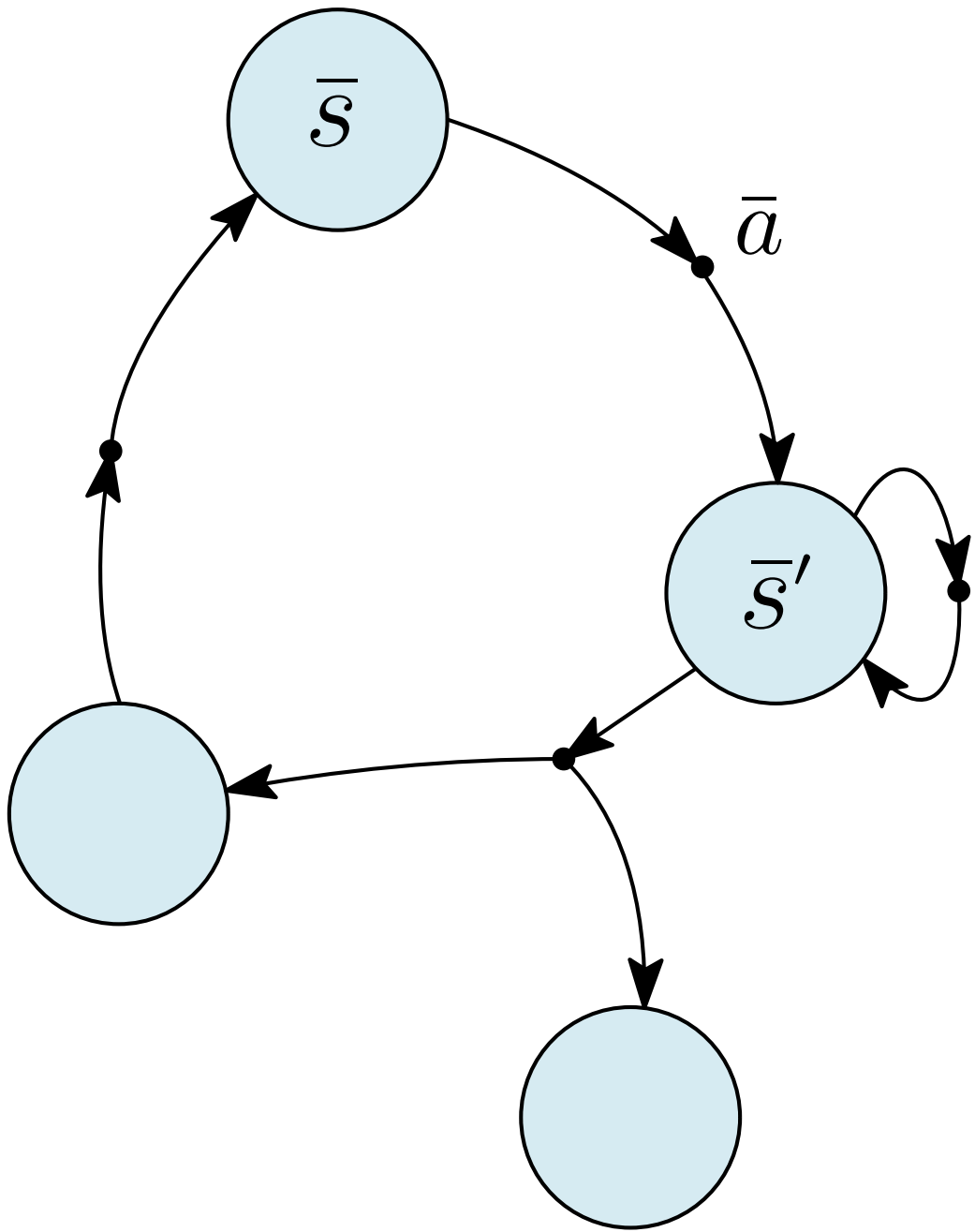
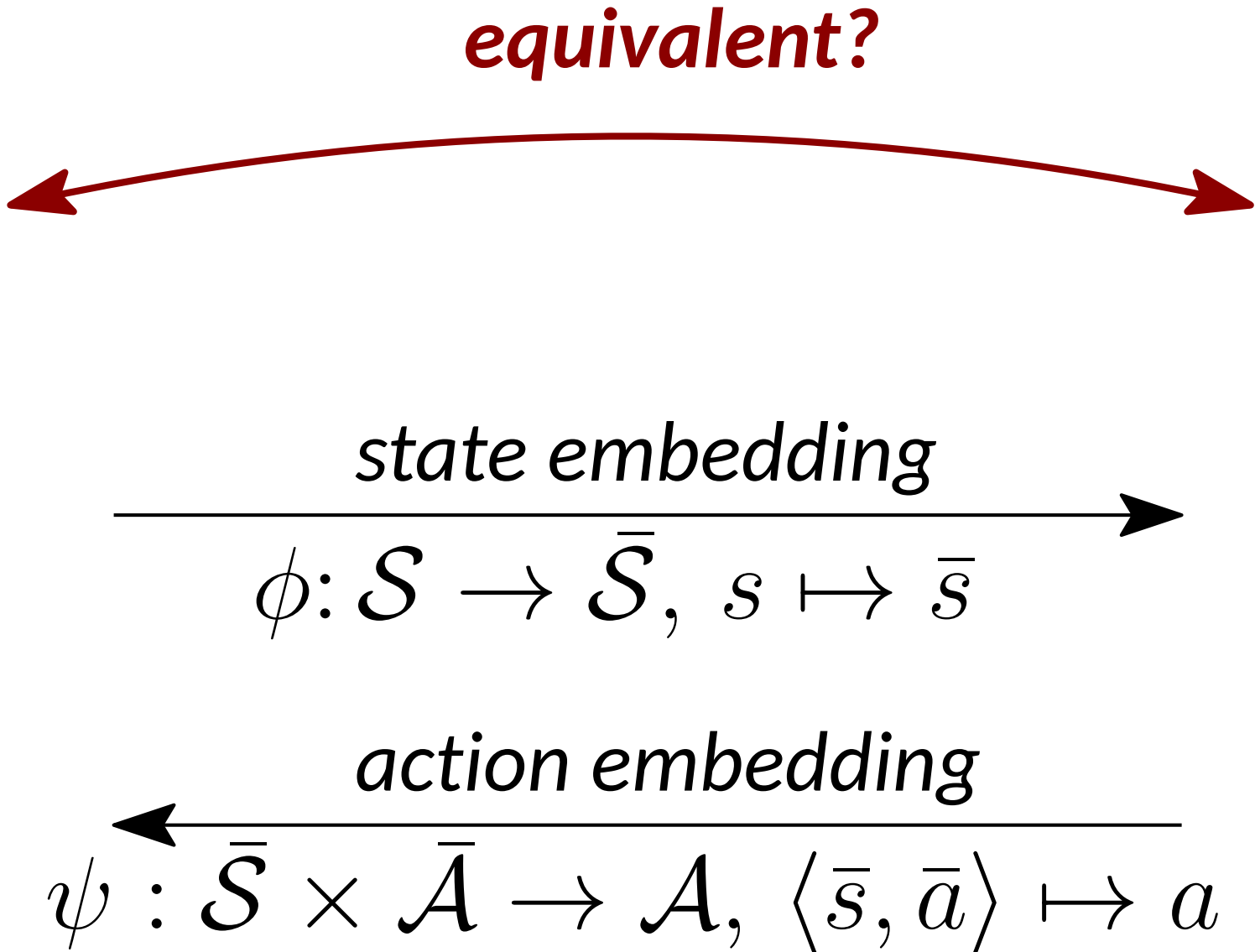
- Full knowledge of the model of the interaction
- Exhaustive exploration of the model
- Sensitive to the state space explosion problem

# Latent Space Model



Continuous-spaces MDP

$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbf{P}, \ell \rangle$$



Discrete latent MDP

$$\bar{\mathcal{M}} = \langle \bar{\mathcal{S}}, \bar{\mathcal{A}}, \bar{\mathcal{R}}, \bar{\mathbf{P}}, \ell \rangle$$

# Latent Space Model

$B \in \mathcal{S}^2$  is a **stochastic bisimulation** iff for all  $s_1, s_2 \in \mathcal{S}, a \in \mathcal{A}, T/B$

$$\ell(s_1) = \ell(s_2) \quad \mathcal{R}(s_1, a) = \mathcal{R}(s_2, a) \quad \text{and} \quad \mathbf{P}(T \mid s_1, a) = \mathbf{P}(T \mid s_2, a)$$

**Largest:**  $\sim$

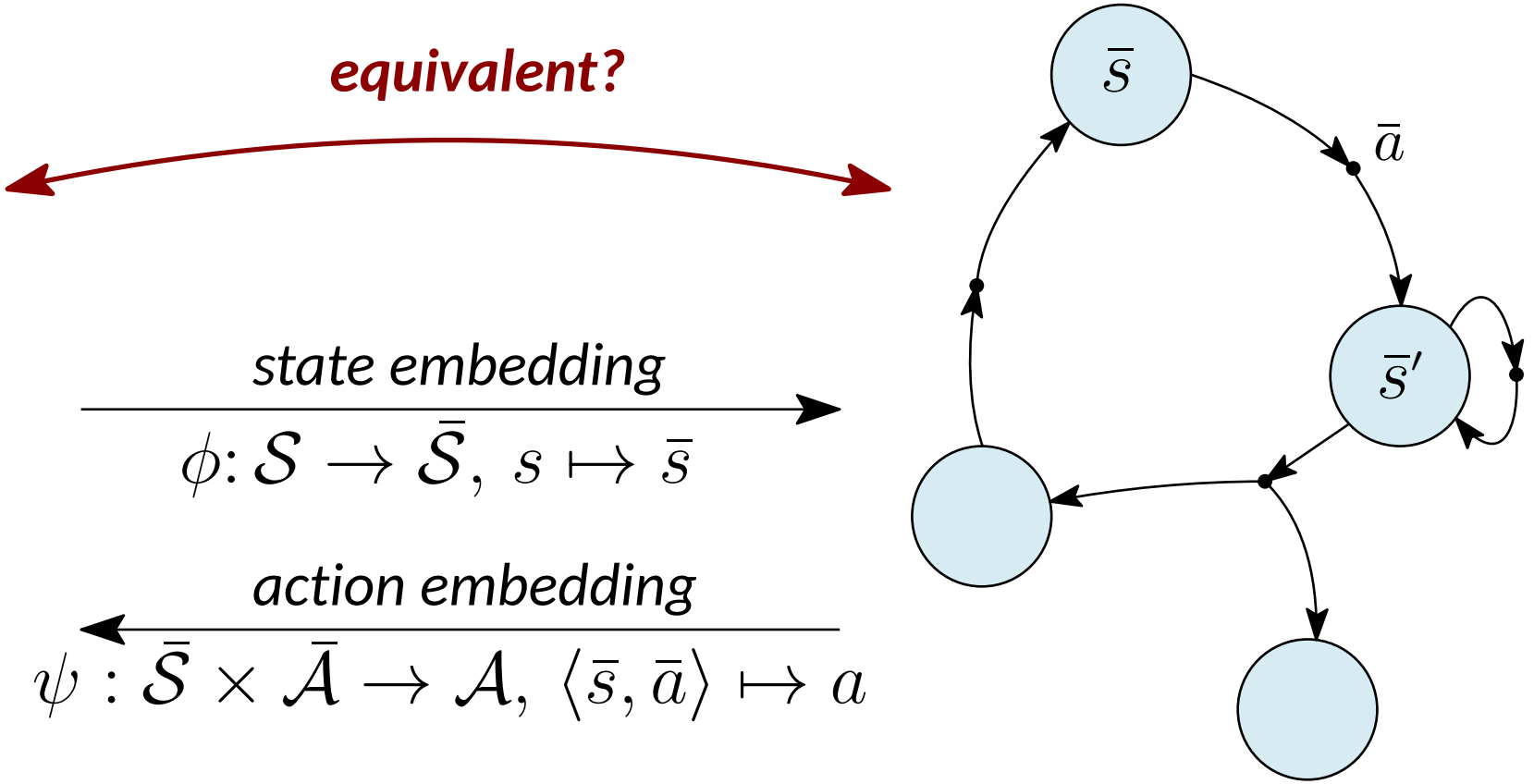
(Larsen and Skou 1989;  
Givan, Dean, and Greig 2003)

- Behavioral equivalence
- ➔ Trajectory, value, and optimal policy equivalence
- ➔ Agents **behave in the same way** in bisimilar models



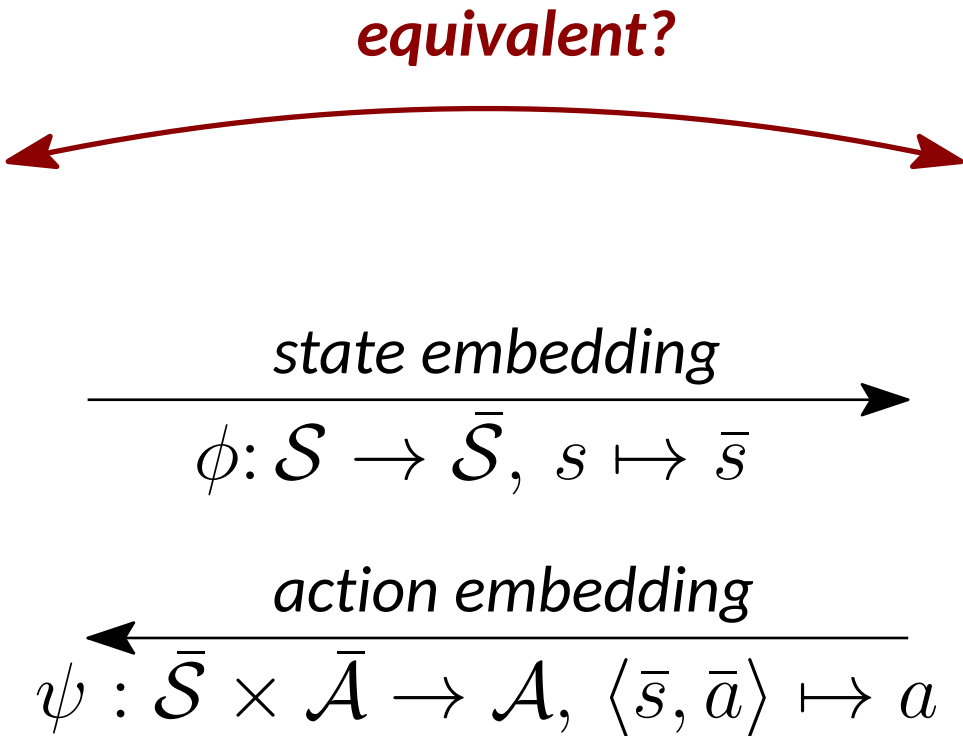
Continuous-spaces MDP

$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbf{P}, \ell \rangle$$



Discrete latent MDP

$$\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$$



# Latent Space Model

$B \in \mathcal{S}^2$  is a ~~stochastic bisimulation~~ iff for all  $s_1, s_2 \in \mathcal{S}, a \in \mathcal{A}, T/B$

$$\ell(s_1) = \ell(s_2) \quad \mathcal{R}(s_1, a) \neq \mathcal{R}(s_2, a) + \epsilon \quad \text{and} \quad \mathbf{P}(T \mid s_1, a) \neq \mathbf{P}(T \mid s_2, a) + \epsilon$$

Largest:  $\sim$

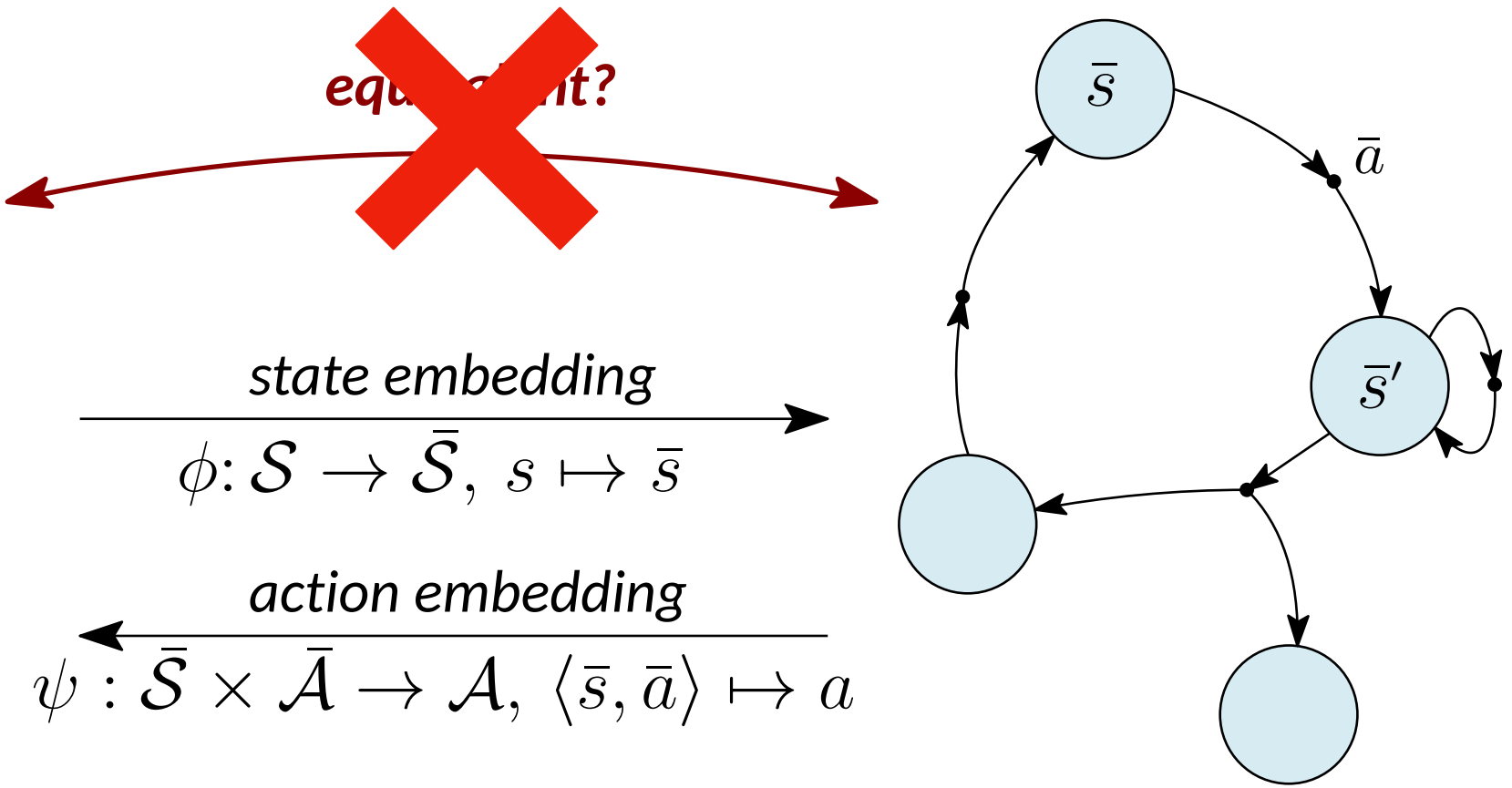
(Larsen and Skou 1989;  
Givan, Dean, and Greig 2003)

- Behavioral equivalence
- ➔ Trajectory, value, and optimal policy equivalence
- ➔ Agents *behave in the same way* in bisimilar models



Continuous-spaces MDP

$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbf{P}, \ell \rangle$$



Discrete latent MDP

$$\bar{\mathcal{M}} = \langle \bar{\mathcal{S}}, \bar{\mathcal{A}}, \bar{\mathcal{R}}, \bar{\mathbf{P}}, \ell \rangle$$

⦿ *All or nothing*: two states *nearly identical* with slight numerical difference  $\epsilon$  are  $\neq$

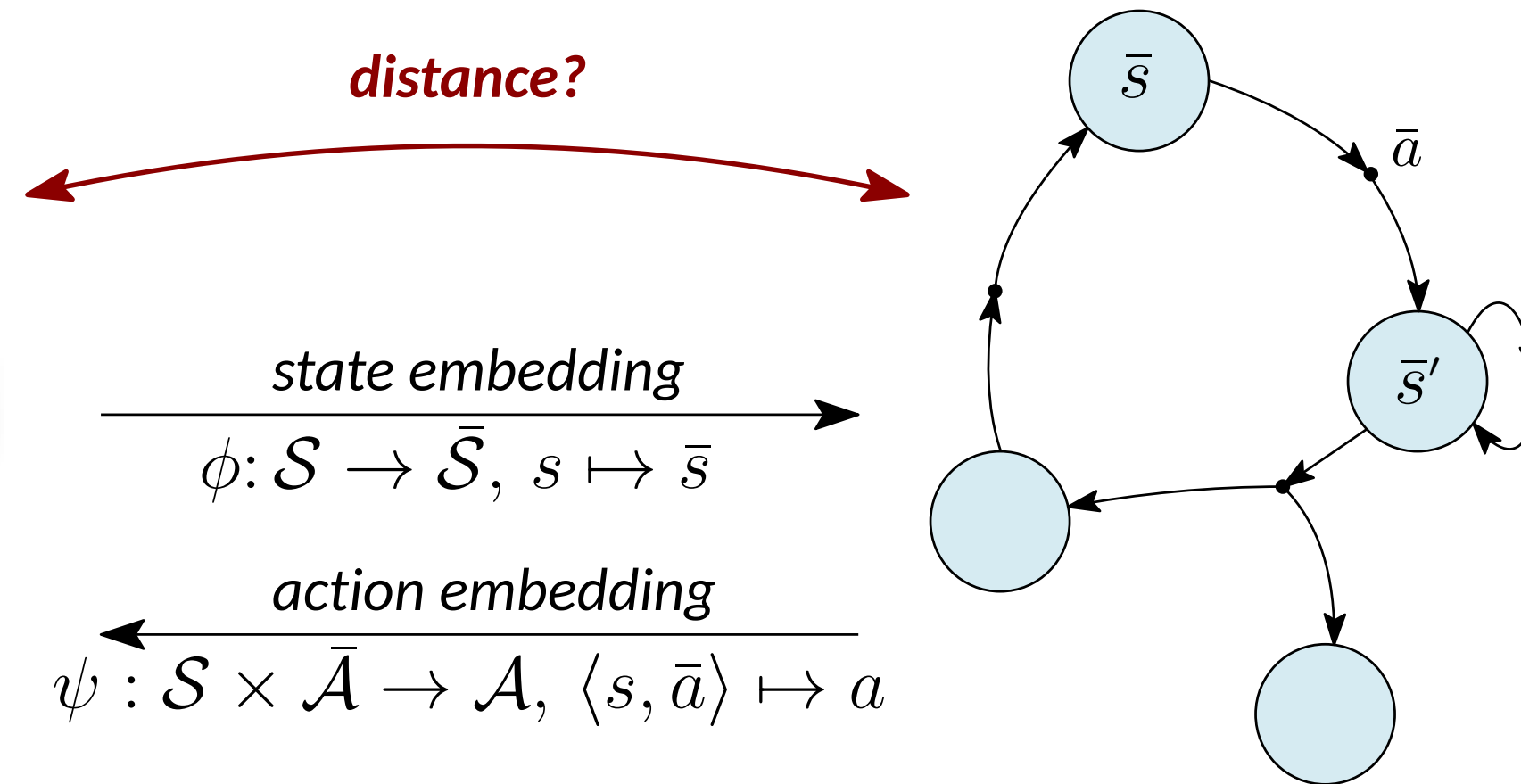
# Bisimulation distance

Continuous-spaces MDP



$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbf{P}, \ell \rangle$$

Discrete latent MDP



$$\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$$

- For policy  $\pi$ ,  $\gamma \in [0, 1[$ , and formal logic  $\mathcal{L}$ :

➔ **Bisimulation distance:** largest behavioral difference (Desharnais et. al, 2004)

$$\tilde{d}_{\pi}(s_1, s_2) = \sup_{V \in \mathcal{F}_{\gamma}^{\mathcal{L}}(\pi)} \left| V_{\pi}(s_1) - V_{\pi}(s_2) \right| \quad \forall s_1, s_2 \in \mathcal{S}$$

where  $\mathcal{F}_{\gamma}^{\mathcal{L}}(\pi)$  is a logical family of functional expressions defining the semantics of  $\mathcal{L}$

➔ **Kernel is bisimilarity:**  $\tilde{d}_{\pi}(s_1, s_2) = 0 \iff s_1 \sim s_2$

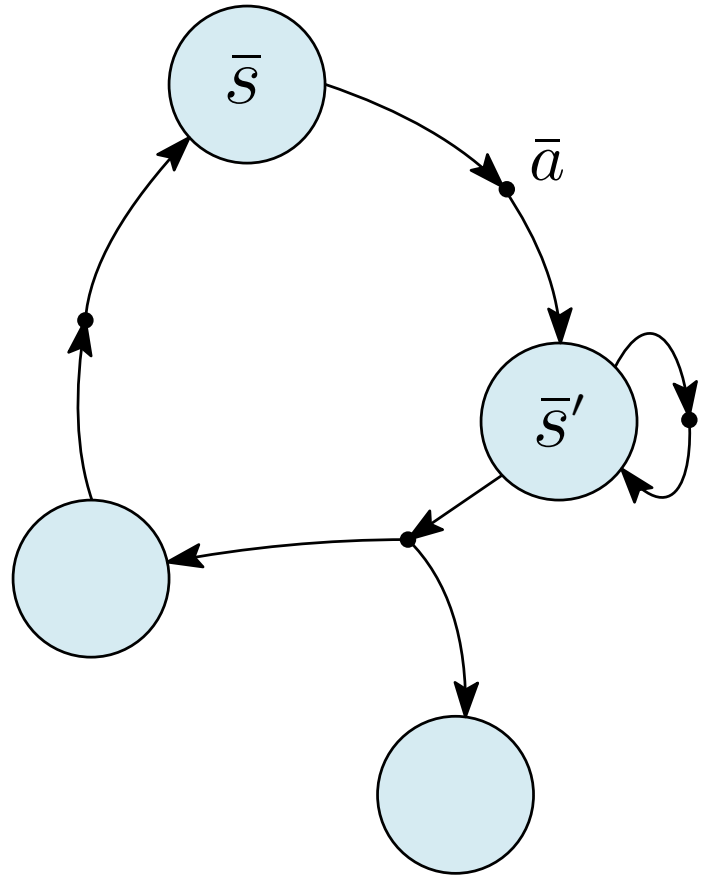
# Bisimulation distance

Continuous-spaces MDP

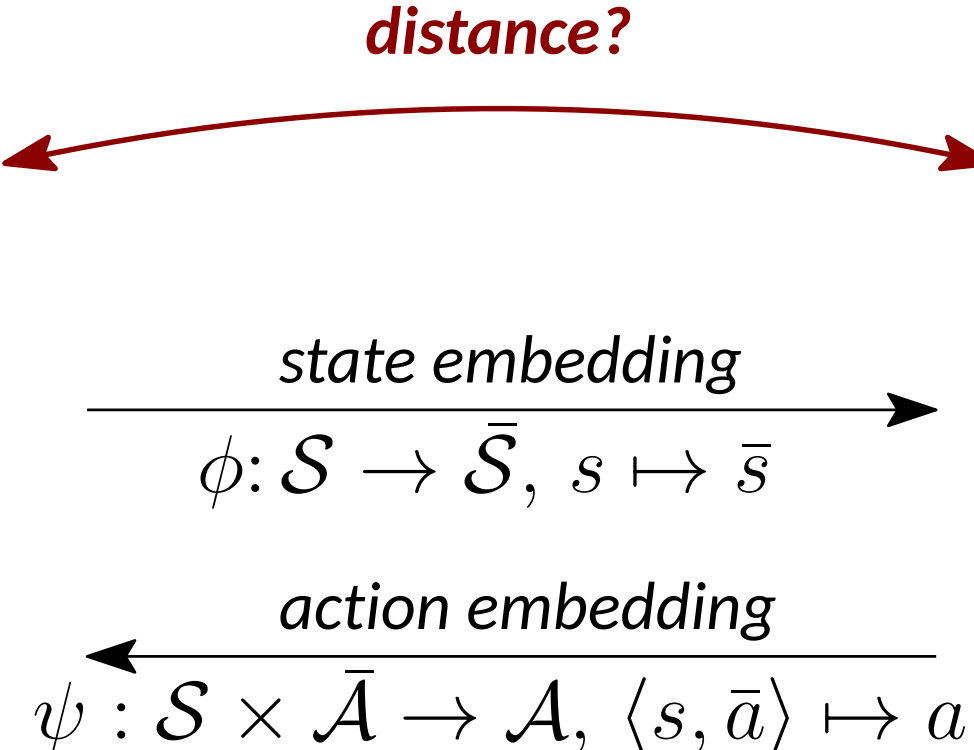


$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbf{P}, \ell \rangle$$

Discrete latent MDP



$$\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$$



- For policy  $\pi$ ,  $\gamma \in [0, 1[$ , and formal logic  $\mathcal{L}$ :

➔ **Bisimulation distance:** largest behavioral difference (Desharnais et. al, 2004)

$$\tilde{d}_\pi(s_1, s_2) = \sup_{V \in \mathcal{F}_\gamma^\mathcal{L}(\pi)} |V_\pi(s_1) - V_\pi(s_2)| \quad \forall s_1, s_2 \in \mathcal{S}$$

We need a policy that can be executed (separately) in the original and latent models

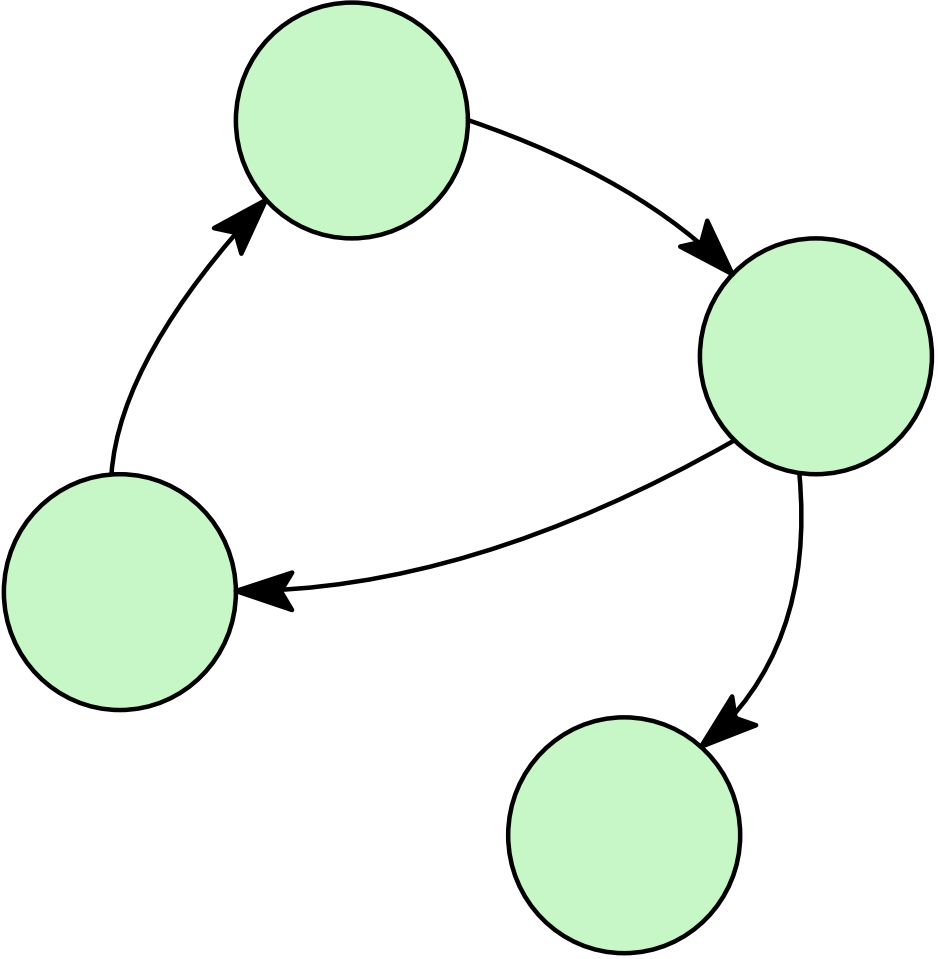
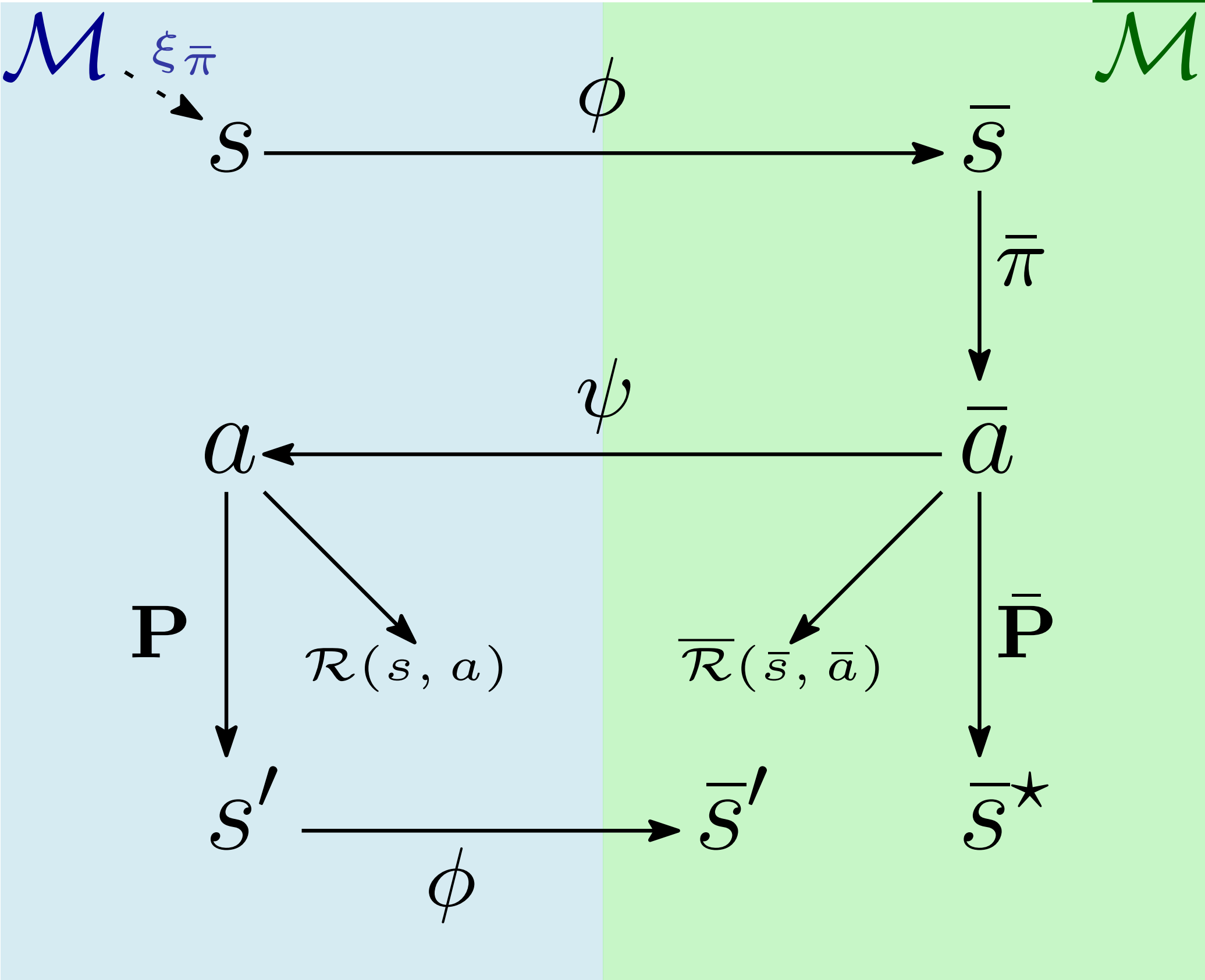
where  $\mathcal{F}_\gamma^\mathcal{L}(\pi)$  is a logical family of functional expressions defining the semantics of  $\mathcal{L}$

➔ **Kernel is bisimilarity:**  $\tilde{d}_\pi(s_1, s_2) = 0 \iff s_1 \sim s_2$



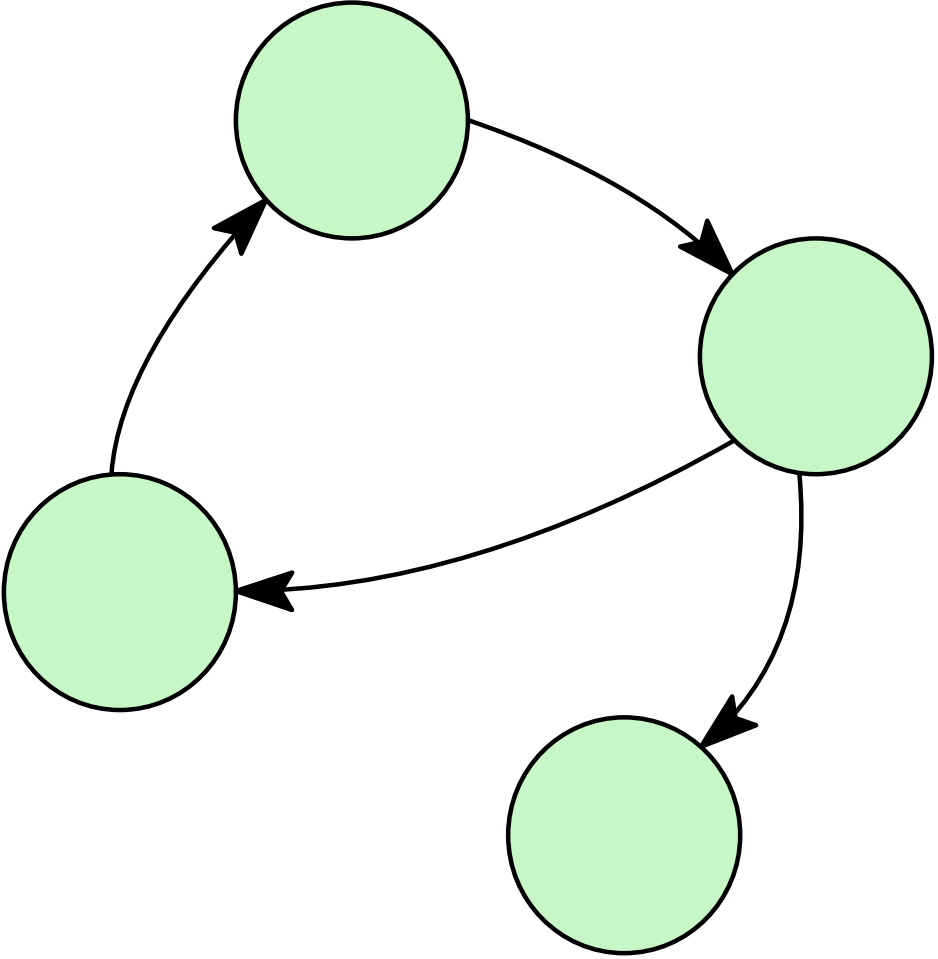
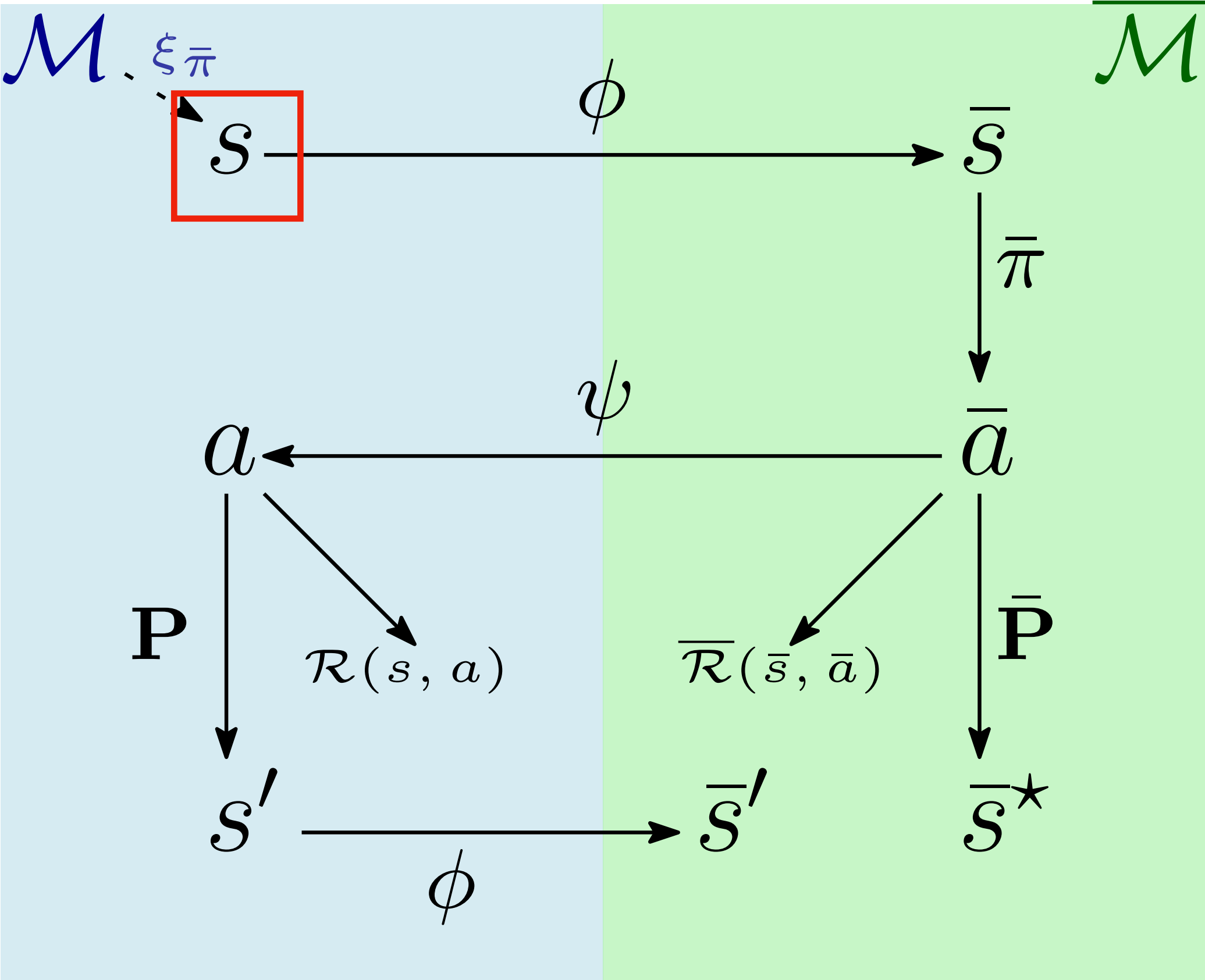
# Latent Flow

Execution of a latent policy  $\bar{\pi}$  in the original model



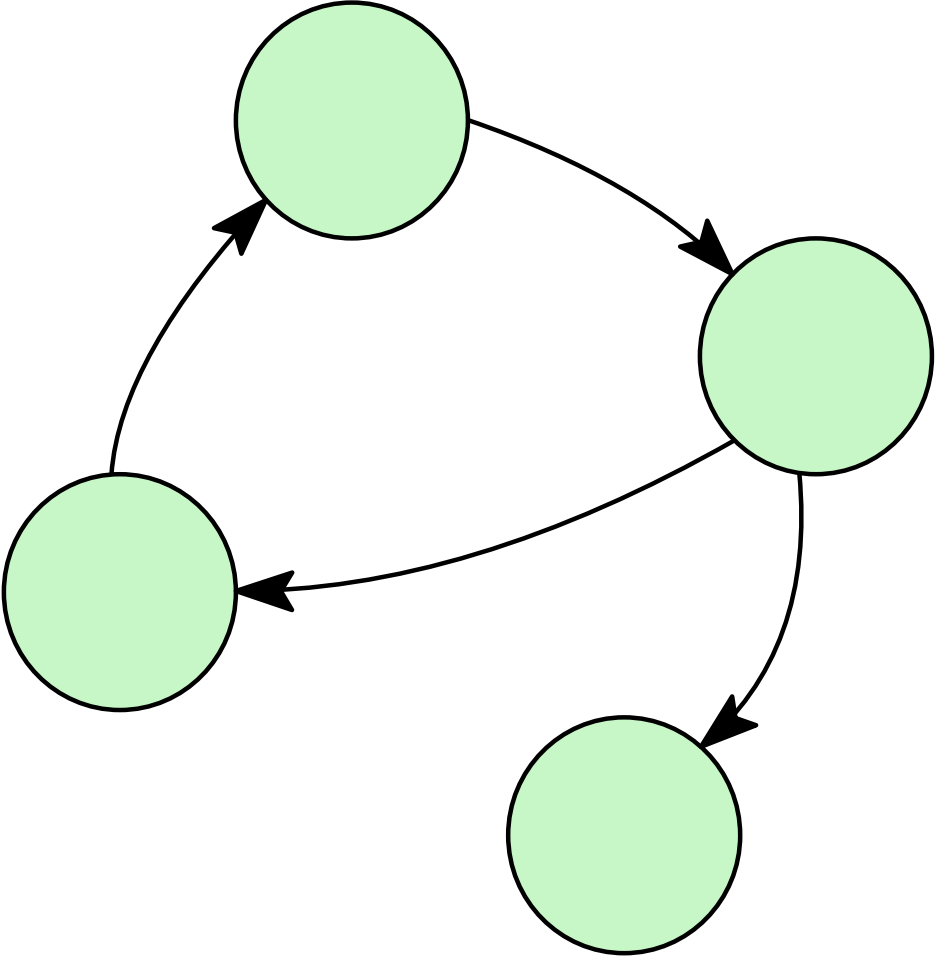
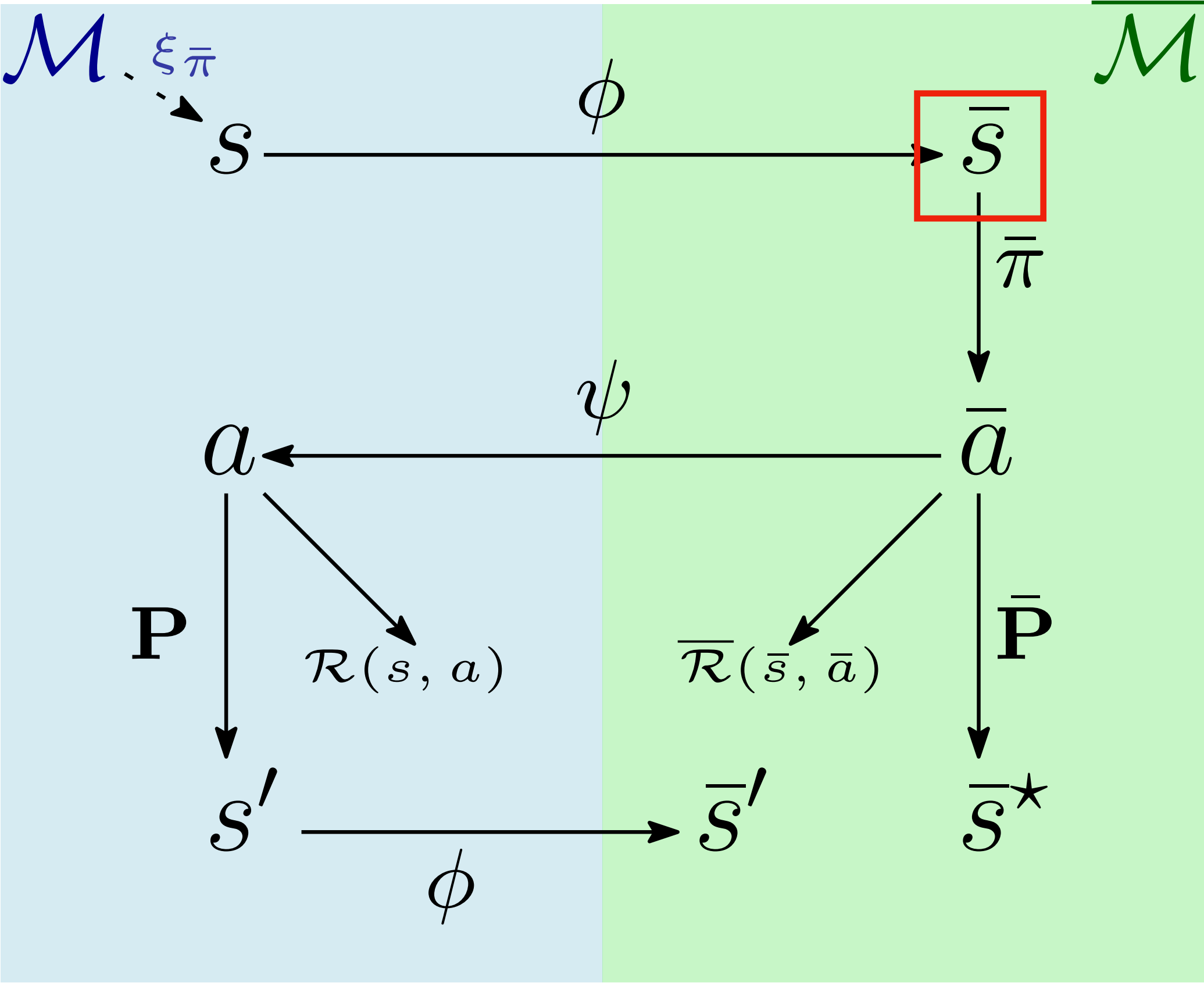
# Latent Flow

Execution of a latent policy  $\bar{\pi}$  in the original model



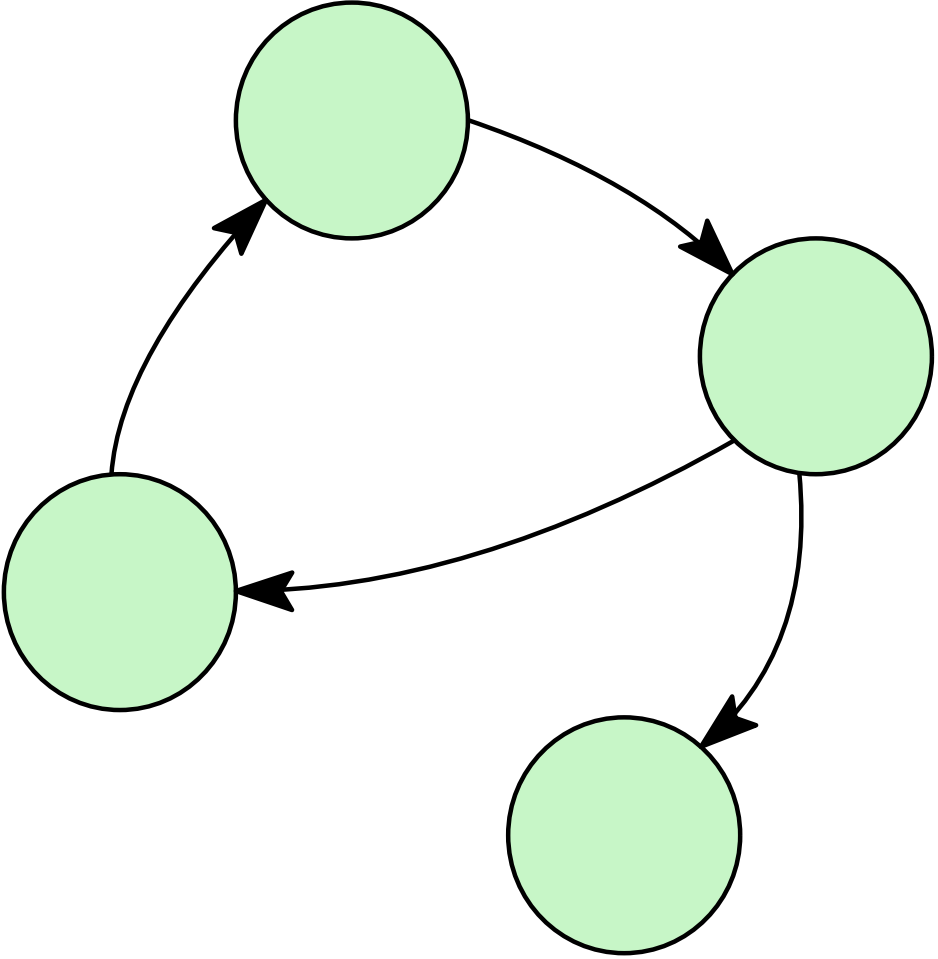
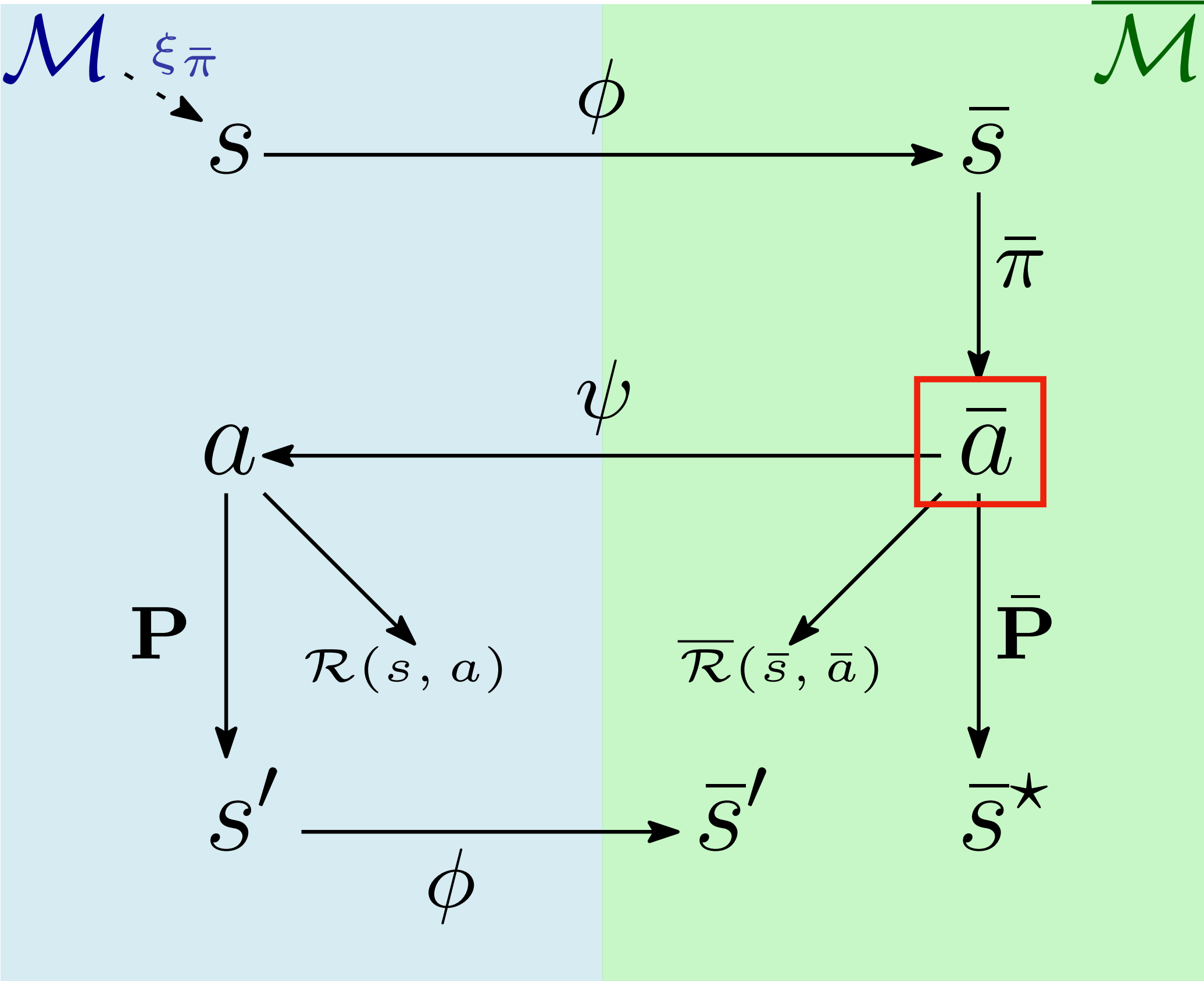
# Latent Flow

Execution of a latent policy  $\bar{\pi}$  in the original model



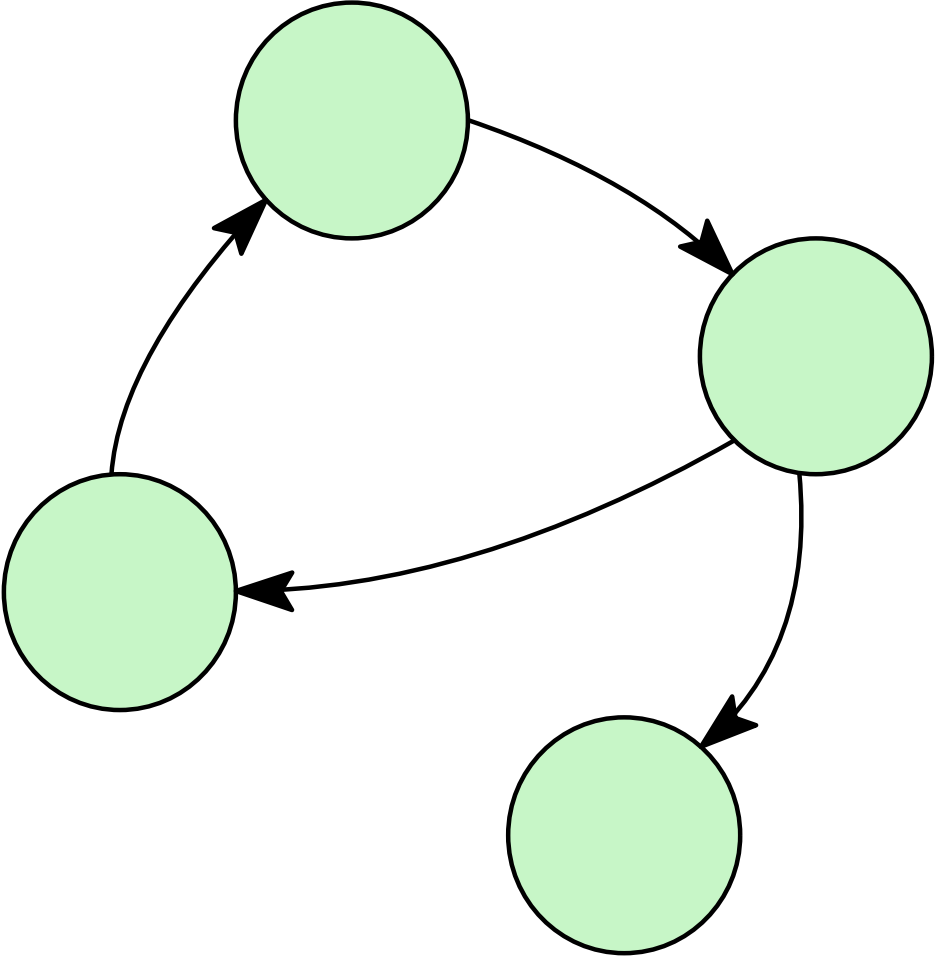
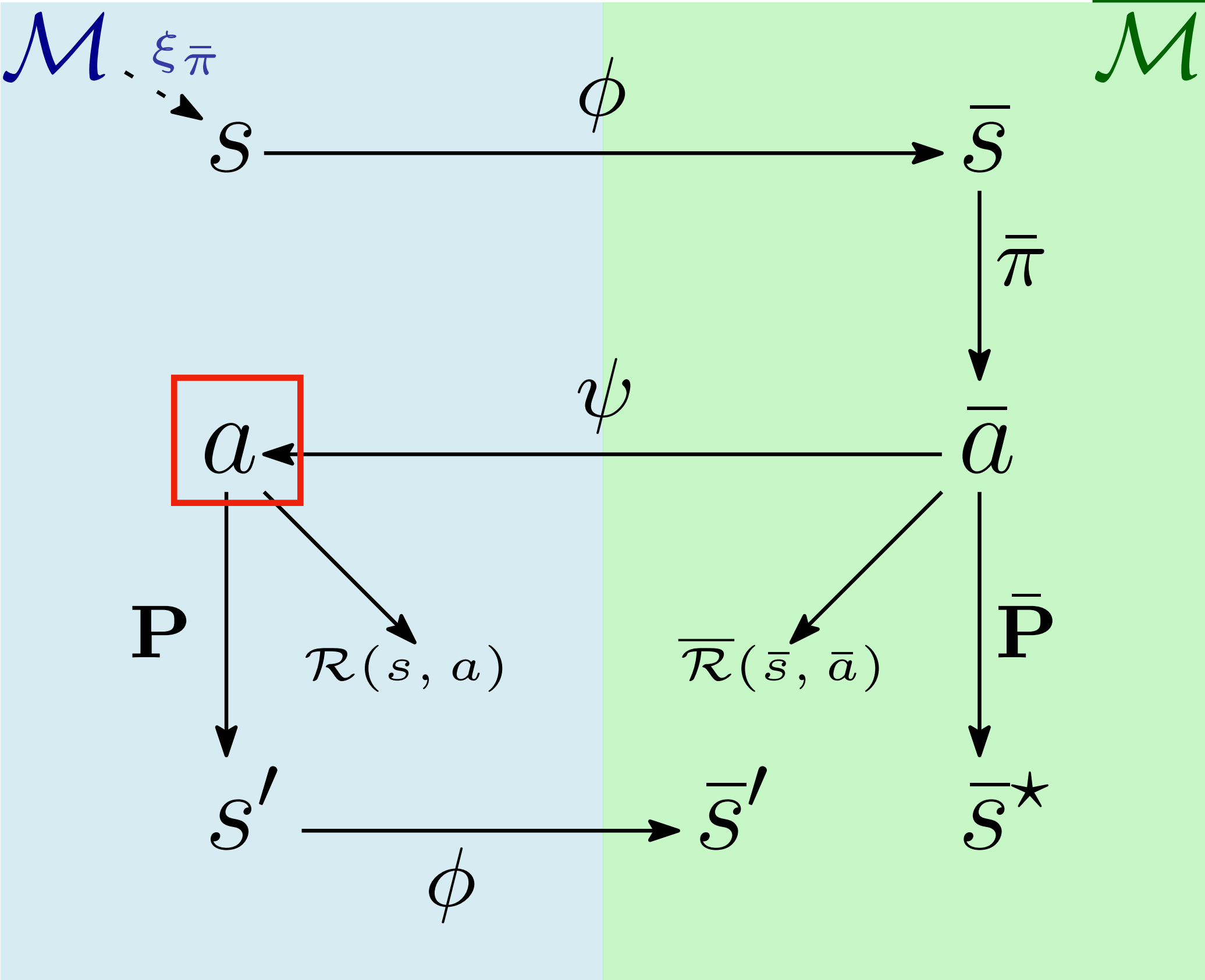
# Latent Flow

Execution of a latent policy  $\bar{\pi}$  in the original model



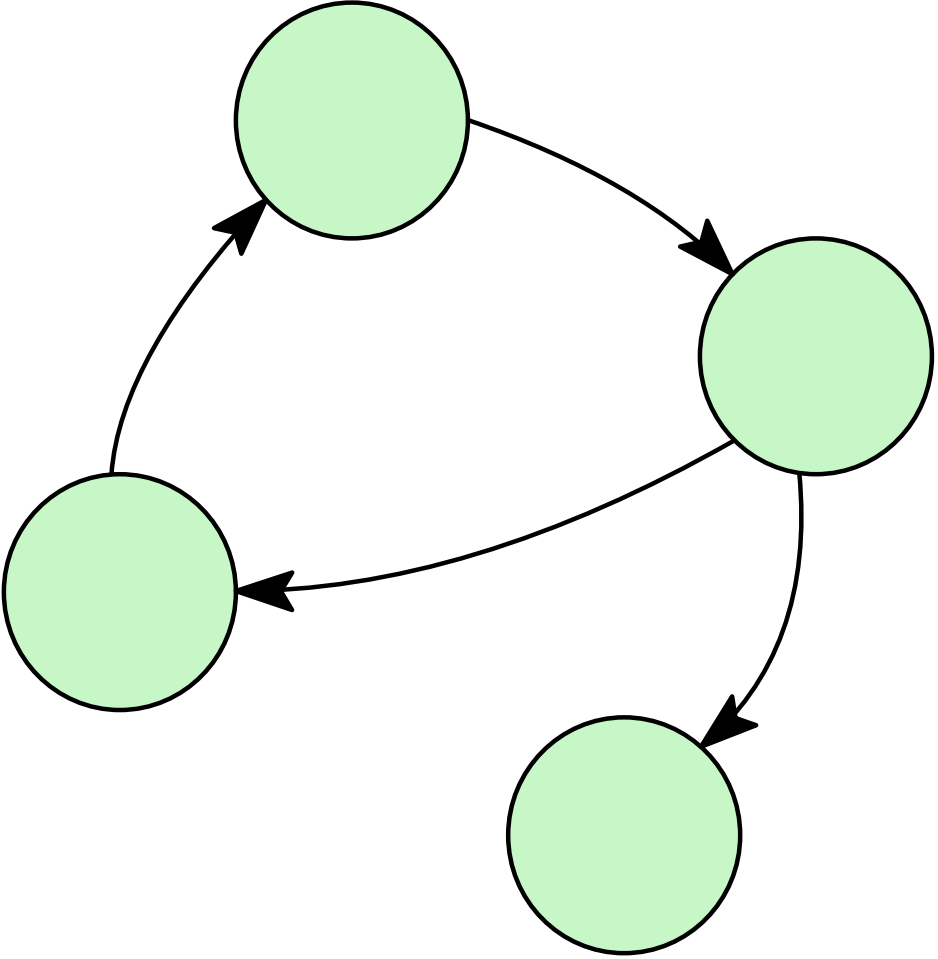
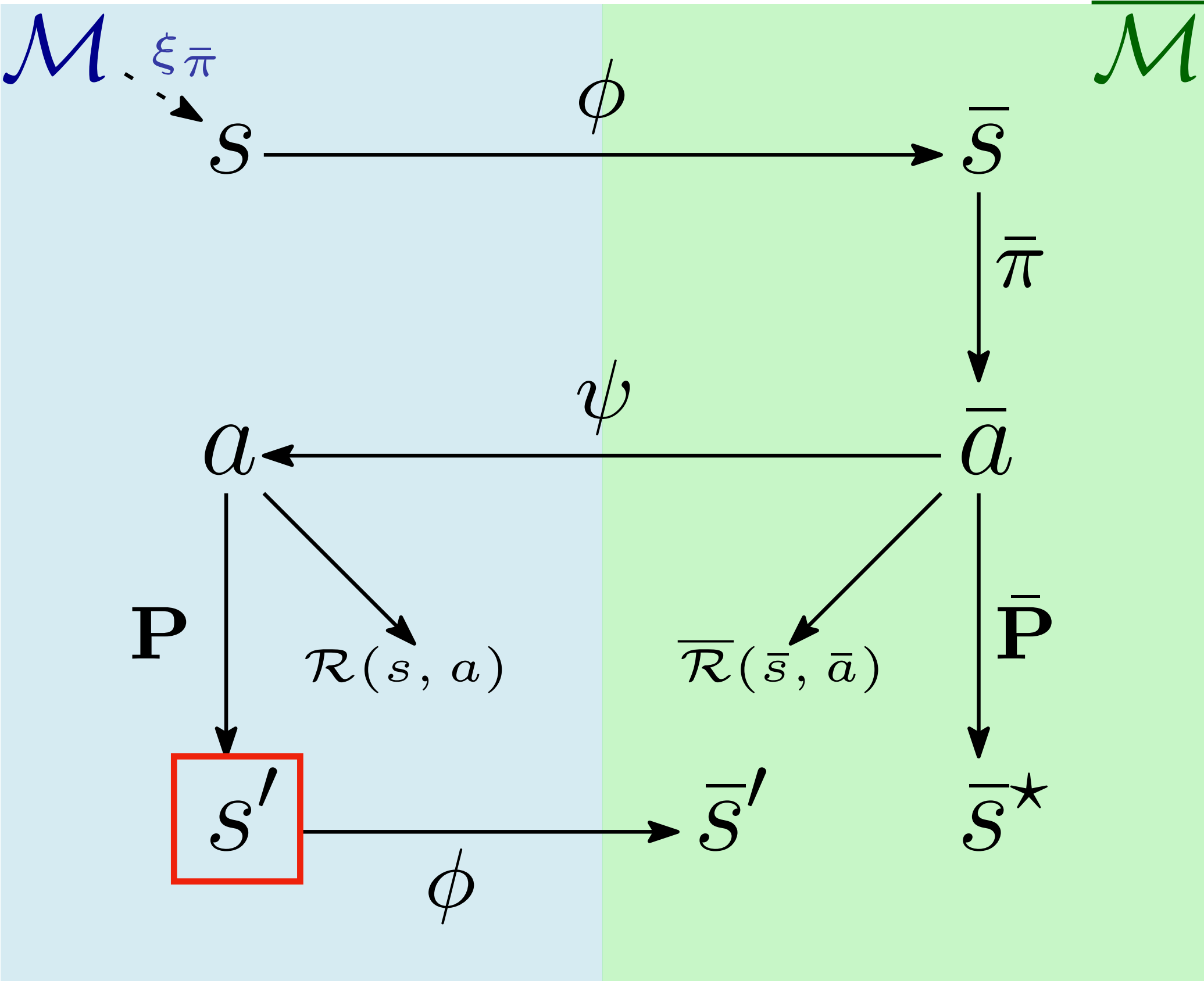
# Latent Flow

Execution of a latent policy  $\bar{\pi}$  in the original model



# Latent Flow

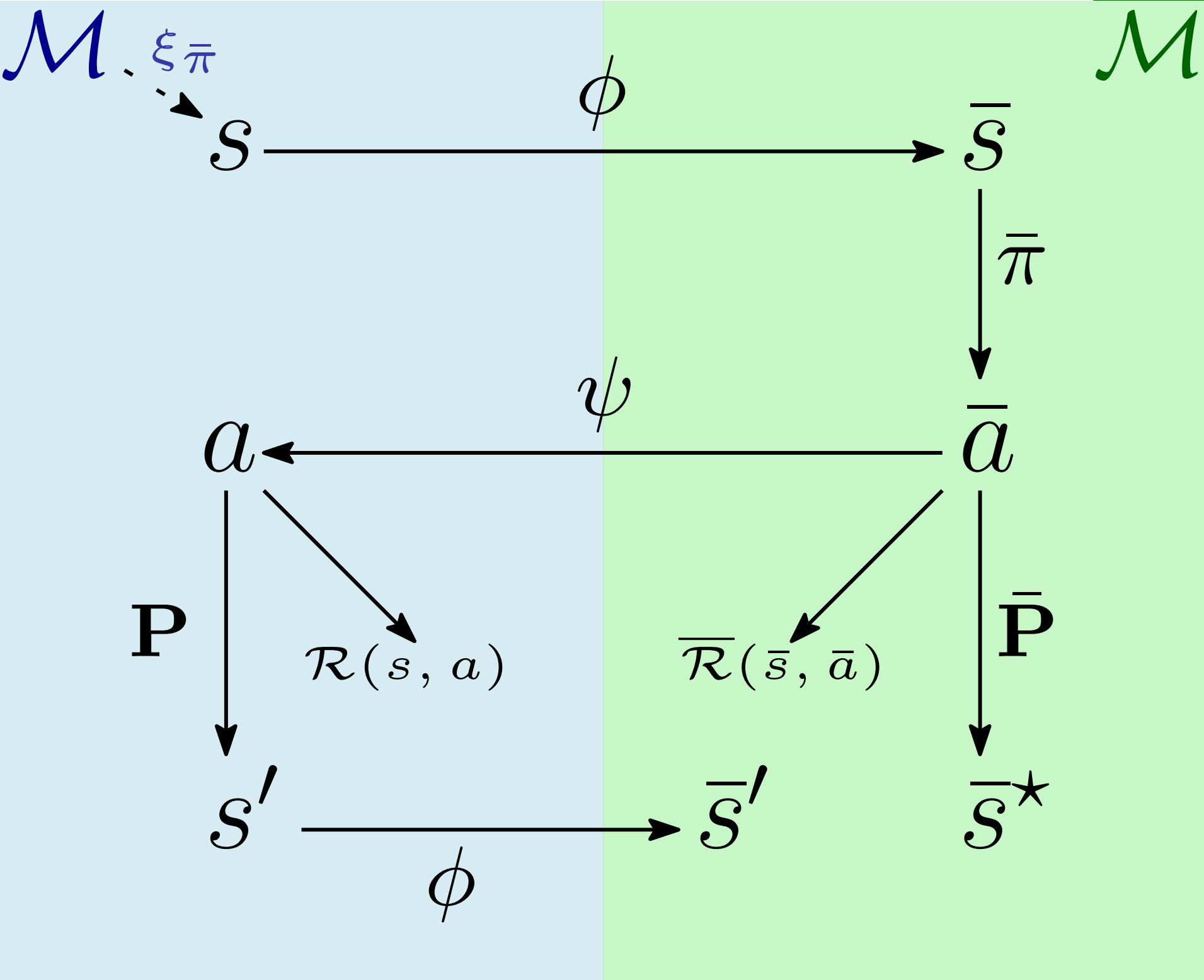
Execution of a latent policy  $\bar{\pi}$  in the original model



# Latent Flow

## Execution of a latent policy $\bar{\pi}$ in the original model: **Local Losses**

- Latent policy  $\bar{\pi}$ , stationary distribution  $\xi_{\bar{\pi}}$

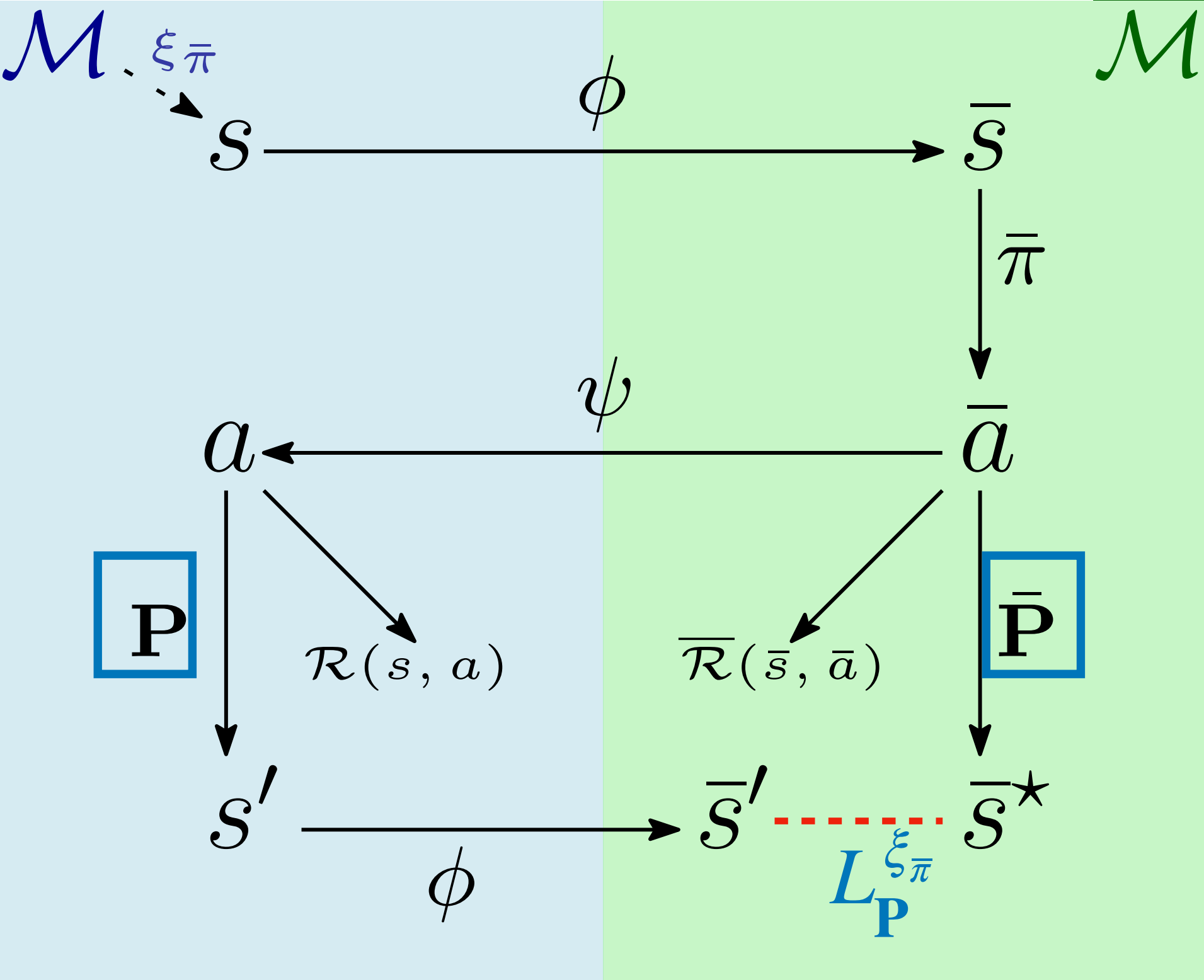


# Latent Flow

## Execution of a latent policy $\bar{\pi}$ in the original model: **Local Losses**

- Latent policy  $\bar{\pi}$ , stationary distribution  $\xi_{\bar{\pi}}$

$$L_{\mathbf{P}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} W_{d_{\bar{S}}}(\phi \mathbf{P}(\cdot | s, \bar{a}), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a}))$$





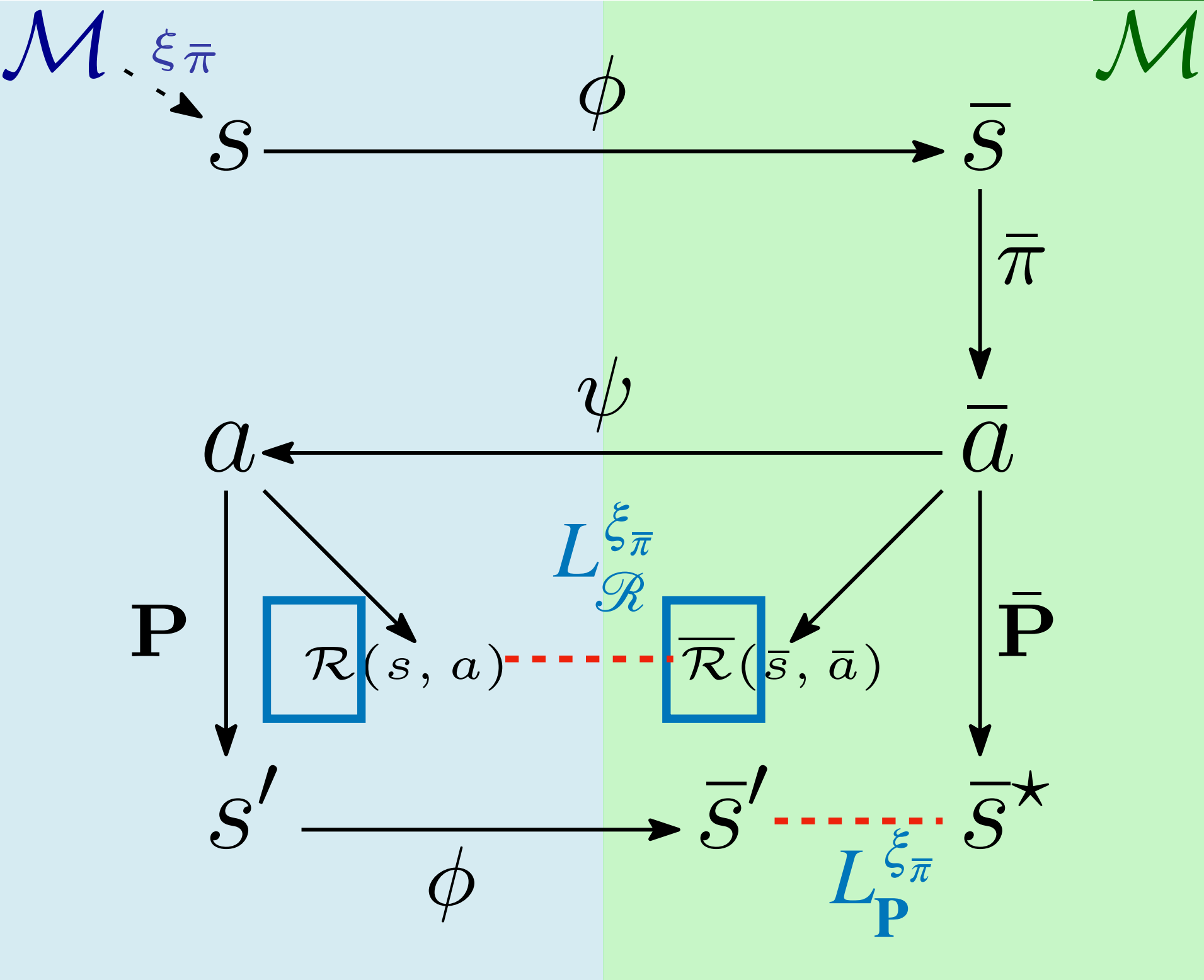
# Latent Flow

## Execution of a latent policy $\bar{\pi}$ in the original model: **Local Losses**

- Latent policy  $\bar{\pi}$ , stationary distribution  $\xi_{\bar{\pi}}$

$$L_{\mathbf{P}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} W_{d_{\bar{S}}}(\phi \mathbf{P}(\cdot | s, \bar{a}), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a}))$$

$$L_{\mathcal{R}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} \left| \mathcal{R}(s, \bar{a}) - \bar{\mathcal{R}}(\phi(s), \bar{a}) \right|$$



# Latent Flow

## Execution of a latent policy $\bar{\pi}$ in the original model: **Local Losses**

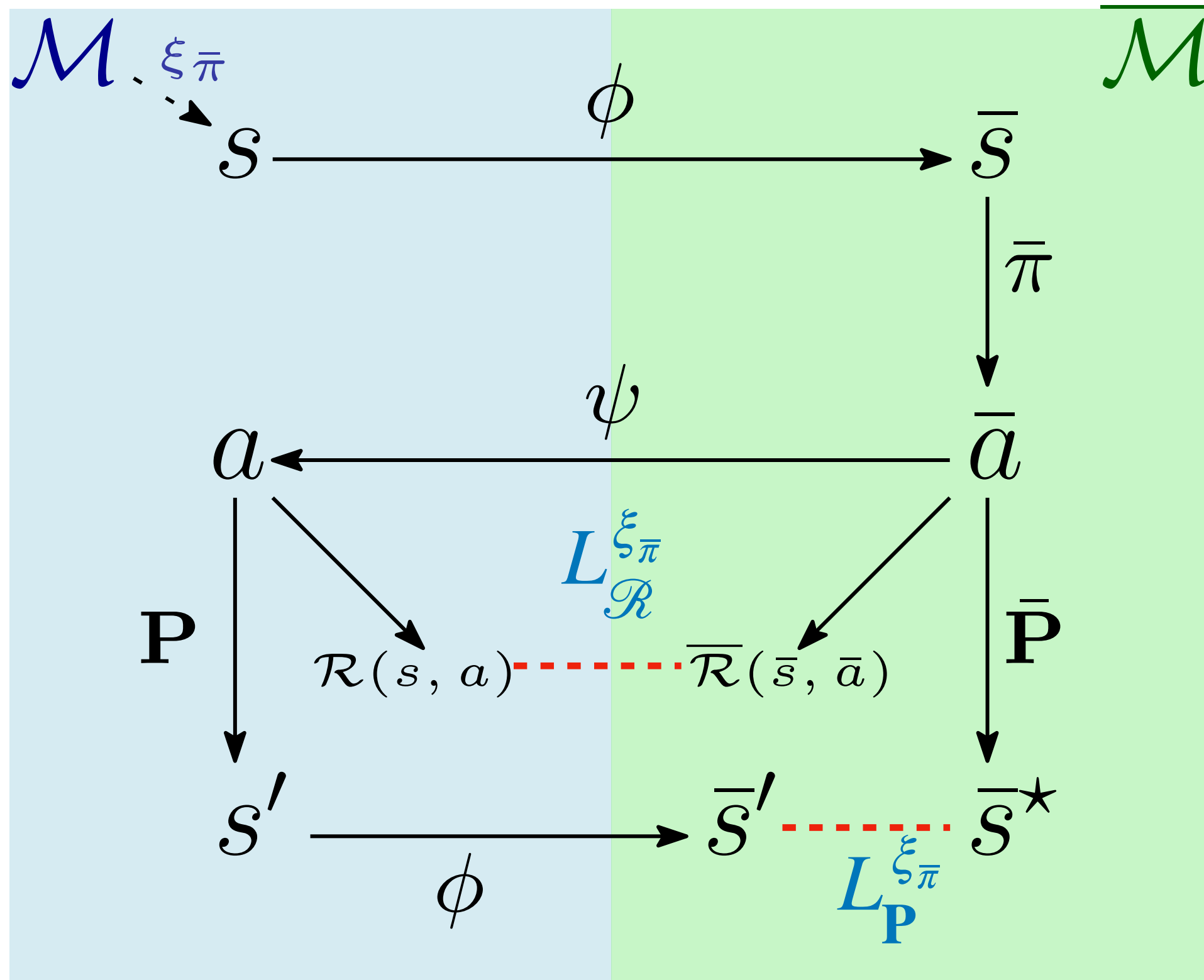
- Latent policy  $\bar{\pi}$ , stationary distribution  $\xi_{\bar{\pi}}$

$$L_{\mathbf{P}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} W_{d_{\bar{S}}} (\phi \mathbf{P} (\cdot | s, \bar{a}), \bar{\mathbf{P}} (\cdot | \phi(s), \bar{a}))$$

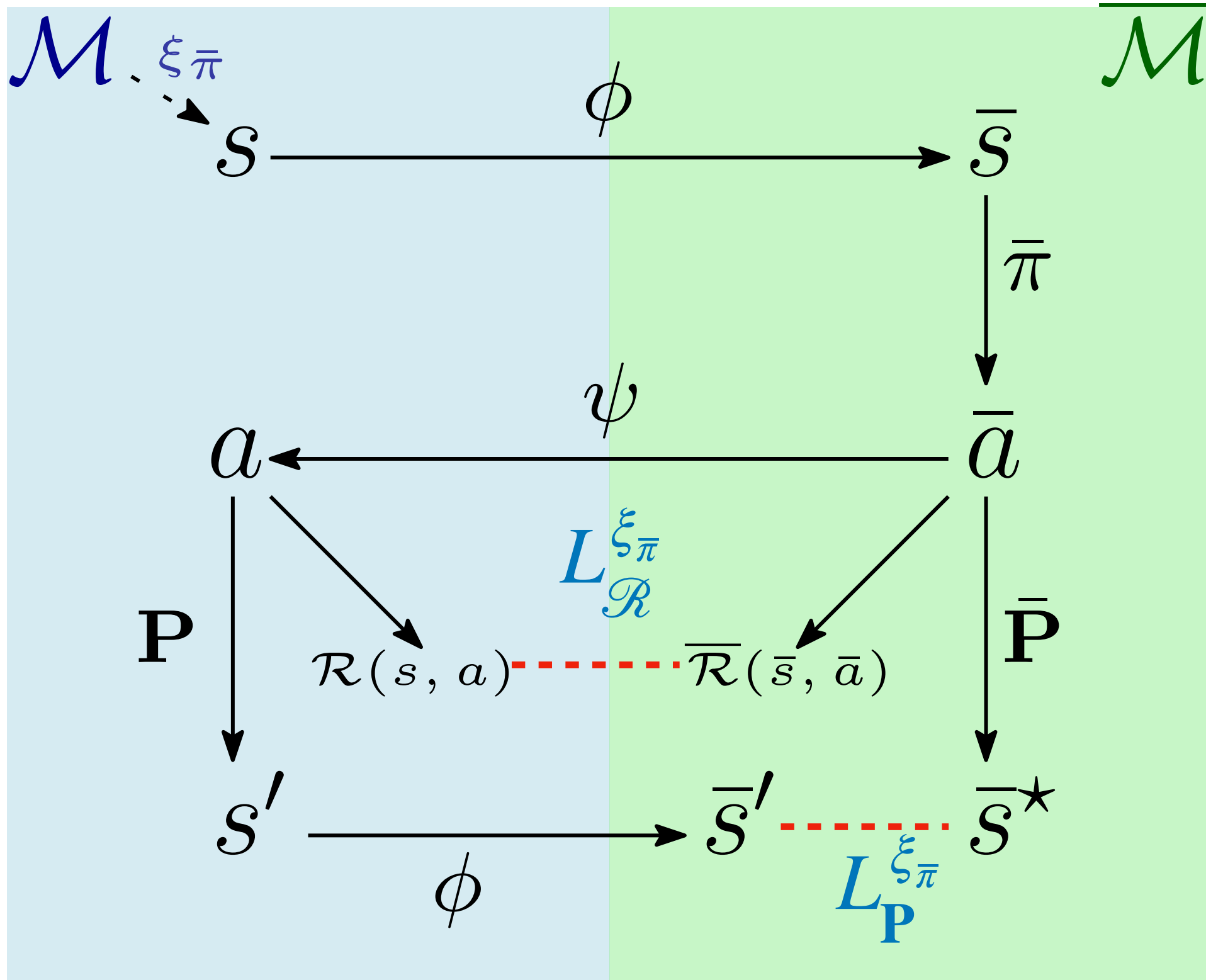
$$L_{\mathcal{R}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} \left| \mathcal{R}(s, \bar{a}) - \bar{\mathcal{R}}(\phi(s), \bar{a}) \right|$$

- Abstraction quality:**  $\mathbb{E}_{s \sim \xi_{\bar{\pi}}} \tilde{d}_{\bar{\pi}}(s, \phi(s)) \leq \frac{L_{\mathcal{R}}^{\xi_{\bar{\pi}}} + \gamma L_{\mathbf{P}}^{\xi_{\bar{\pi}}}}{1 - \gamma}$
- Representation quality:** for all  $s_1, s_2 \in \mathcal{S}$  such that  $\phi(s_1) = \phi(s_2)$

$$\tilde{d}_{\bar{\pi}}(s_1, s_2) \leq \left( \frac{L_{\mathcal{R}}^{\xi_{\bar{\pi}}} + \gamma L_{\mathbf{P}}^{\xi_{\bar{\pi}}}}{1 - \gamma} \right) \cdot (\xi_{\bar{\pi}}^{-1}(s_1) + \xi_{\bar{\pi}}^{-1}(s_2))$$



## Execution of a latent policy $\bar{\pi}$ in the original model: **Local Losses**



- Latent policy  $\bar{\pi}$ , stationary distribution  $\xi_{\bar{\pi}}$

$$L_{\mathbf{P}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} W_{d_{\bar{S}}}(\phi \mathbf{P}(\cdot | s, \bar{a}), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a}))$$

$$L_{\mathcal{R}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} \left| \mathcal{R}(s, \bar{a}) - \bar{\mathcal{R}}(\phi(s), \bar{a}) \right|$$

- **Abstraction quality:**  $\mathbb{E}_{s \sim \xi_{\bar{\pi}}} \tilde{d}_{\bar{\pi}}(s, \phi(s)) \leq \frac{L_{\mathcal{R}}^{\xi_{\bar{\pi}}} + \gamma L_{\mathbf{P}}^{\xi_{\bar{\pi}}}}{1 - \gamma}$
- **Representation quality:** for all  $s_1, s_2 \in \mathcal{S}$  such that  $\phi(s_1) = \phi(s_2)$

$$\tilde{d}_{\bar{\pi}}(s_1, s_2) \leq \left( \frac{L_{\mathcal{R}}^{\xi_{\bar{\pi}}} + \gamma L_{\mathbf{P}}^{\xi_{\bar{\pi}}}}{1 - \gamma} \right) \cdot (\xi_{\bar{\pi}}^{-1}(s_1) + \xi_{\bar{\pi}}^{-1}(s_2))$$

- **PAC scheme from samples:** let trace  $\langle s_{0:T}, \bar{a}_{0:T-1}, r_{0:T-1} \rangle \sim \xi_{\bar{\pi}}$ ,  $\epsilon, \delta \in ]0, 1[$  and

$$T \geq \left\lceil \frac{-\log(\delta/4)}{2\epsilon^2} \right\rceil:$$

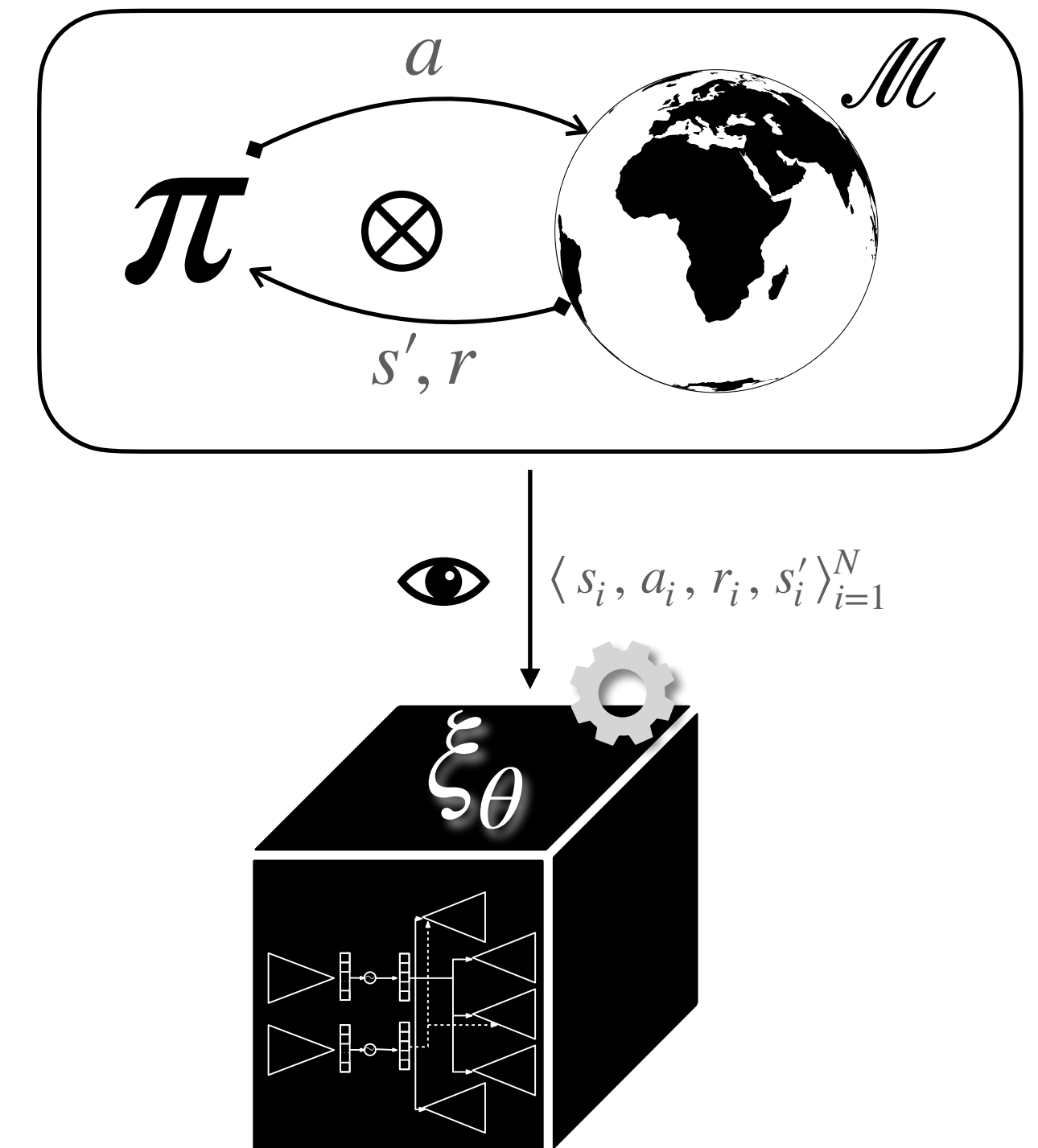
$$\hat{L}_{\mathcal{R}}^{\xi_{\bar{\pi}}} = \frac{1}{T} \sum_{t=0}^{T-1} \left| r_t - \bar{\mathcal{R}}(\phi(s_t), \bar{a}_t) \right| \quad \text{and} \quad \hat{L}_{\mathbf{P}}^{\xi_{\bar{\pi}}} = \frac{1}{T} \sum_{t=0}^{T-1} \left[ 1 - \bar{\mathbf{P}}(\phi(s_{t+1}) | \phi(s_t), \bar{a}_t) \right]$$

Then,  $\left| L_{\mathcal{R}}^{\xi_{\bar{\pi}}} - \hat{L}_{\mathcal{R}}^{\xi_{\bar{\pi}}} \right| \leq \epsilon$  and  $\left| L_{\mathbf{P}}^{\xi_{\bar{\pi}}} - \hat{L}_{\mathbf{P}}^{\xi_{\bar{\pi}}} \right| \leq \epsilon$  with probability  $1 - \delta$

# Learning the Latent Space Model

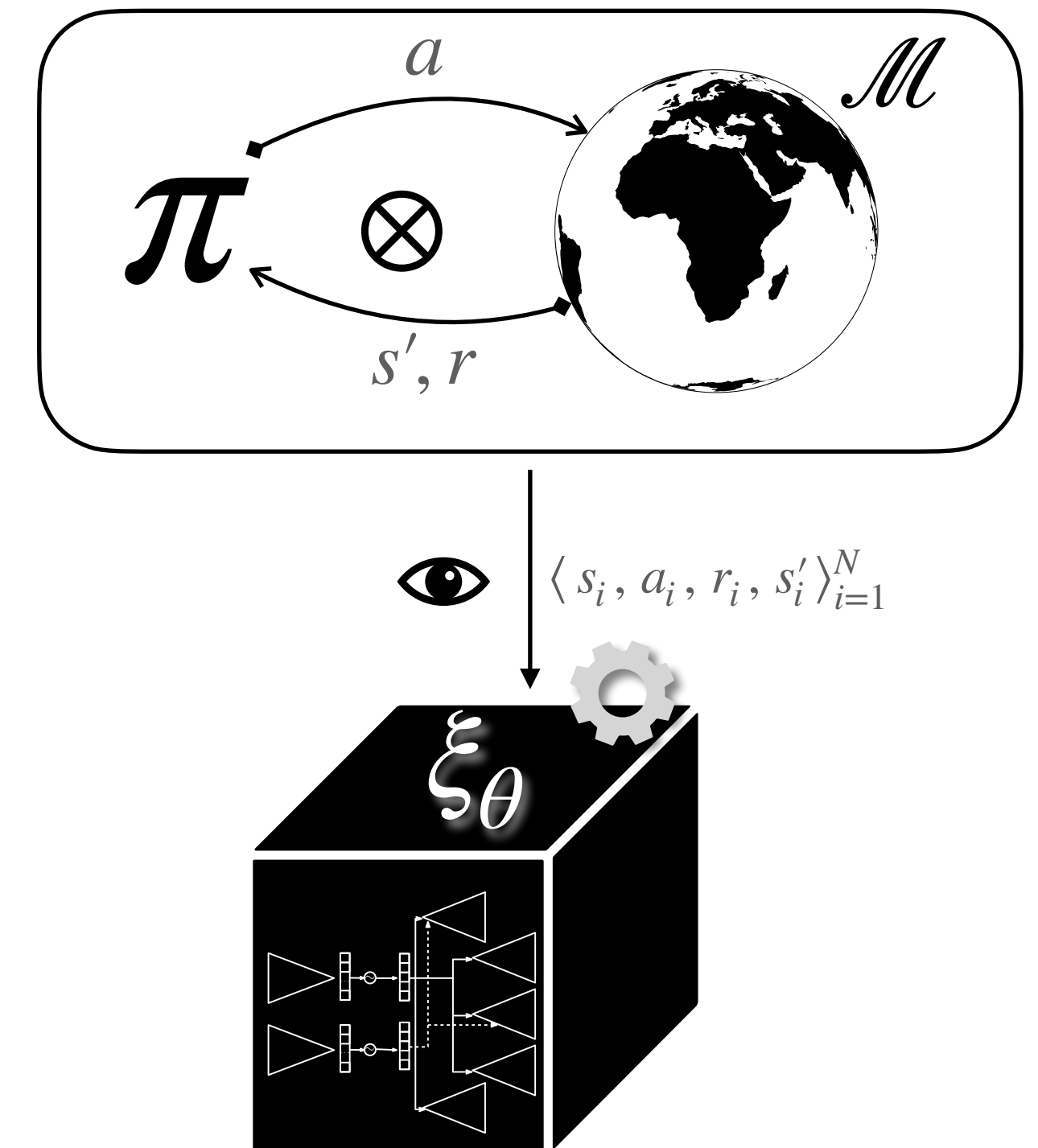
# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$



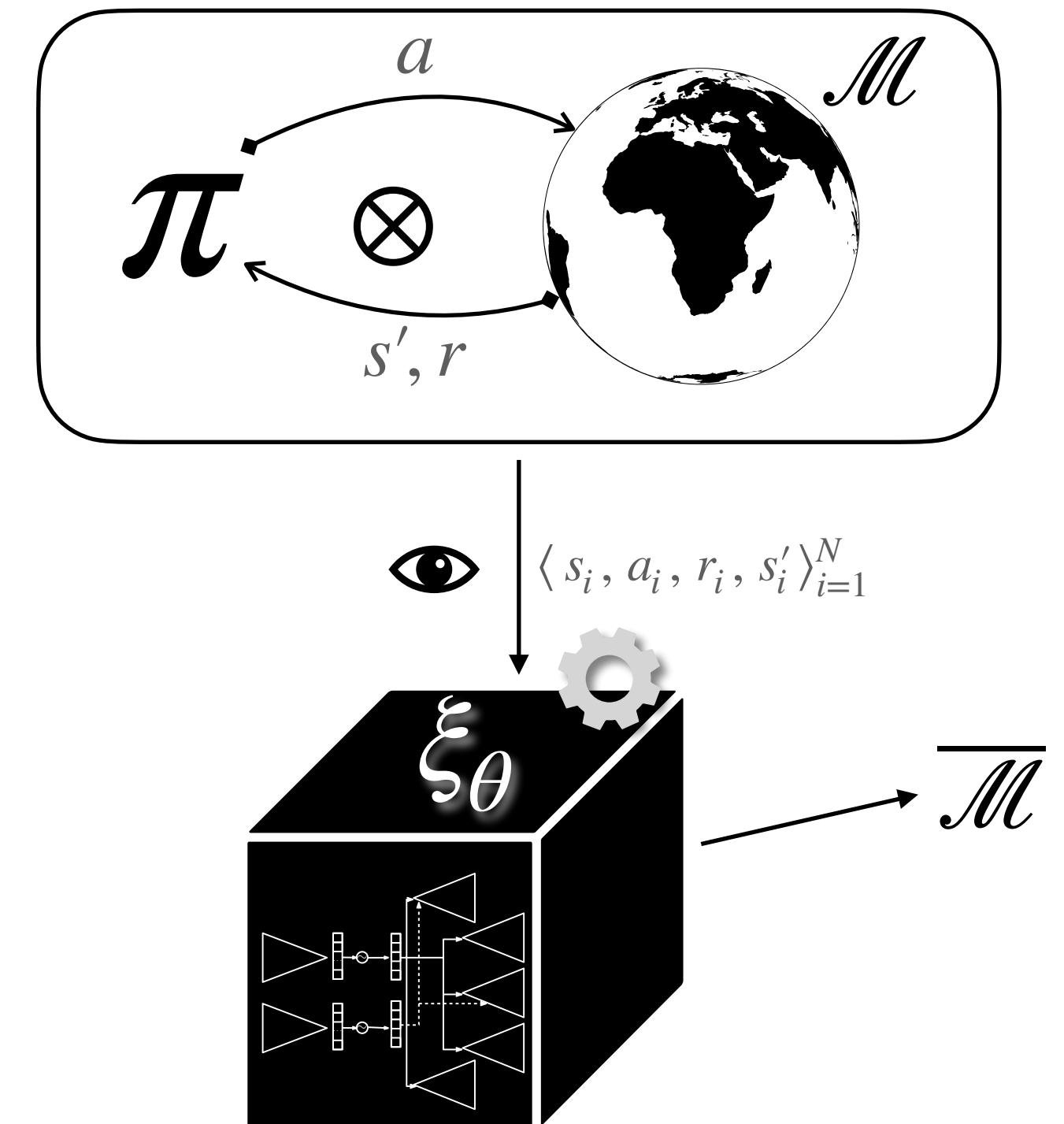
# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$
- **Goal:** learn  $\xi_\theta$  so that we can retrieve:



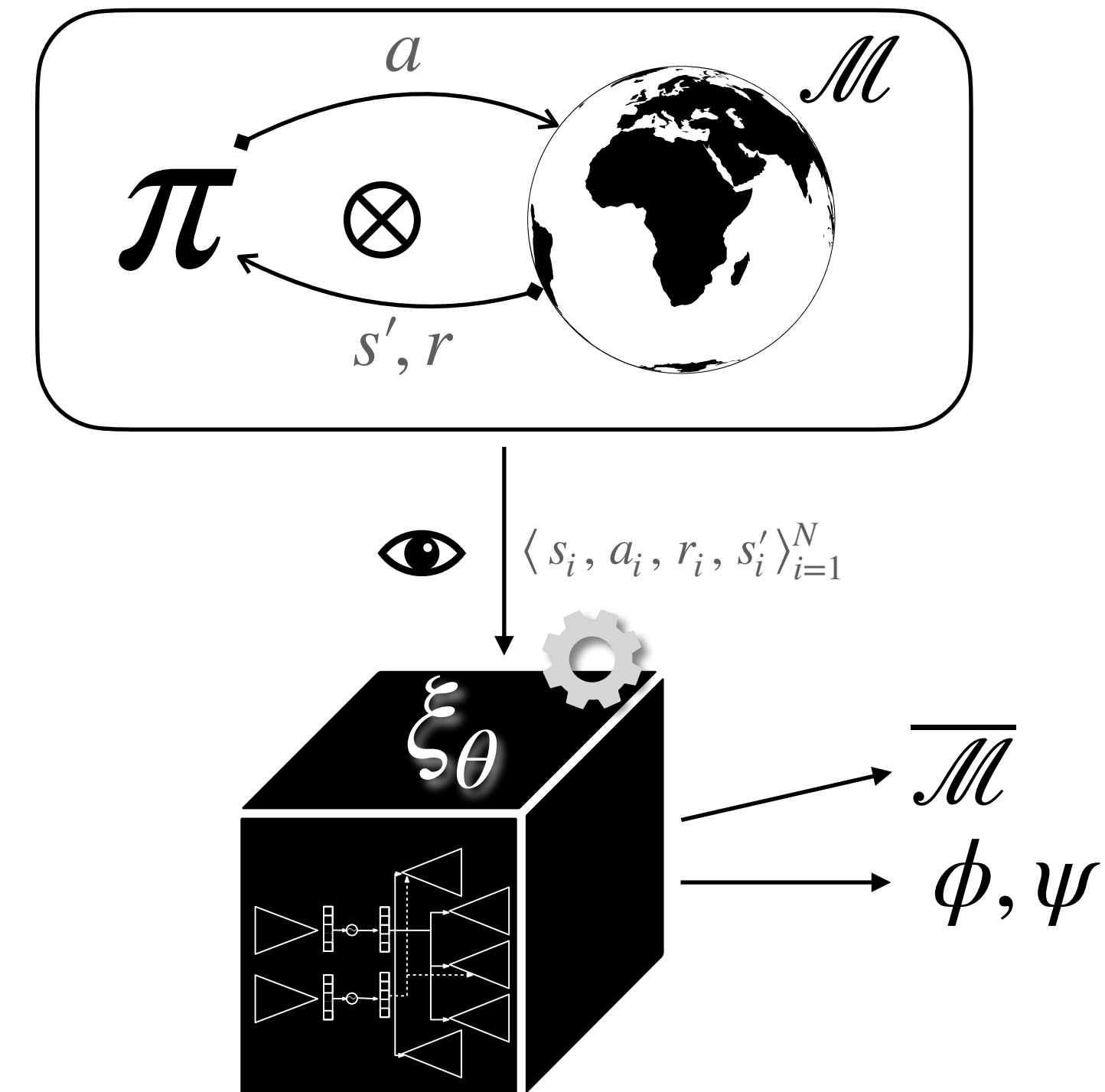
# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$
- **Goal:** learn  $\xi_\theta$  so that we can retrieve:
  - The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$



# Learning the Latent Space Model

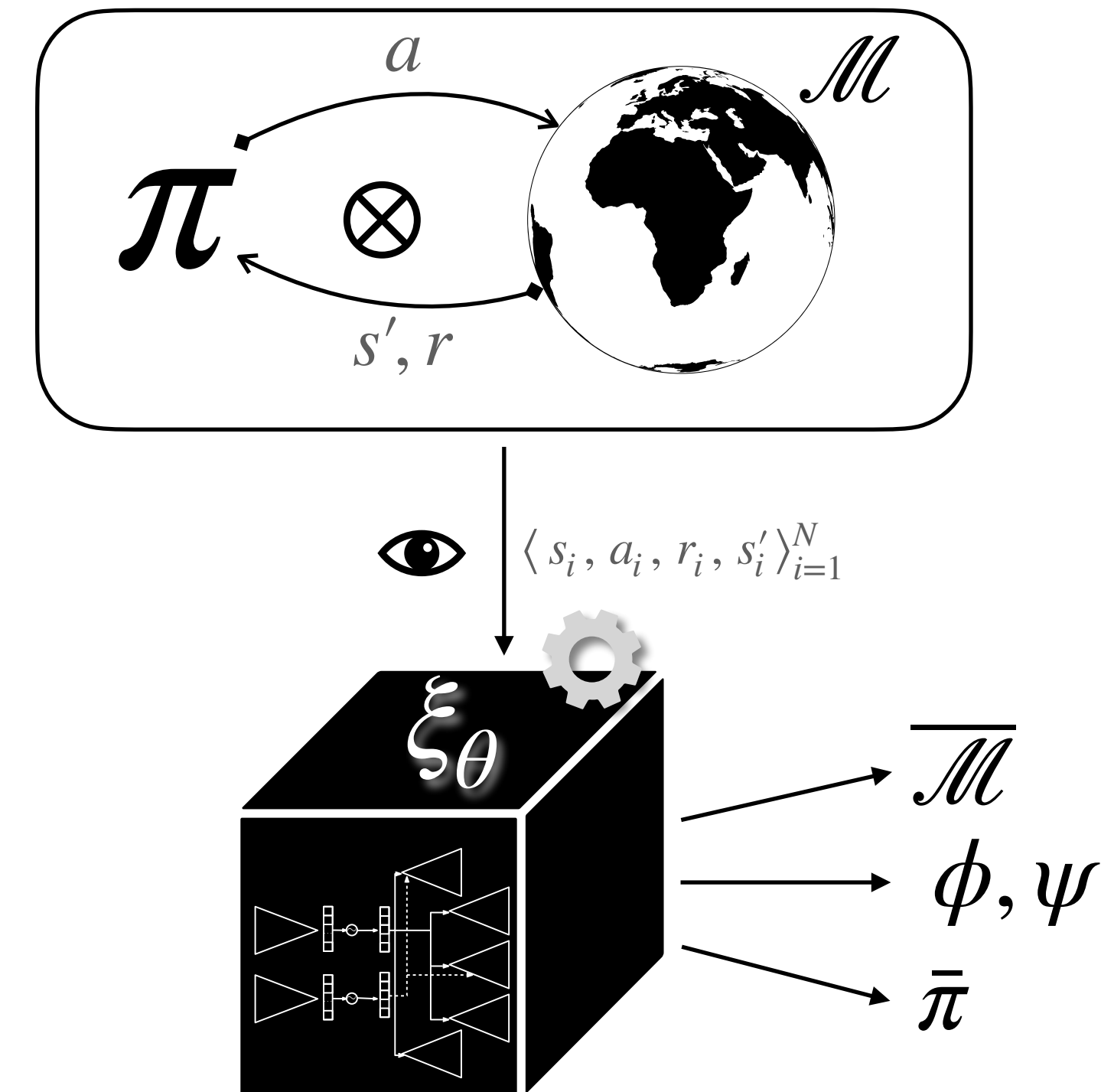
- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$
- **Goal:** learn  $\xi_\theta$  so that we can retrieve:
  - The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$
  - The embedding functions  $\phi, \psi$





# Learning the Latent Space Model

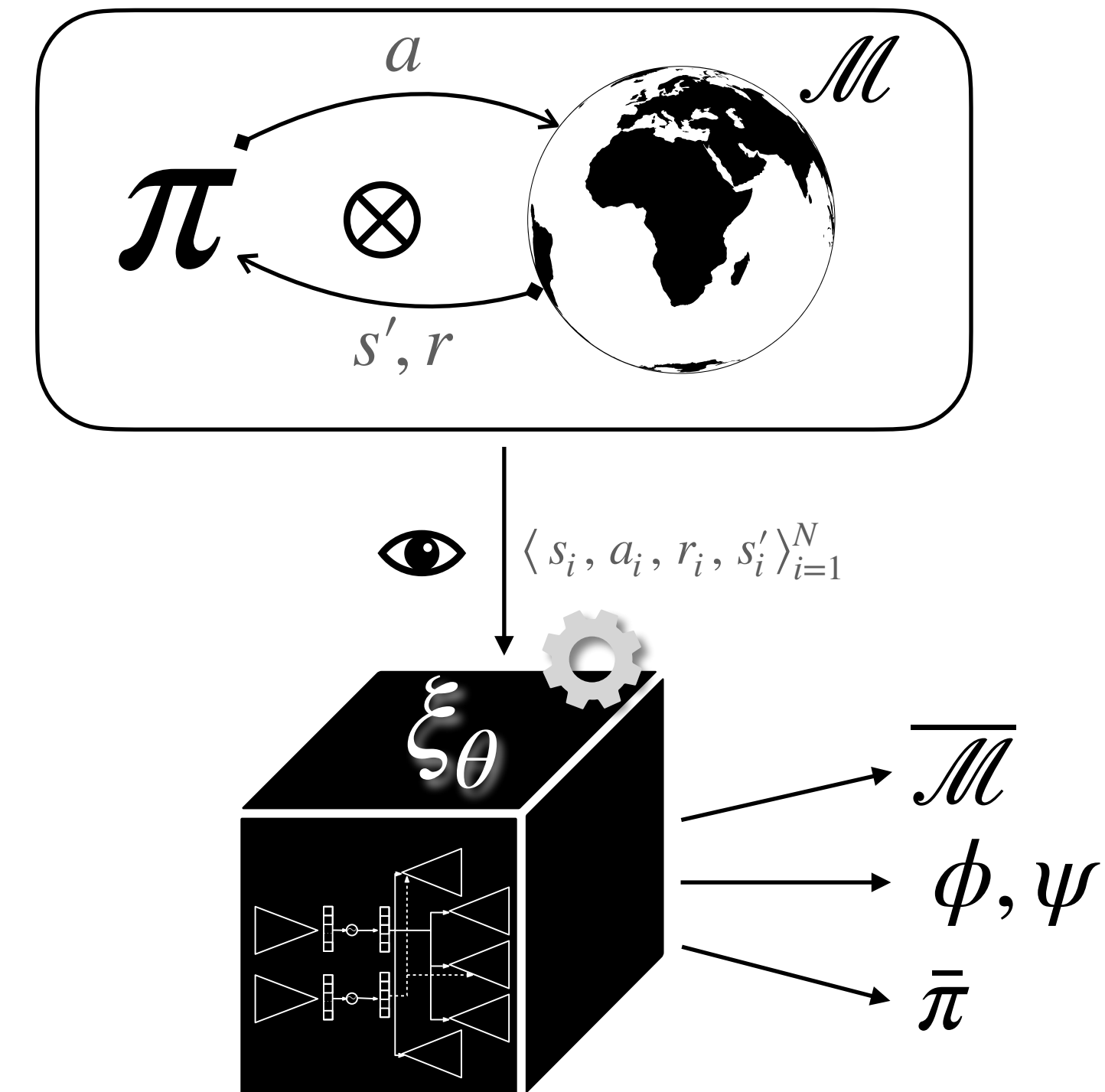
- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$
- **Goal:** learn  $\xi_\theta$  so that we can retrieve:
  - The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$
  - The embedding functions  $\phi, \psi$
  - A latent policy  $\bar{\pi}$  distilled from  $\pi$



# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$
- **Goal:** learn  $\xi_\theta$  so that we can retrieve:
  - The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$
  - The embedding functions  $\phi, \psi$
  - A latent policy  $\bar{\pi}$  distilled from  $\pi$
- Minimize a *discrepancy*  $D$  between  $\mathcal{M} \otimes \pi$  and  $\xi_\theta$

$$\min_{\theta} D(\mathcal{M} \otimes \pi, \xi_\theta)$$



# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$

- **Goal:** learn  $\xi_\theta$  so that we can retrieve:

- The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$

- The embedding functions  $\phi, \psi$

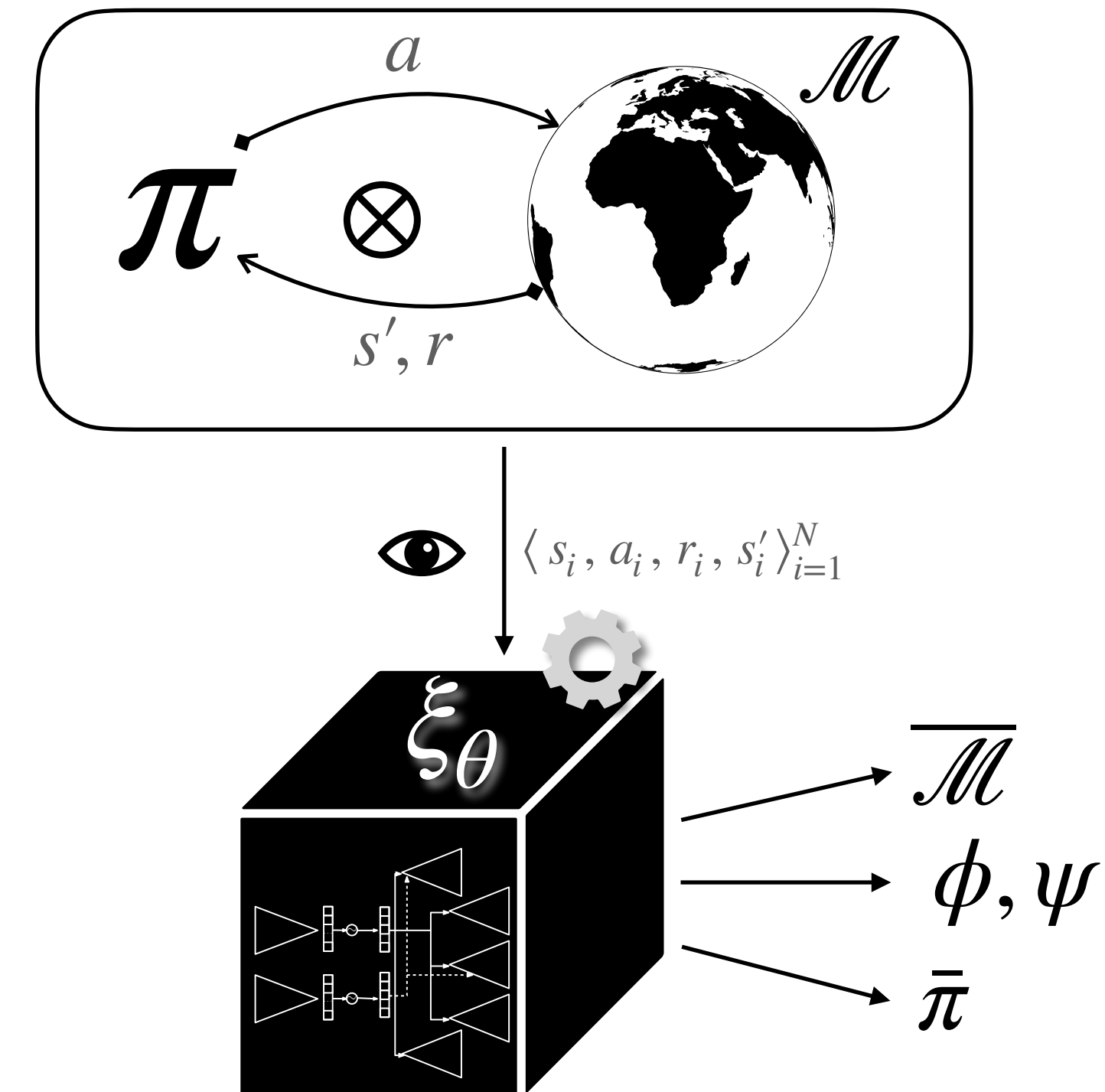
- A latent policy  $\bar{\pi}$  distilled from  $\pi$

- Minimize a *discrepancy*  $D$  between  $\mathcal{M} \otimes \pi$  and  $\xi_\theta$

$$\min_{\theta} D_{KL} (\mathcal{M} \otimes \pi, \xi_\theta)$$

- Choose the *Kullback-Leibler divergence*

$$D_{KL} (P, Q) = \mathbb{E}_{x \sim P} \left[ \log \left( \frac{P(x)}{Q(x)} \right) \right]$$



# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$

- **Goal:** learn  $\xi_\theta$  so that we can retrieve:

- The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$
- The embedding functions  $\phi, \psi$
- A latent policy  $\bar{\pi}$  distilled from  $\pi$

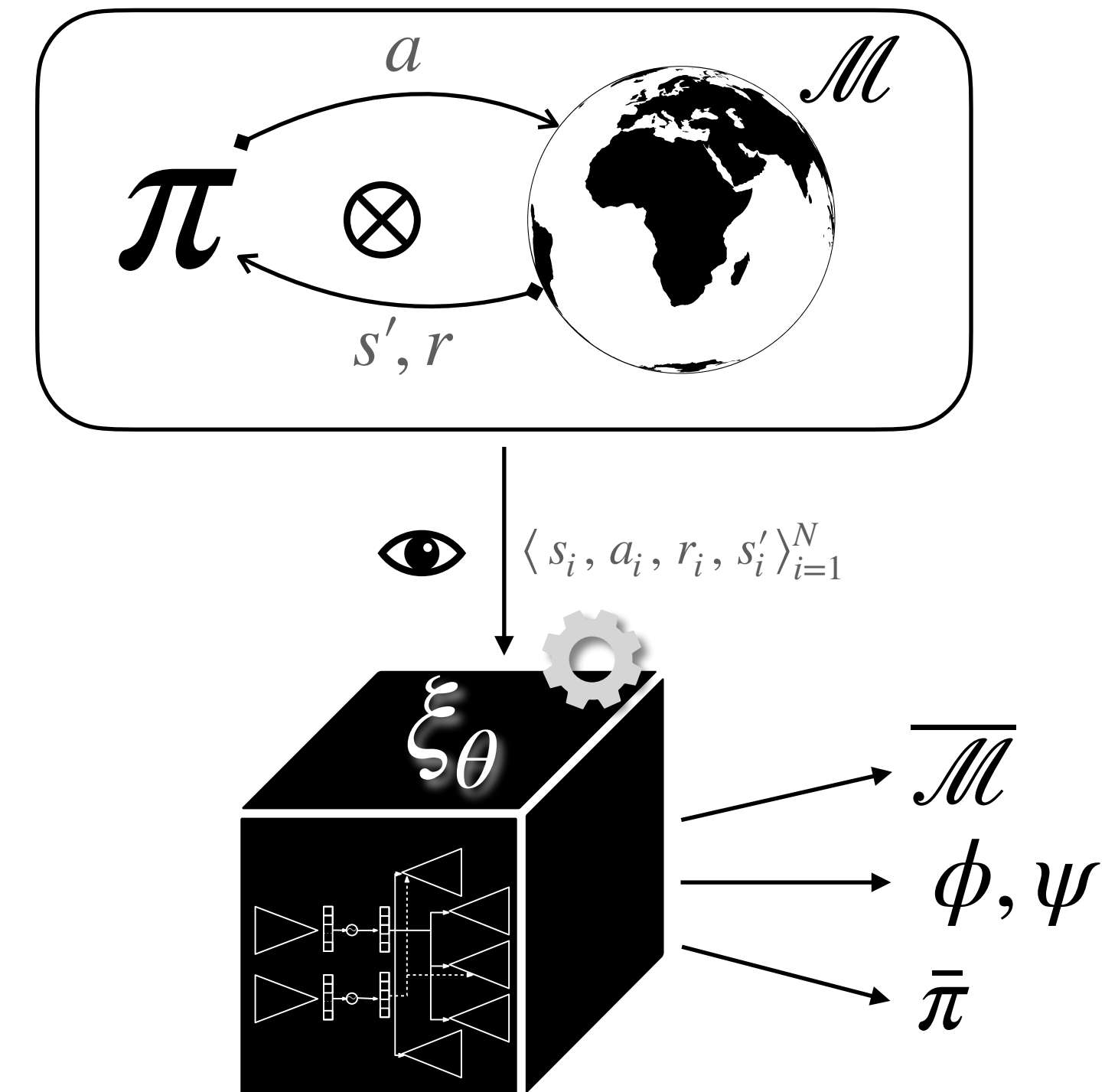
- Minimize a *discrepancy*  $D$  between  $\mathcal{M} \otimes \pi$  and  $\xi_\theta$

$$\min_{\theta} D_{KL} (\mathcal{M} \otimes \pi, \xi_\theta)$$

$$\equiv \max_{\theta} \mathbb{E}_{\tau \sim \mathcal{M} \otimes \pi} [\log \xi_\theta(\tau)]$$

- Choose the *Kullback-Leibler divergence*

$$D_{KL} (P, Q) = \mathbb{E}_{x \sim P} \left[ \log \left( \frac{P(x)}{Q(x)} \right) \right]$$



# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$

- **Goal:** learn  $\xi_\theta$  so that we can retrieve:

- The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$
- The embedding functions  $\phi, \psi$
- A latent policy  $\bar{\pi}$  distilled from  $\pi$

- Minimize a *discrepancy*  $D$  between  $\mathcal{M} \otimes \pi$  and  $\xi_\theta$

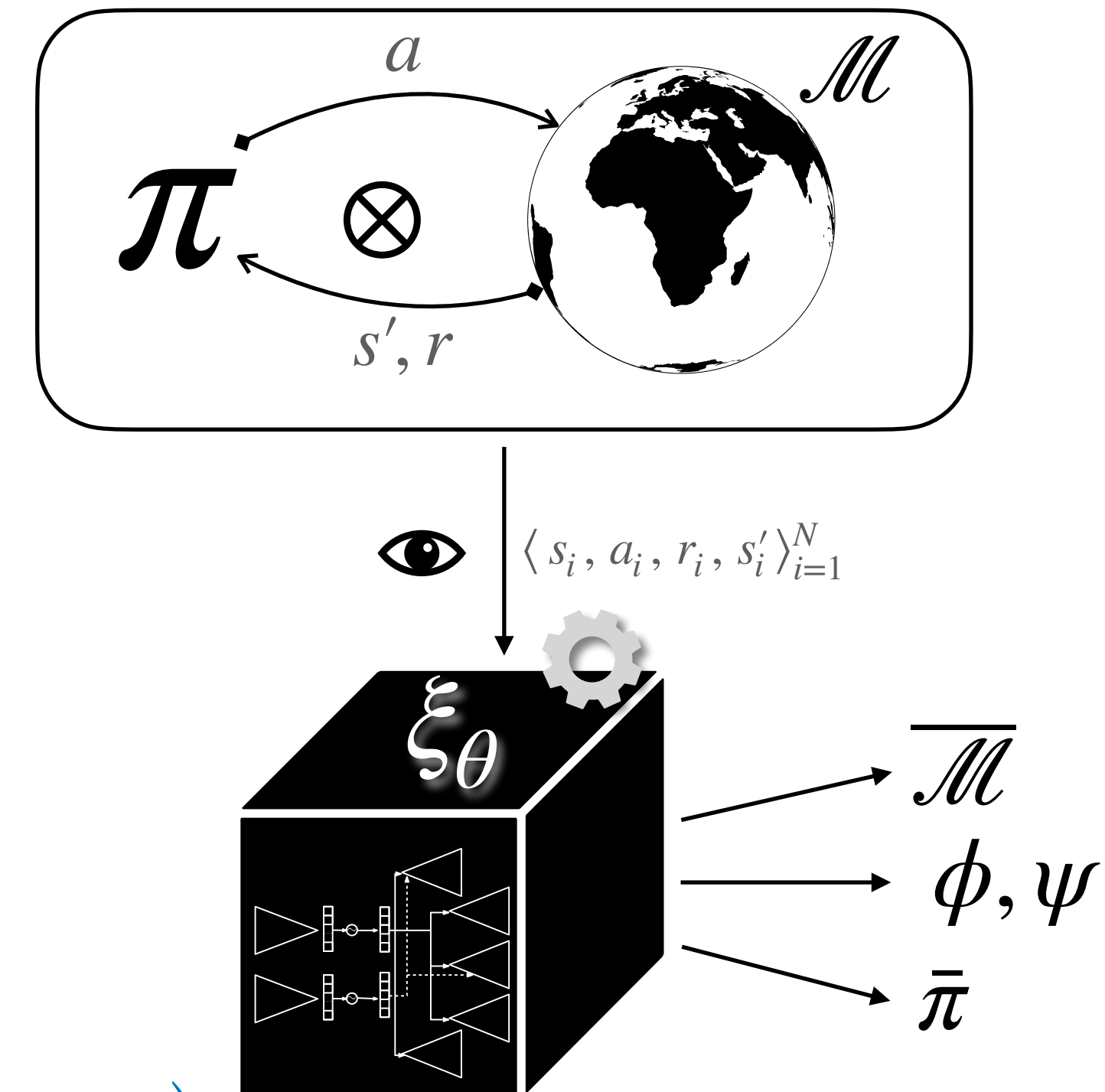
$$\min_{\theta} D_{KL} (\mathcal{M} \otimes \pi, \xi_\theta)$$

$$\equiv \max_{\theta} \mathbb{E}_{\tau \sim \mathcal{M} \otimes \pi} [\log \xi_\theta(\tau)] \geq \max_{\nu, \theta} ELBO (\overline{\mathcal{M}}_\nu, \phi_\nu, \psi_\nu)$$

- Choose the *Kullback-Leibler divergence*

$$D_{KL} (P, Q) = \mathbb{E}_{x \sim P} \left[ \log \left( \frac{P(x)}{Q(x)} \right) \right]$$

(Kingma & Welling, 2014; Hoffman et al., 2013)



$$\max_{\iota, \theta} ELBO(\overline{\mathcal{M}}_{\theta}, \phi_{\iota}, \psi_{\theta})$$

$$\max_{\iota, \theta} ELBO(\overline{\mathcal{M}}_{\theta}, \phi_{\iota}, \psi_{\theta}) = -\min_{\iota, \theta} \{\mathbf{D}_{\iota, \theta} + \mathbf{R}_{\iota, \theta}\}$$

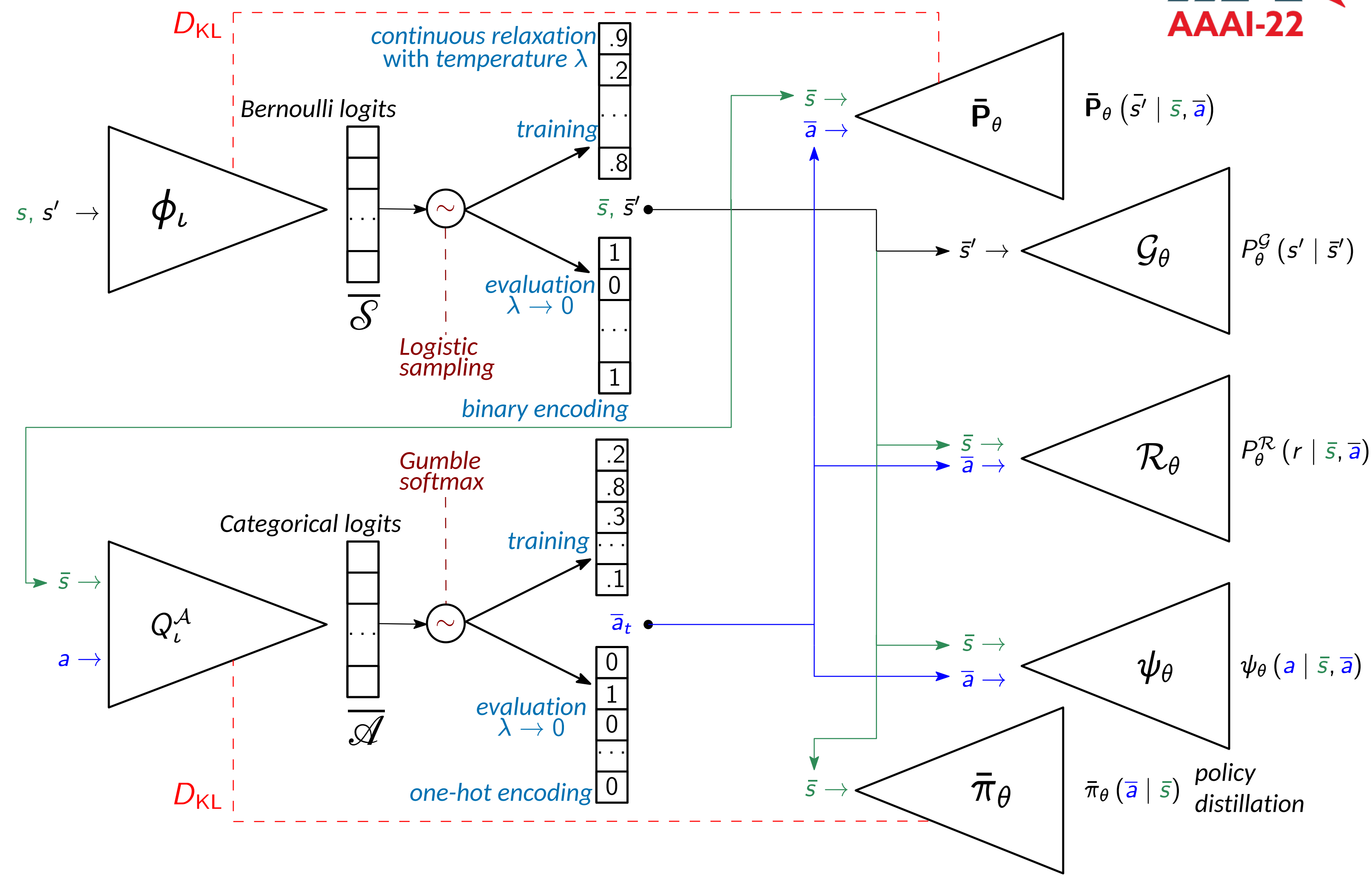
# Variational Markov Decision Process



$$\max_{\iota, \theta} ELBO(\bar{\mathcal{M}}_{\theta}, \phi_{\iota}, \psi_{\theta}) = -\min_{\iota, \theta} \{ \mathbf{D}_{\iota, \theta} + \mathbf{R}_{\iota, \theta} \}$$

$$\mathbf{D}_{\iota, \theta} = - \mathbb{E}_{\substack{s, a, r, s' \sim \xi_{\pi} \\ \bar{s}, \bar{s}' \sim \phi_{\iota}(\cdot | s, s') \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [\log P_{\theta}^G(s' | \bar{s}') + \log \psi_{\theta}(a | \bar{s}, \bar{a}) + \log P_{\theta}^R(r | \bar{s}, \bar{a})]$$

$$\mathbf{R}_{\iota, \theta} = \mathbb{E}_{\substack{s, a, s' \sim \xi_{\pi} \\ \bar{s} \sim \phi_{\iota}(\cdot | s) \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [D_{KL}(\phi_{\iota}(\cdot | s') \| \bar{\mathbf{P}}_{\theta}(\cdot | \bar{s}, \bar{a})) + D_{KL}(Q_{\iota}^A(\cdot | \bar{s}, a) \| \bar{\pi}_{\theta}(\cdot | \bar{s}))]$$





# Variational Markov Decision Process



$$\max_{\iota, \theta} ELBO(\bar{\mathcal{M}}_{\theta}, \phi_{\iota}, \psi_{\theta}) = -\min_{\iota, \theta} \{ \mathbf{D}_{\iota, \theta} + \mathbf{R}_{\iota, \theta} \}$$

$$\mathbf{D}_{\iota, \theta} = - \mathbb{E}_{\substack{s, a, r, s' \sim \xi_{\pi} \\ \bar{s}, \bar{s}' \sim \phi_{\iota}(\cdot | s, s') \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [\log P_{\theta}^G(s' | \bar{s}') + \log \psi_{\theta}(a | \bar{s}, \bar{a}) + \log P_{\theta}^R(r | \bar{s}, \bar{a})]$$

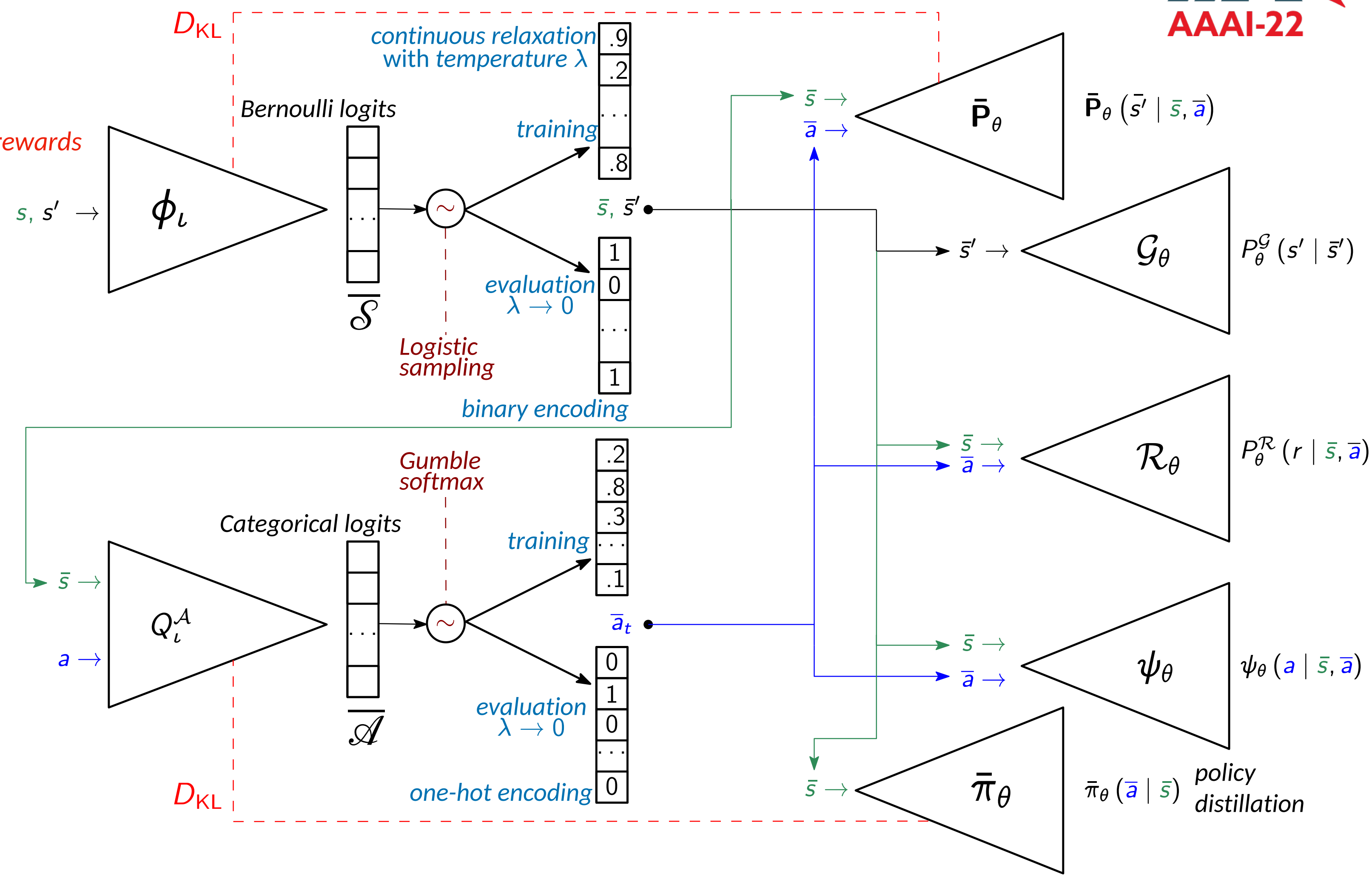
Log-likelihood of rewards

$$\mathbf{R}_{\iota, \theta} = \mathbb{E}_{\substack{s, a, s' \sim \xi_{\pi} \\ \bar{s} \sim \phi_{\iota}(\cdot | s) \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [D_{KL}(\phi_{\iota}(\cdot | s') \| \bar{\mathbf{P}}_{\theta}(\cdot | \bar{s}, \bar{a})) + D_{KL}(Q_{\iota}^A(\cdot | \bar{s}, a) \| \bar{\pi}_{\theta}(\cdot | \bar{s}))]$$

Variational version of local transition loss

$$L_{\mathbf{P}}^{\xi_{\pi}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\pi}} W_{d_{\bar{s}}}(\phi_{\mathbf{P}}(\cdot | s, \bar{a}), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a})) \leq \mathbb{E}_{s, \bar{a}, s' \sim \xi_{\pi}} W_{d_{\bar{s}}}(\phi(\cdot | s'), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a}))$$

- Variational proxies to local losses



# Variational Markov Decision Process



$$\max_{\iota, \theta} ELBO(\bar{\mathcal{M}}_{\theta}, \phi_{\iota}, \psi_{\theta}) = -\min_{\iota, \theta} \{ \mathbf{D}_{\iota, \theta} + \mathbf{R}_{\iota, \theta} \}$$

$$\mathbf{D}_{\iota, \theta} = - \mathbb{E}_{\substack{s, a, r, s' \sim \xi_{\pi} \\ \bar{s}, \bar{s}' \sim \phi_{\iota}(\cdot | s, s') \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [\log P_{\theta}^G(s' | \bar{s}') + \log \psi_{\theta}(a | \bar{s}, \bar{a}) + \log P_{\theta}^R(r | \bar{s}, \bar{a})]$$

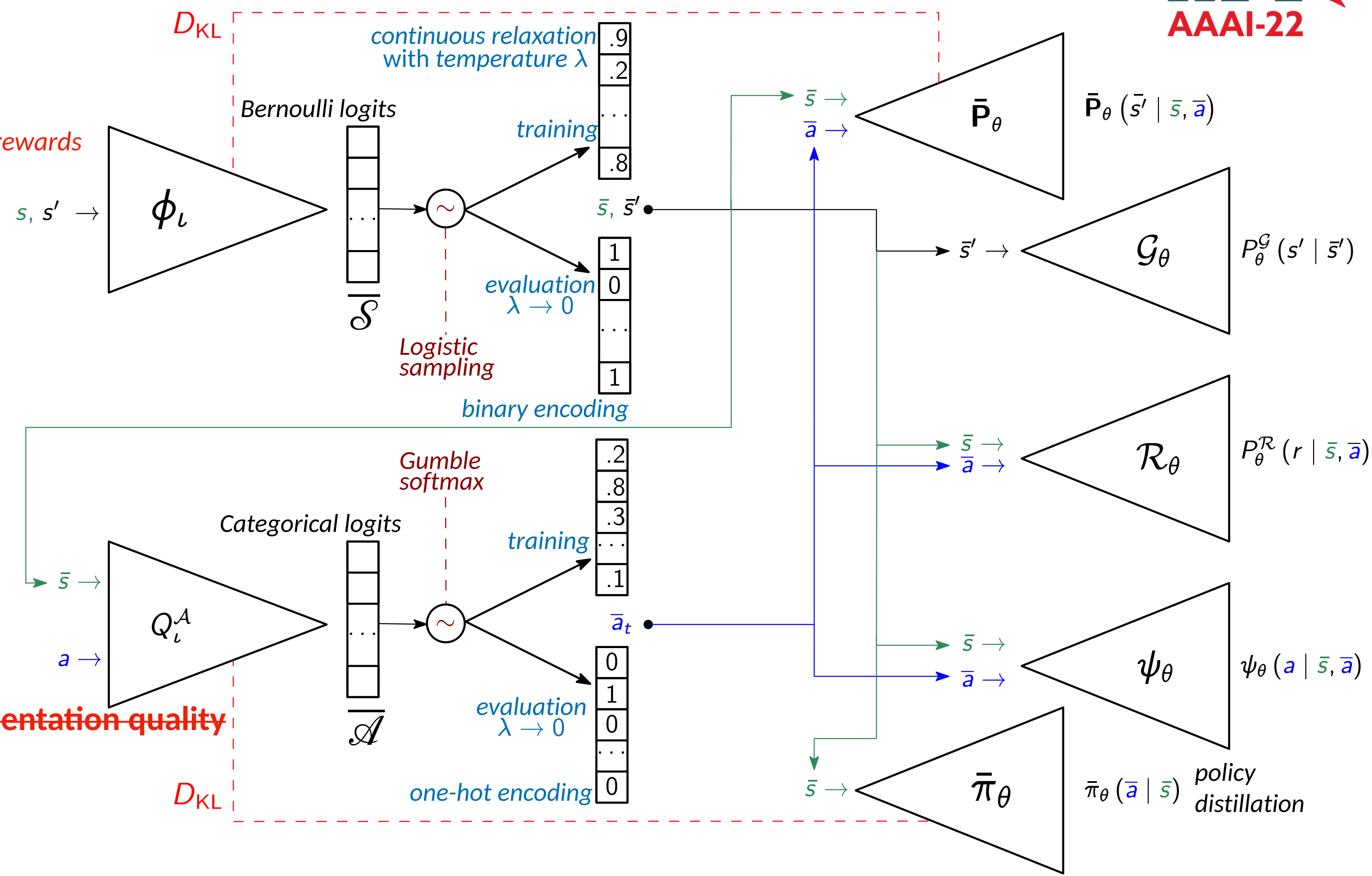
Log-likelihood of rewards

$$\mathbf{R}_{\iota, \theta} = \mathbb{E}_{\substack{s, a, s' \sim \xi_{\pi} \\ \bar{s} \sim \phi_{\iota}(\cdot | s) \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [D_{KL}(\phi_{\iota}(\cdot | s') \parallel \bar{\mathbf{P}}_{\theta}(\cdot | \bar{s}, \bar{a})) + D_{KL}(Q_{\iota}^A(\cdot | \bar{s}, a) \parallel \bar{\pi}_{\theta}(\cdot | \bar{s}))]$$

Variational version of local transition loss

$$L_{\mathbf{P}}^{\xi_{\pi}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\pi}} W_{d_{\bar{s}}}(\phi_{\mathbf{P}}(\cdot | s, \bar{a}), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a})) \leq \mathbb{E}_{s, \bar{a}, s' \sim \xi_{\pi}} W_{d_{\bar{s}}}(\phi(\cdot | s'), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a}))$$

- Variational proxies to local losses
  - No learning guarantee: abstraction-quality, representation-quality



# Variational Markov Decision Process



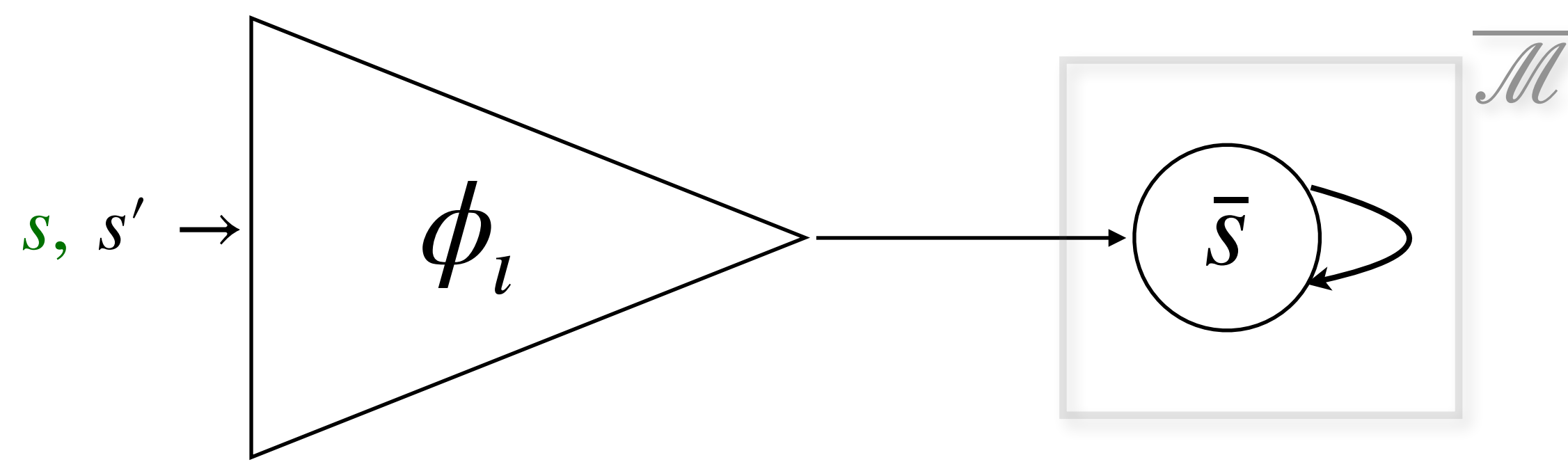
$$\max_{\iota, \theta} ELBO(\overline{\mathcal{M}}_{\theta}, \phi_{\iota}, \psi_{\theta}) = -\min_{\iota, \theta} \{ \mathbf{D}_{\iota, \theta} + \mathbf{R}_{\iota, \theta} \}$$

$$\mathbf{D}_{\iota, \theta} = - \mathbb{E}_{\substack{s, a, r, s' \sim \xi_{\pi} \\ \bar{s}, \bar{s}' \sim \phi_{\iota}(\cdot | s, s') \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [\log P_{\theta}^G(s' | \bar{s}') + \log \psi_{\theta}(a | \bar{s}, \bar{a}) + \log P_{\theta}^R(r | \bar{s}, \bar{a})]$$

$$\mathbf{R}_{\iota, \theta} = \mathbb{E}_{\substack{s, a, s' \sim \xi_{\pi} \\ \bar{s} \sim \phi_{\iota}(\cdot | s) \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [D_{KL}(\phi_{\iota}(\cdot | s') \parallel \bar{\mathbf{P}}_{\theta}(\cdot | \bar{s}, \bar{a})) + D_{KL}(Q_{\iota}^A(\cdot | \bar{s}, a) \parallel \bar{\pi}_{\theta}(\cdot | \bar{s}))]$$

Variational version of local transition loss

$$L_{\mathbf{P}}^{\xi_{\pi}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\pi}} W_{d_{\bar{s}}}(\phi_{\mathbf{P}}(\cdot | s, \bar{a}), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a})) \leq \mathbb{E}_{s, \bar{a}, s' \sim \xi_{\pi}} W_{d_{\bar{s}}}(\phi(\cdot | s'), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a}))$$



- Variational proxies to local losses
  - ➔ No learning guarantee: abstraction-quality, representation-quality
  - ➔ Mode collapse

# Variational Markov Decision Process

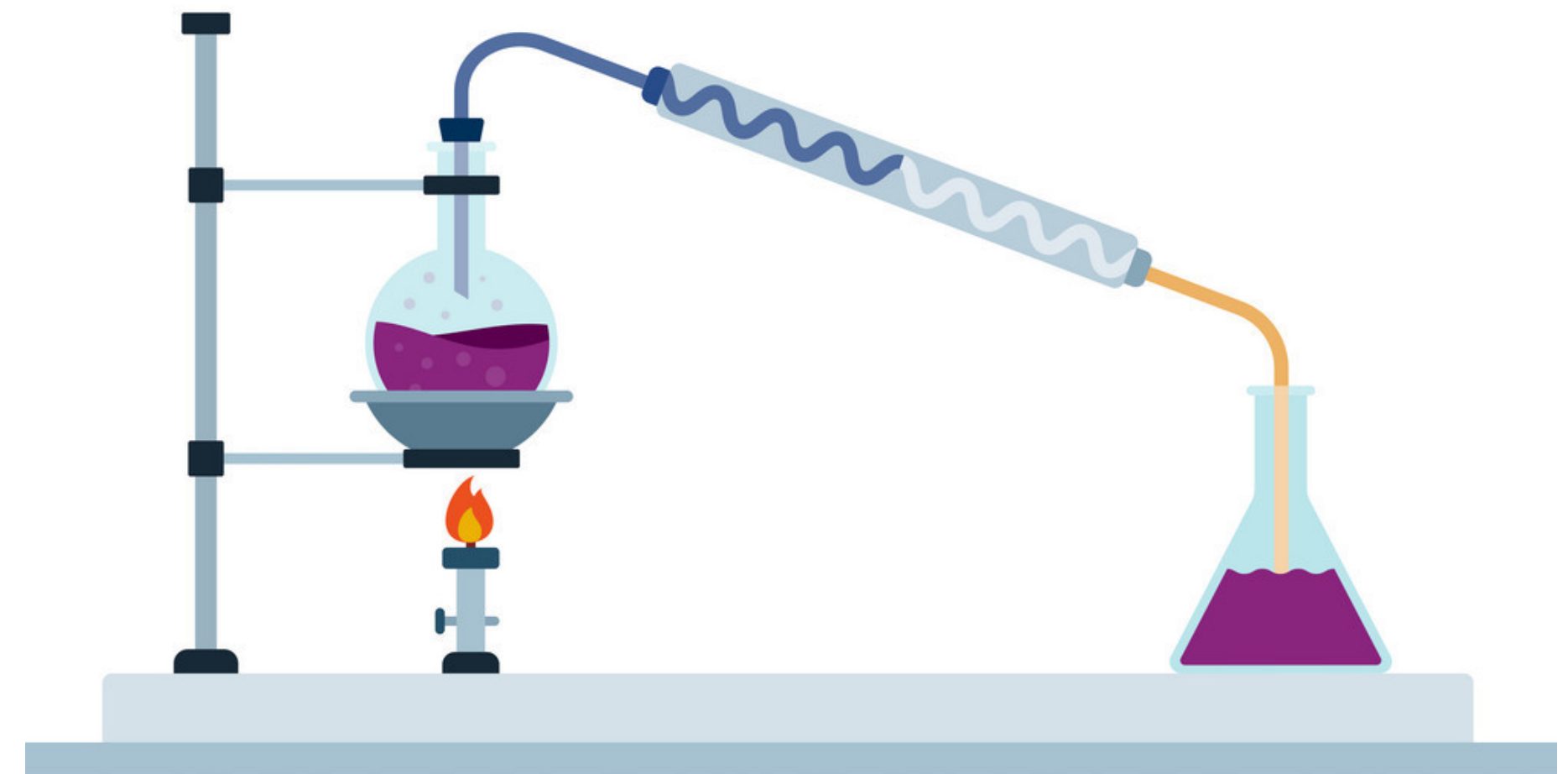
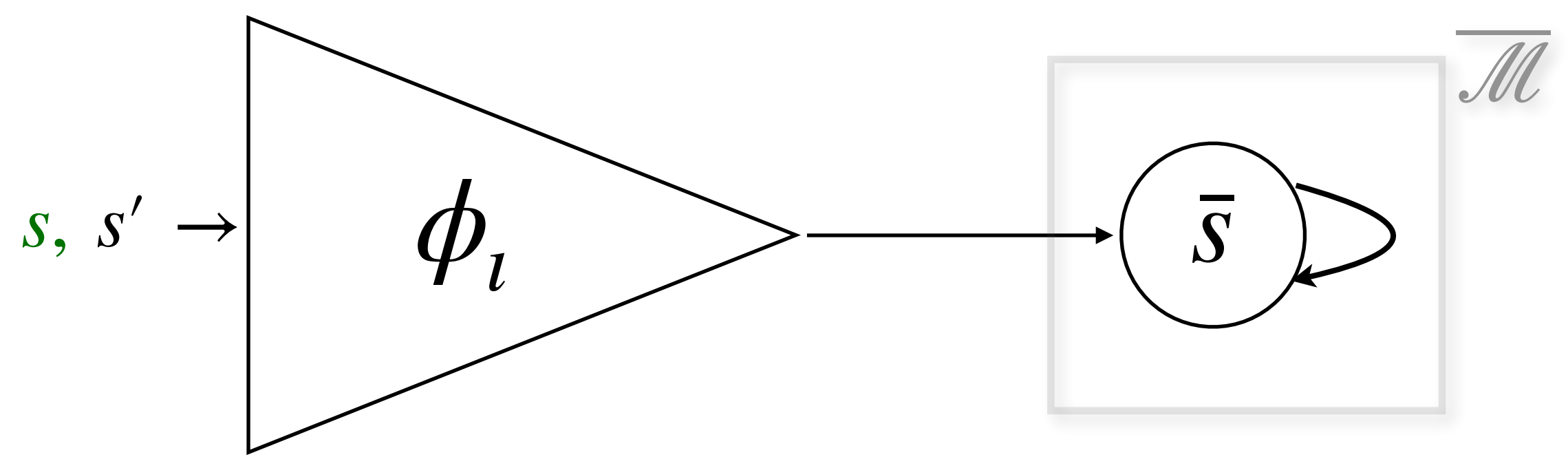


$$\max_{\iota, \theta} ELBO(\overline{\mathcal{M}}_{\theta}, \phi_{\iota}, \psi_{\theta}) = -\min_{\iota, \theta} \{ \mathbf{D}_{\iota, \theta} + \mathbf{R}_{\iota, \theta} \}$$

$$\mathbf{D}_{\iota, \theta} = - \mathbb{E}_{\substack{s, a, r, s' \sim \xi_{\pi} \\ \bar{s}, \bar{s}' \sim \phi_{\iota}(\cdot | s, s') \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [\log P_{\theta}^G(s' | \bar{s}') + \log \psi_{\theta}(a | \bar{s}, \bar{a}) + \log P_{\theta}^R(r | \bar{s}, \bar{a})]$$

$$\mathbf{R}_{\iota, \theta} = \mathbb{E}_{\substack{s, a, s' \sim \xi_{\pi} \\ \bar{s} \sim \phi_{\iota}(\cdot | s) \\ \bar{a} \sim Q_{\iota}^A(\cdot | \bar{s}, a)}} [D_{KL}(\phi_{\iota}(\cdot | s') \parallel \bar{\mathbf{P}}_{\theta}(\cdot | \bar{s}, \bar{a})) + D_{KL}(Q_{\iota}^A(\cdot | \bar{s}, a) \parallel \bar{\pi}_{\theta}(\cdot | \bar{s}))]$$

$$L_{\mathbf{P}}^{\xi_{\pi}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\pi}} W_{d_{\bar{s}}}(\phi_{\mathbf{P}}(\cdot | s, \bar{a}), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a})) \leq \mathbb{E}_{s, \bar{a}, s' \sim \xi_{\pi}} W_{d_{\bar{s}}}(\phi(\cdot | s'), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a}))$$

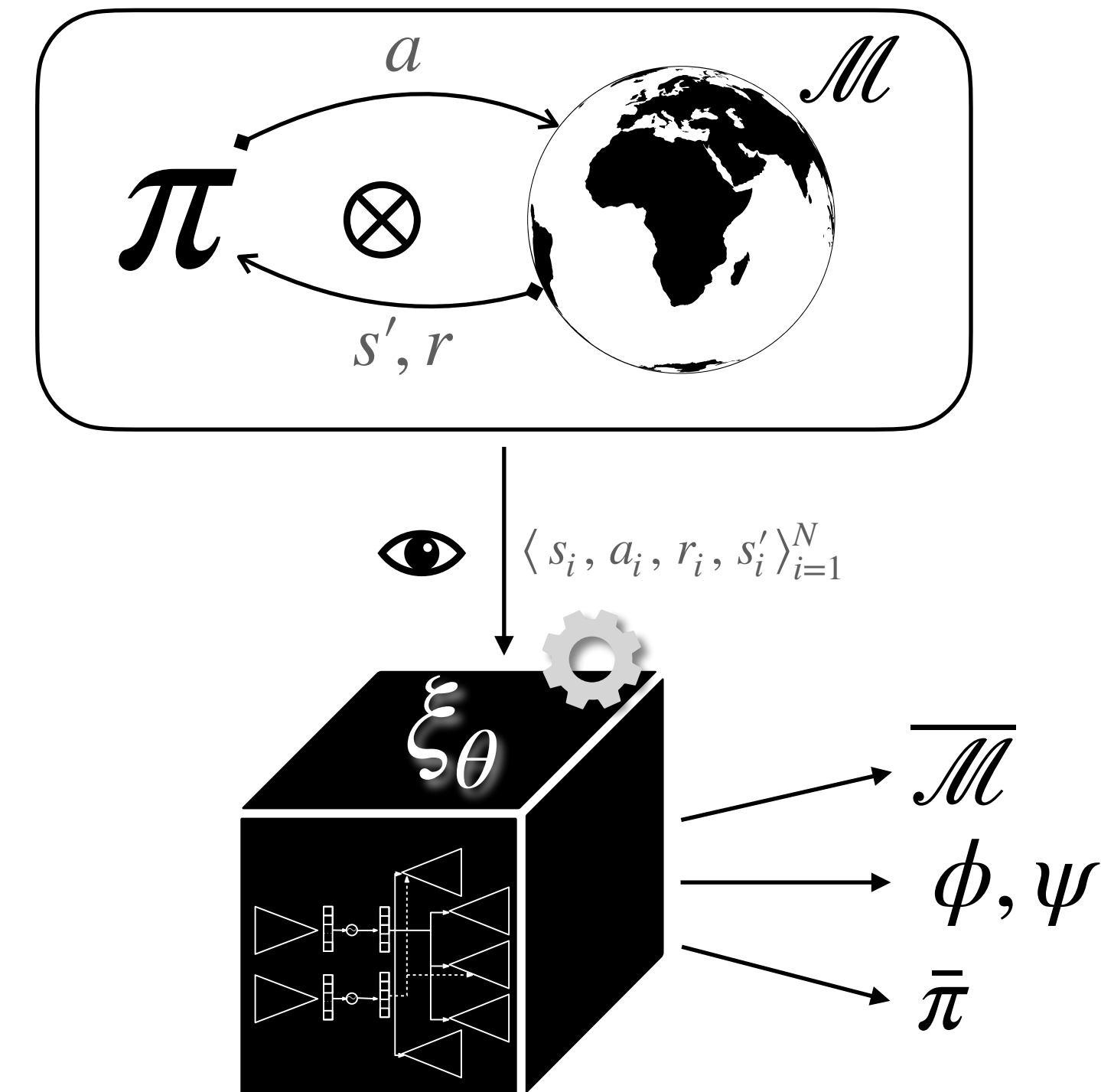


- Variational proxies to local losses
  - ➔ No learning guarantee: abstraction-quality, representation-quality
  - ➔ Mode collapse
  - ➔ Fix: annealing scheme, extra entropy regularization term, prioritized experience replay, ...

# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$
- **Goal:** learn  $\xi_\theta$  so that we can retrieve:
  - The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$
  - The embedding functions  $\phi, \psi$
  - A latent policy  $\bar{\pi}$  distilled from  $\pi$
- Minimize a *discrepancy*  $D$  between  $\mathcal{M} \otimes \pi$  and  $\xi_\theta$

$$\min_{\theta} D(\mathcal{M} \otimes \pi, \xi_\theta)$$



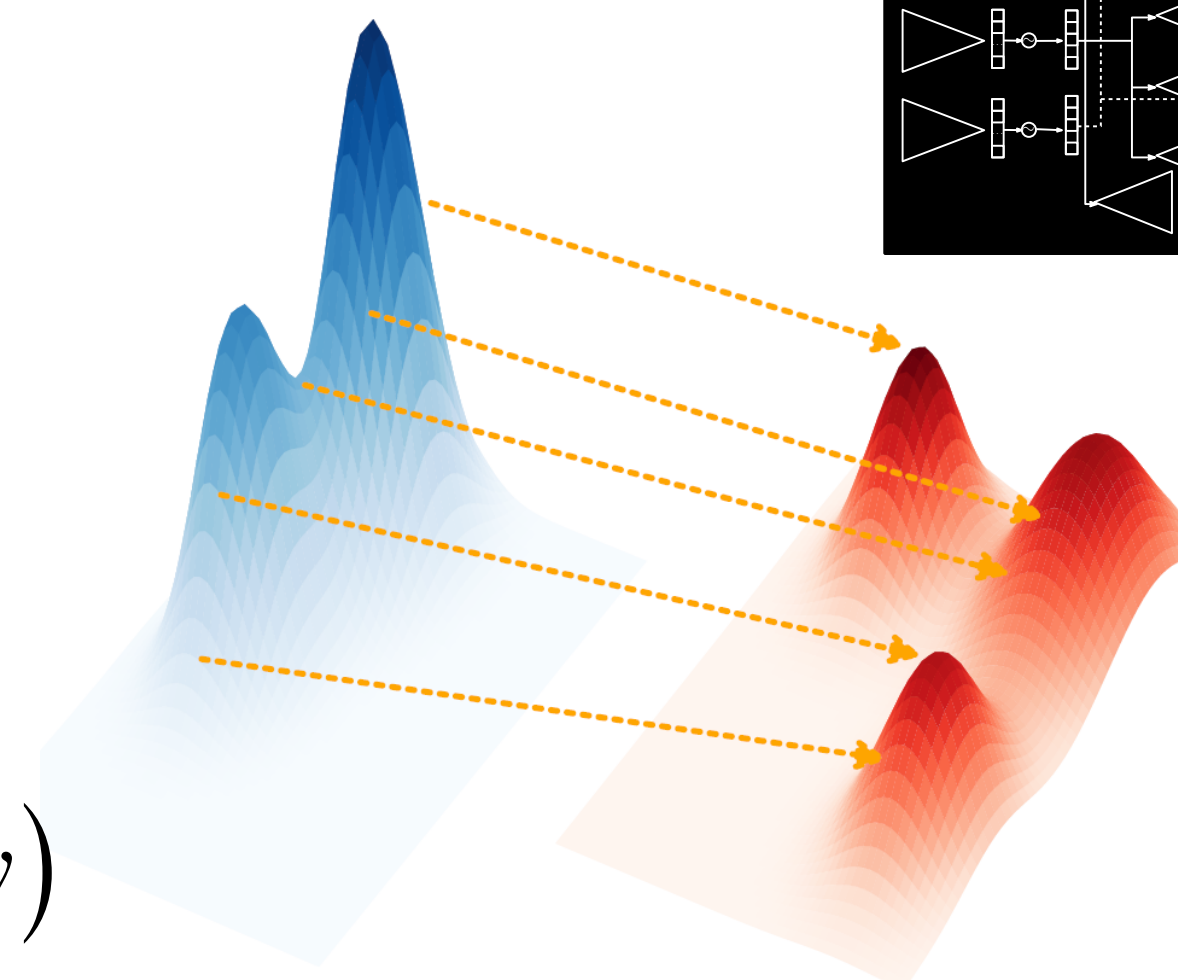
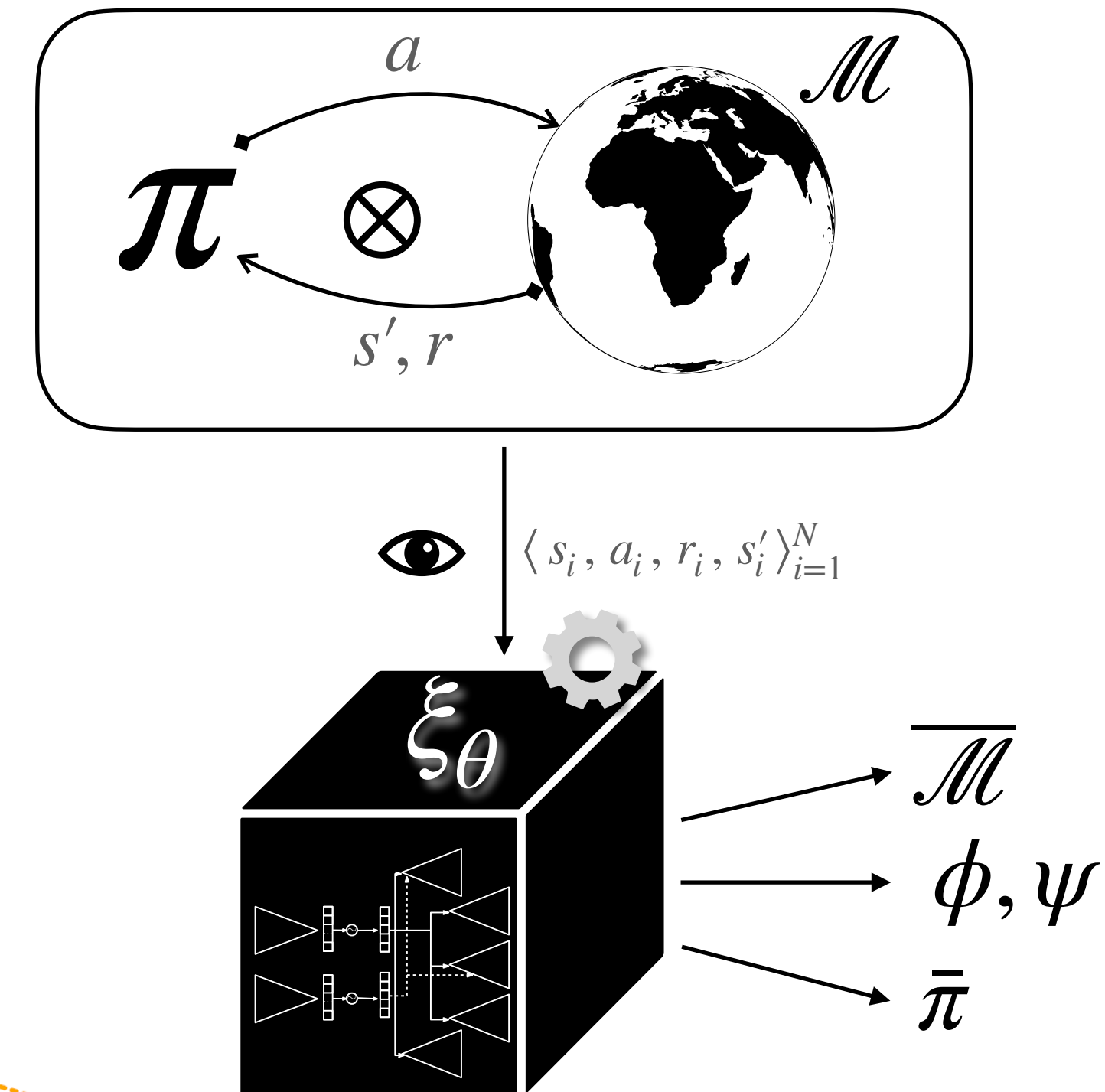
# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$
- **Goal:** learn  $\xi_\theta$  so that we can retrieve:
  - The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$
  - The embedding functions  $\phi, \psi$
  - A latent policy  $\bar{\pi}$  distilled from  $\pi$
- Minimize a *discrepancy*  $D$  between  $\mathcal{M} \otimes \pi$  and  $\xi_\theta$

$$\min_{\theta} W(\mathcal{M} \otimes \pi, \xi_\theta)$$

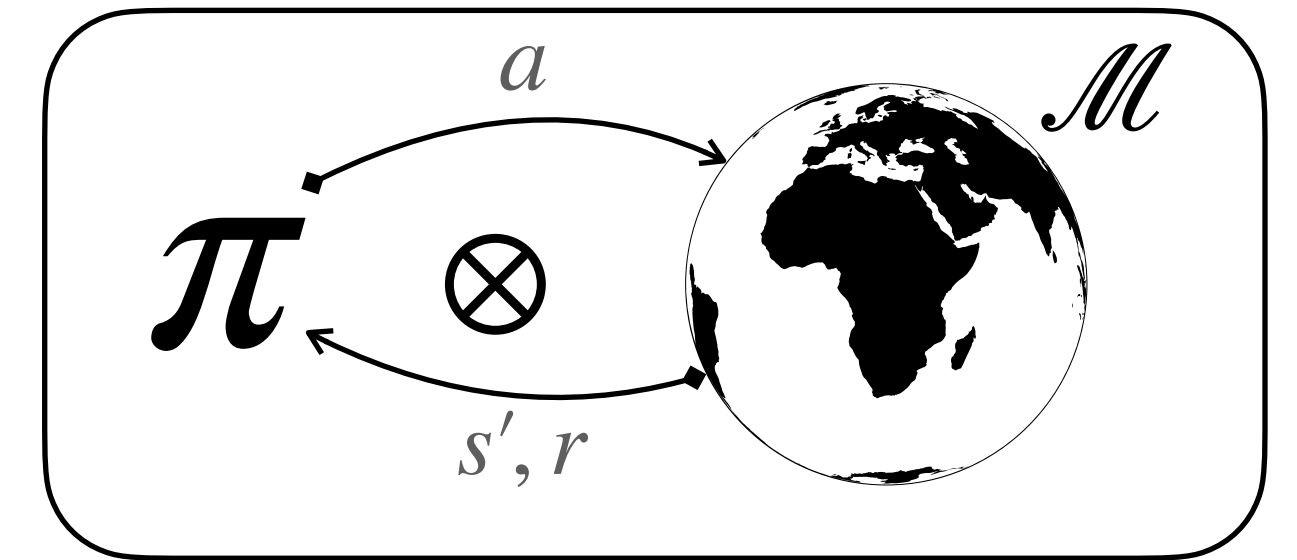
- Choose the *Wasserstein Distance*

$$W(P, Q) = \inf_{\lambda \in \Lambda(P, Q)} \mathbb{E}_{x, y \sim \lambda} d(x, y) = \sup_{\|f\| \leq 1} \mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{y \sim Q} f(y)$$

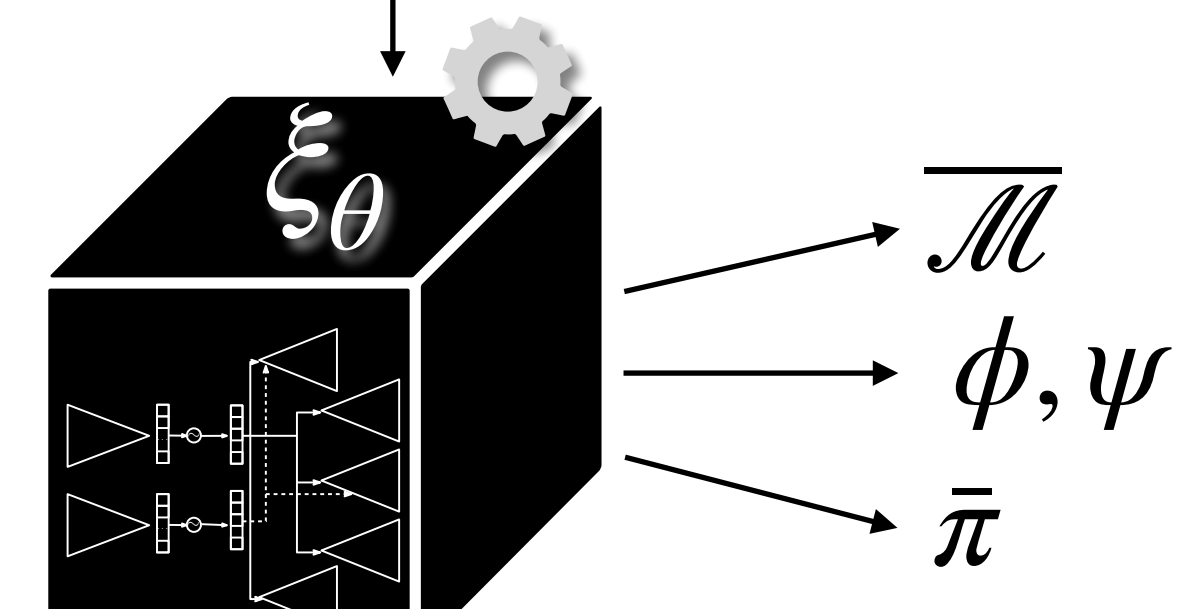


# Learning the Latent Space Model

- Train a *behavioral model*  $\xi_\theta$  by learning from traces produced by executing the RL policy  $\pi$  in the original model  $\mathcal{M}$
- **Goal:** learn  $\xi_\theta$  so that we can retrieve:
  - The latent MDP  $\overline{\mathcal{M}} = \langle \overline{\mathcal{S}}, \overline{\mathcal{A}}, \overline{\mathcal{R}}, \overline{\mathbf{P}}, \ell \rangle$
  - The embedding functions  $\phi, \psi$
  - A latent policy  $\bar{\pi}$  distilled from  $\pi$
- Minimize a *discrepancy*  $D$  between  $\mathcal{M} \otimes \pi$  and  $\xi_\theta$



$\langle s_i, a_i, r_i, s'_i \rangle_{i=1}^N$



$$\min_{\theta} W(\mathcal{M} \otimes \pi, \xi_\theta)$$

$$\leq \min_{\iota, \theta} \mathbb{E}_{s, a, s' \sim \xi_\pi} \mathbb{E}_{\bar{s}, \bar{a}, \bar{s}' \sim \phi_\iota(\cdot | s, a, s')} \|\langle s, a, s' \rangle - \langle \mathcal{G}_\theta(\bar{s}), \psi_\theta(\bar{s}, \bar{a}), \mathcal{G}_\theta(\bar{s}') \rangle\| + L_{\mathcal{R}}^{\xi_\pi} + \beta \left( \mathcal{W}_{\xi_\pi} + L_{\mathbf{P}}^{\xi_\pi} \right)$$

- Choose the *Wasserstein Distance*

$$W(P, Q) = \inf_{\lambda \in \Lambda(P, Q)} \mathbb{E}_{x, y \sim \lambda} d(x, y) = \sup_{\|f\| \leq 1} \mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{y \sim Q} f(y)$$

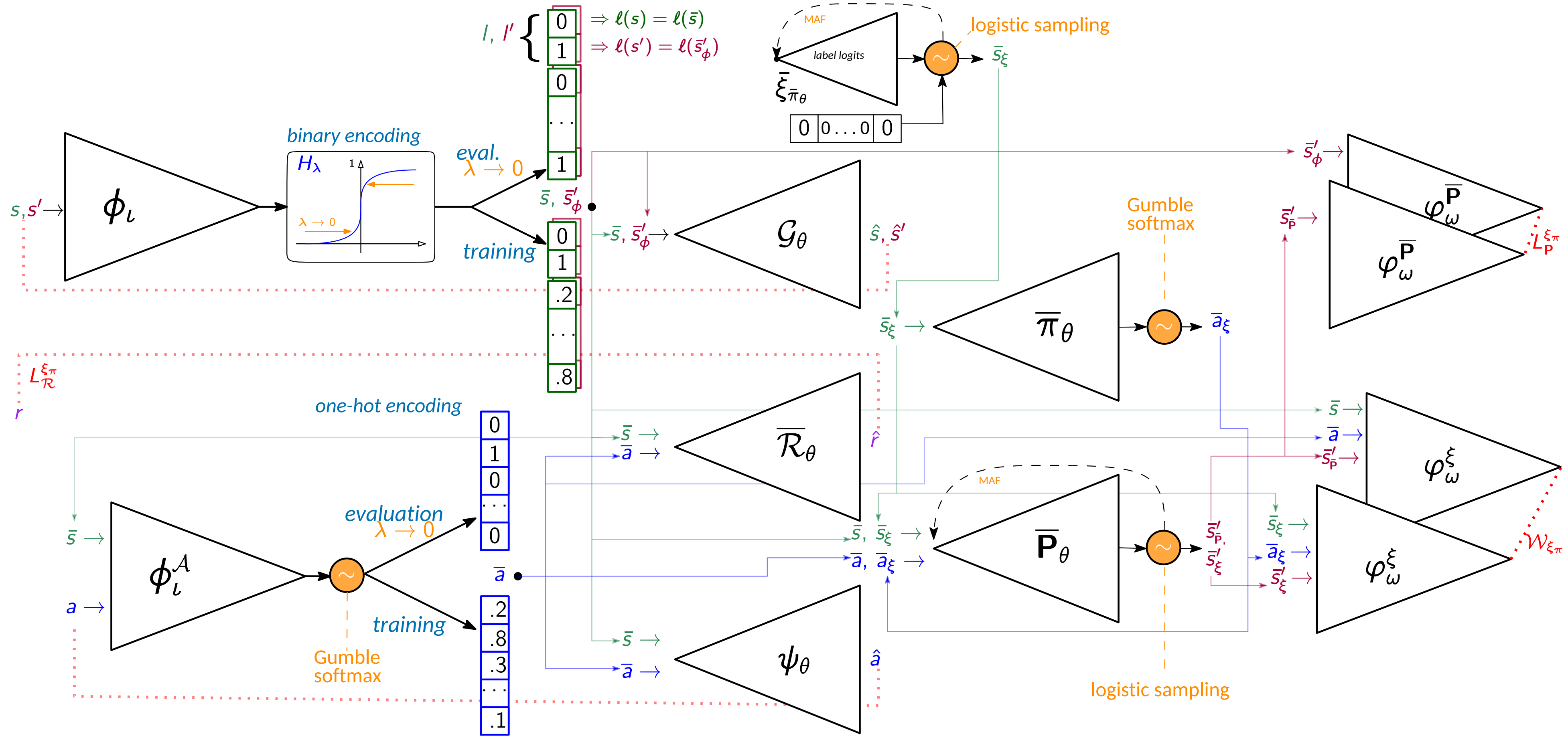
$$\min_{\iota, \theta} \mathbb{E}_{s, a, s' \sim \xi_{\pi}} \mathbb{E}_{\bar{s}, \bar{a}, \bar{s}' \sim \phi_{\iota}(\cdot | s, a, s')} \|\langle s, a, s' \rangle - \langle \mathcal{G}_{\theta}(\bar{s}), \psi_{\theta}(\bar{s}, \bar{a}), \mathcal{G}_{\theta}(\bar{s}') \rangle\| + L_{\mathcal{R}}^{\xi\pi} + \beta \left( \mathcal{W}_{\xi\pi} + L_{\mathbf{P}}^{\xi\pi} \right)$$



$$\min_{\iota, \theta} \mathbb{E}_{s, a, s' \sim \xi_{\pi}} \mathbb{E}_{\bar{s}, \bar{a}, \bar{s}' \sim \phi_{\iota}(\cdot | s, a, s')} \|\langle s, a, s' \rangle - \langle \mathcal{G}_{\theta}(\bar{s}), \psi_{\theta}(\bar{s}, \bar{a}), \mathcal{G}_{\theta}(\bar{s}') \rangle\| + L_{\mathcal{R}}^{\xi\pi} + \beta \left( \mathcal{W}_{\xi\pi} + L_{\mathbf{P}}^{\xi\pi} \right)$$

# Wasserstein Auto-encoded Markov Decision Process

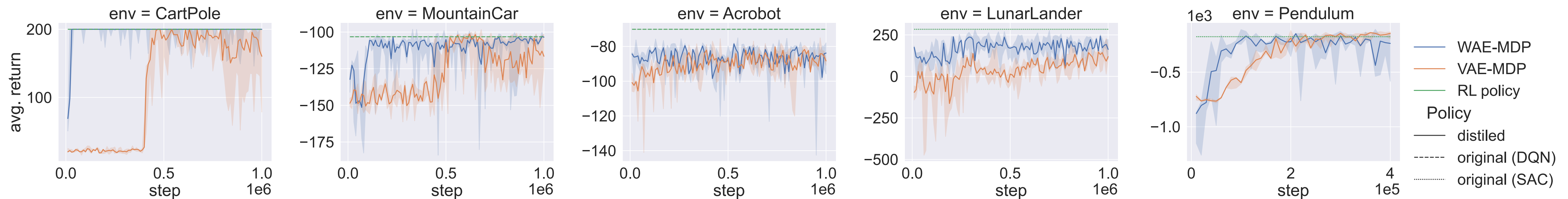
$$\min_{\iota, \theta} \mathbb{E}_{s, a, s' \sim \xi_{\pi}} \mathbb{E}_{\bar{s}, \bar{a}, \bar{s}' \sim \phi_{\iota}(\cdot | s, a, s')} \|\langle s, a, s' \rangle - \langle \mathcal{G}_{\theta}(\bar{s}), \psi_{\theta}(\bar{s}, \bar{a}), \mathcal{G}_{\theta}(\bar{s}') \rangle\| + L_{\mathcal{R}}^{\xi_{\pi}} + \beta \left( \mathcal{W}_{\xi_{\pi}} + L_{\mathbf{P}}^{\xi_{\pi}} \right)$$



- $\mathcal{W}_{\xi_{\pi}} = \max_{\omega: \|\varphi_{\omega}^{\xi}\| \leq 1} \mathbb{E}_{s, a \sim \xi_{\pi}} \mathbb{E}_{\bar{a} \sim \phi_{\iota}^{\mathcal{A}}(\cdot | \phi_{\iota}(s), a)} \mathbb{E}_{\bar{s}' \sim \bar{\mathbf{P}}_{\theta}(\cdot | \bar{s}, \bar{a})} \varphi_{\omega}^{\xi}(\phi_{\iota}(s), \bar{a}, \bar{s}') - \mathbb{E}_{\bar{s}, \bar{a}, \bar{s}' \sim \xi_{\pi}} \varphi_{\omega}^{\xi}(\bar{s}, \bar{a}, \bar{s}')$

- $L_{\mathbf{P}}^{\xi_{\pi}} = \max_{\omega: \|\varphi_{\omega}^{\mathbf{P}}\| \leq 1} \mathbb{E}_{s, a, s' \sim \xi_{\pi}} \mathbb{E}_{\bar{s}, \bar{a}, \bar{s}' \sim \phi_{\iota}(\cdot | s, a, s')} \left[ \varphi_{\omega}^{\mathbf{P}}(s, a, \bar{s}, \bar{a}, \bar{s}') - \mathbb{E}_{\bar{s}' \sim \bar{\mathbf{P}}_{\theta}(\cdot | \bar{s}, \bar{a})} \varphi_{\omega}^{\mathbf{P}}(s, a, \bar{s}, \bar{a}, \bar{s}') \right]$

## Distillation: performance of $\bar{\pi}$



**WAE-MDPs** distill policies up to 10 times faster than **VAE-MDPs**

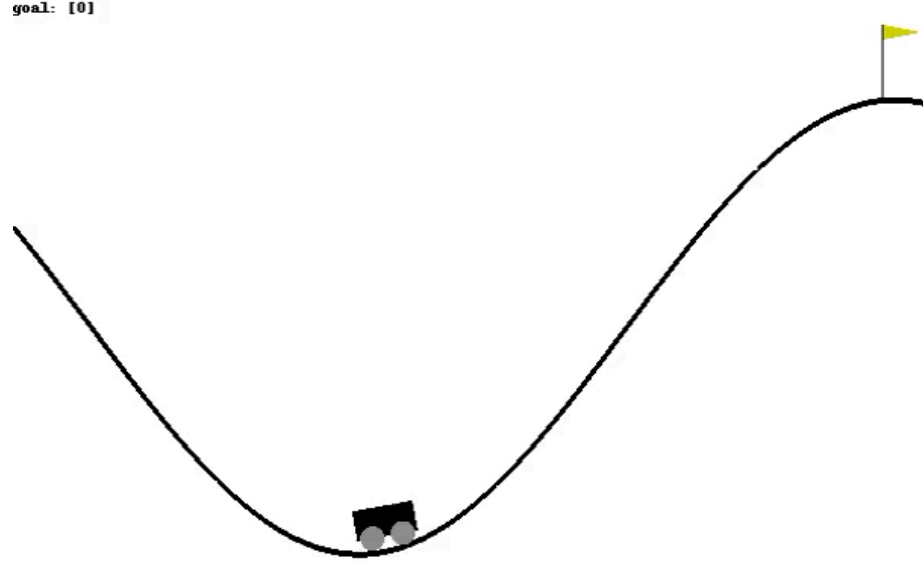
- *Faster*
- *Better performance*
- *Learning guarantees*
- *Similar or even better model quality*

# Evaluation

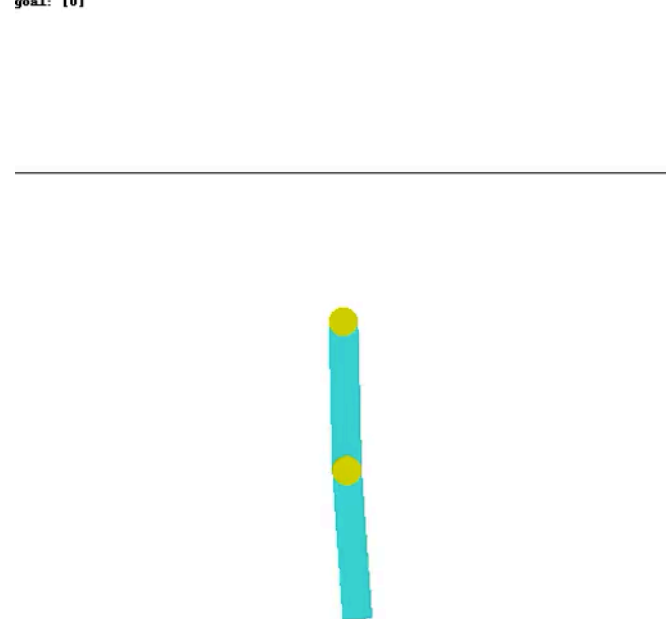
CartPole



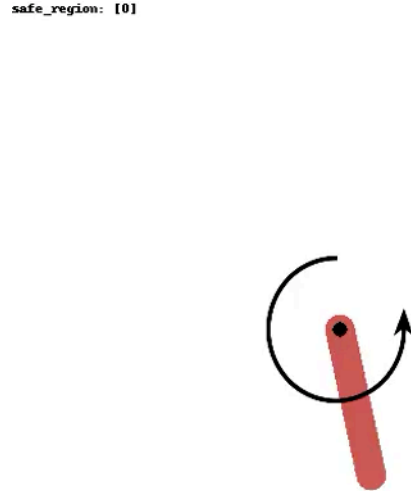
MountainCar



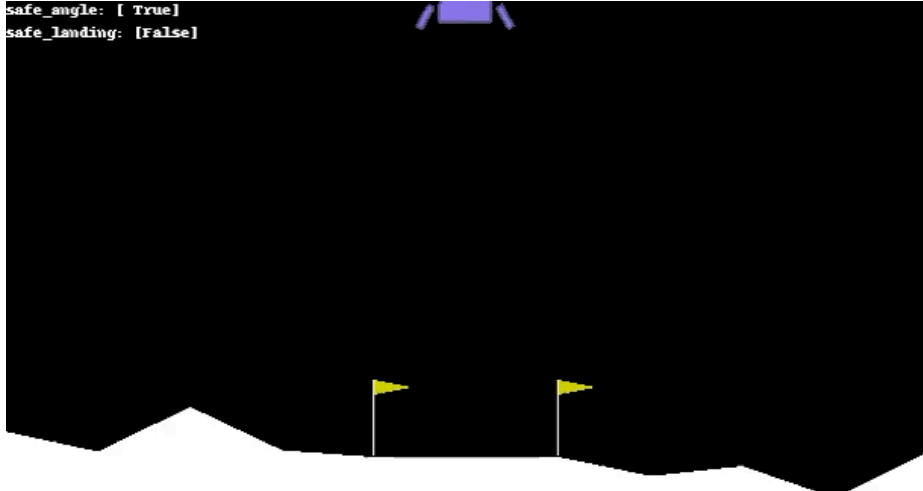
Acrobot



Pendulum



LunarLander



Time-to-failure properties (lower is better)

$$\varphi = \neg \text{Reset} \mathcal{U} \neg \text{Safe}$$

$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0.032$$

$$\varphi = \neg \text{Goal} \mathcal{U} \text{Reset}$$

$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0$$

$$\varphi = \neg \text{Goal} \mathcal{U} \text{Reset}$$

$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0.0022$$

$$\varphi = \diamond(\neg \text{Safe} \wedge \bigcirc \text{Reset})$$

$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0.037$$

$$\varphi = \neg \text{SafeLanding} \mathcal{U} \text{Reset}$$

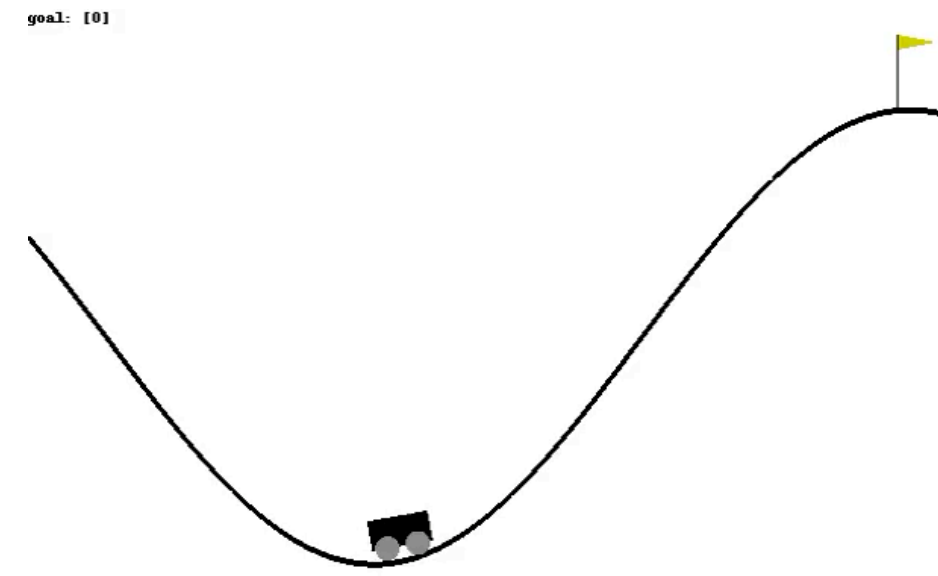
$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0.0702$$

# Evaluation

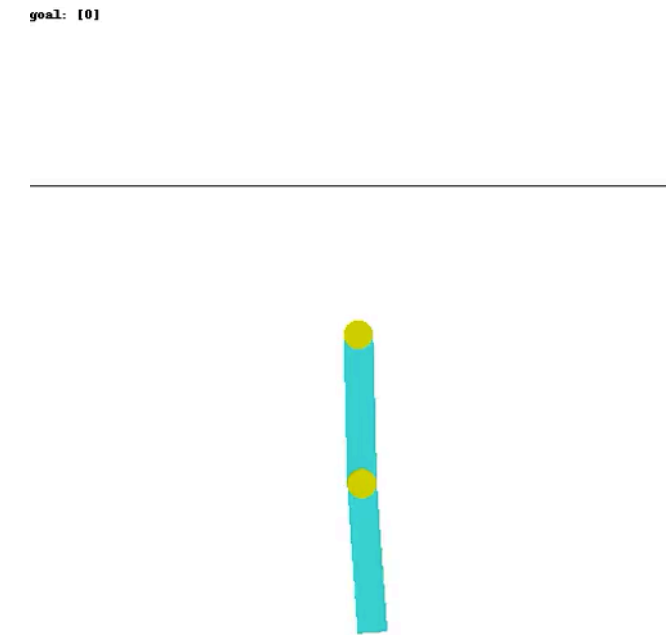
CartPole



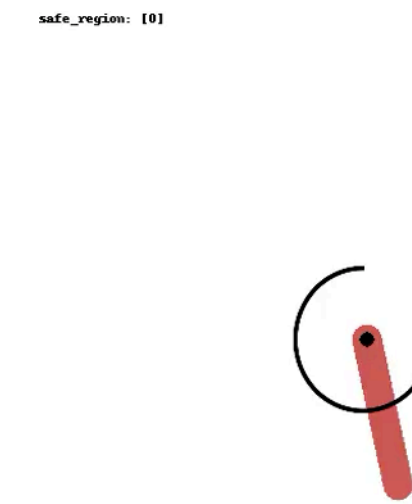
MountainCar



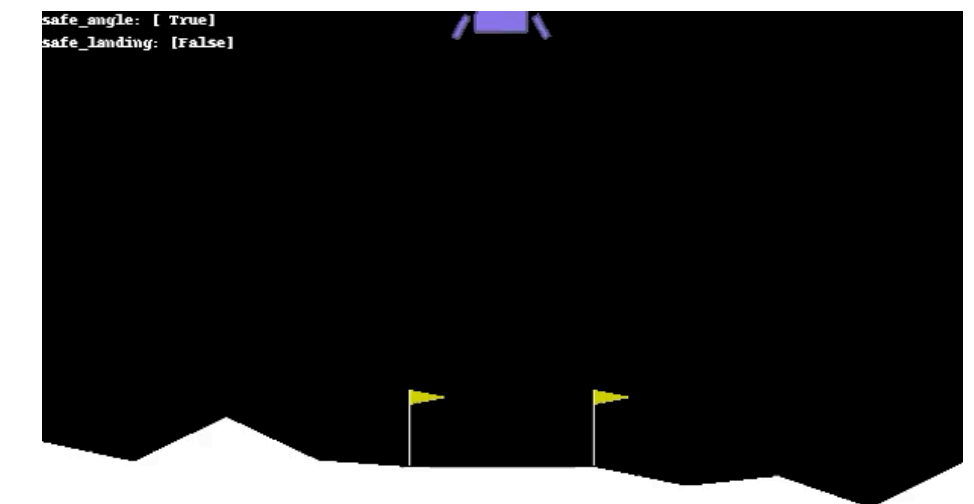
Acrobot



Pendulum



LunarLander



Time-to-failure properties (lower is better)

$$\varphi = \neg \text{Reset} \mathcal{U} \neg \text{Safe}$$

$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0.032$$

$$\varphi = \neg \text{Goal} \mathcal{U} \text{Reset}$$

$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0$$

$$\varphi = \neg \text{Goal} \mathcal{U} \text{Reset}$$

$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0.0022$$

$$\varphi = \diamond(\neg \text{Safe} \wedge \bigcirc \text{Reset})$$

$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0.037$$

$$\varphi = \neg \text{SafeLanding} \mathcal{U} \text{Reset}$$

$$\bar{V}_{\bar{\pi}_\theta}^\varphi(\bar{s}_I) = 0.0702$$