

Wasserstein Auto-encoded MDPs

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

Florent Delgrange, Ann Nowé, Guillermo A. Pérez



ICLR



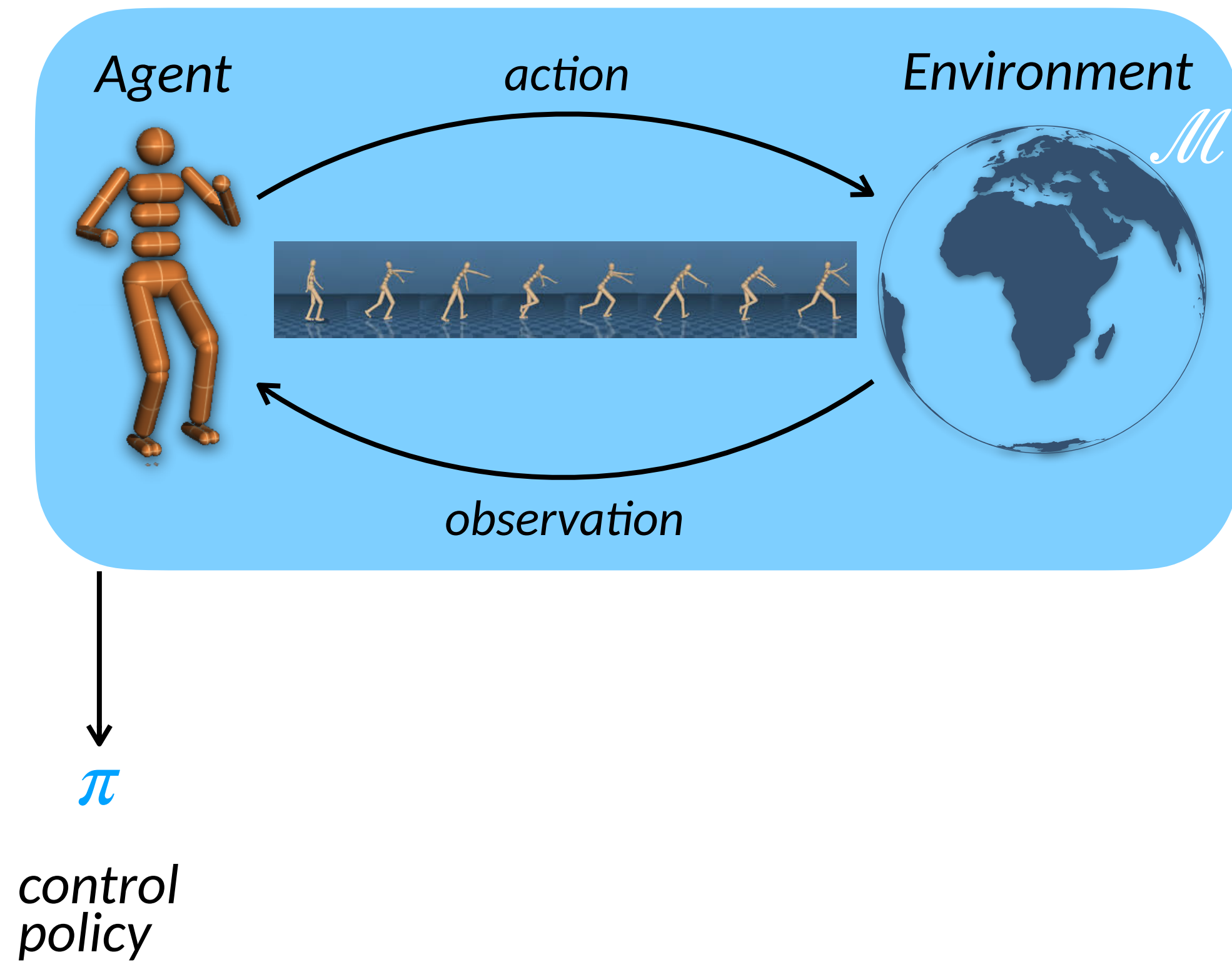
ARTIFICIAL
INTELLIGENCE
RESEARCH GROUP



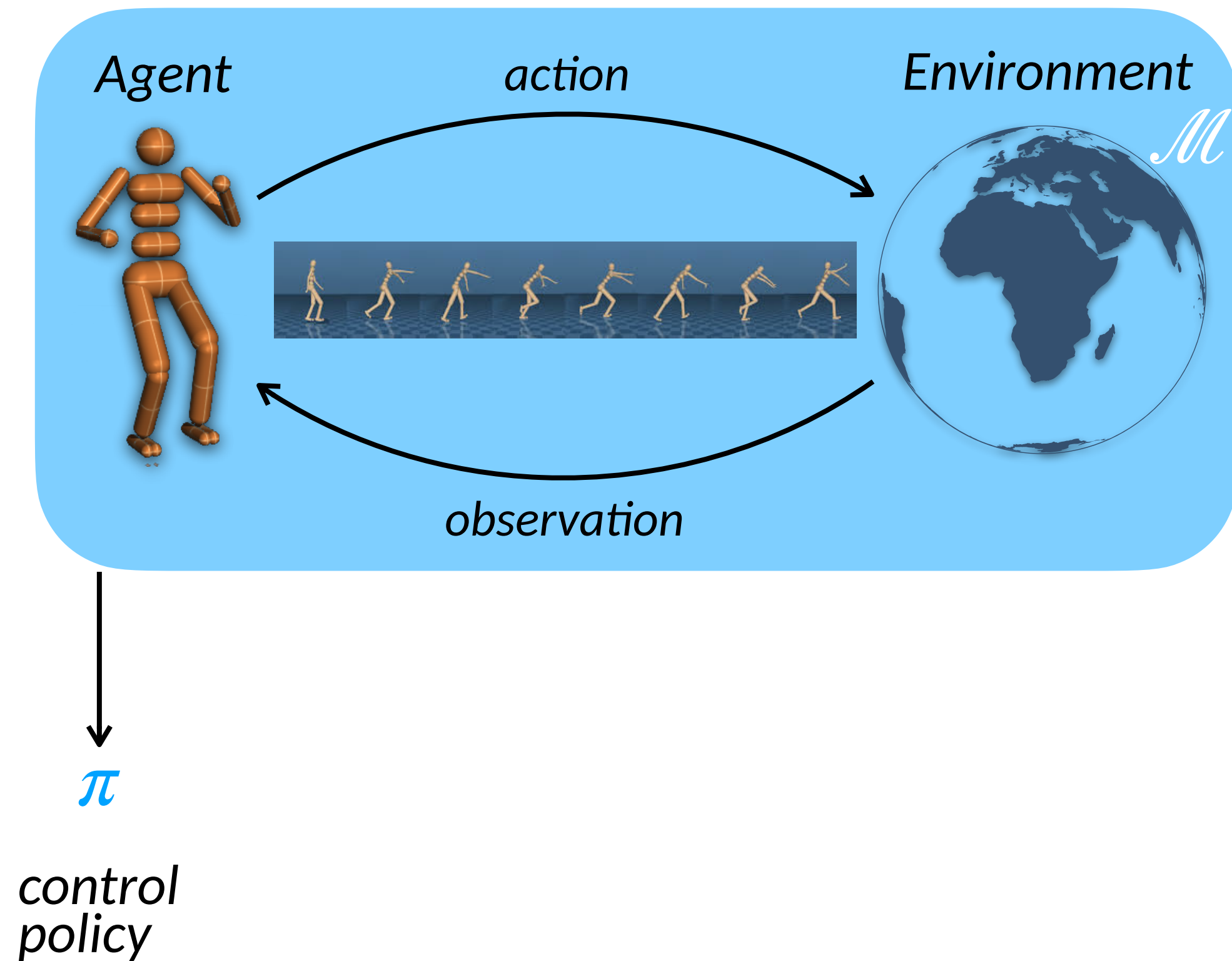
Universiteit
Antwerpen



Reinforcement Learning

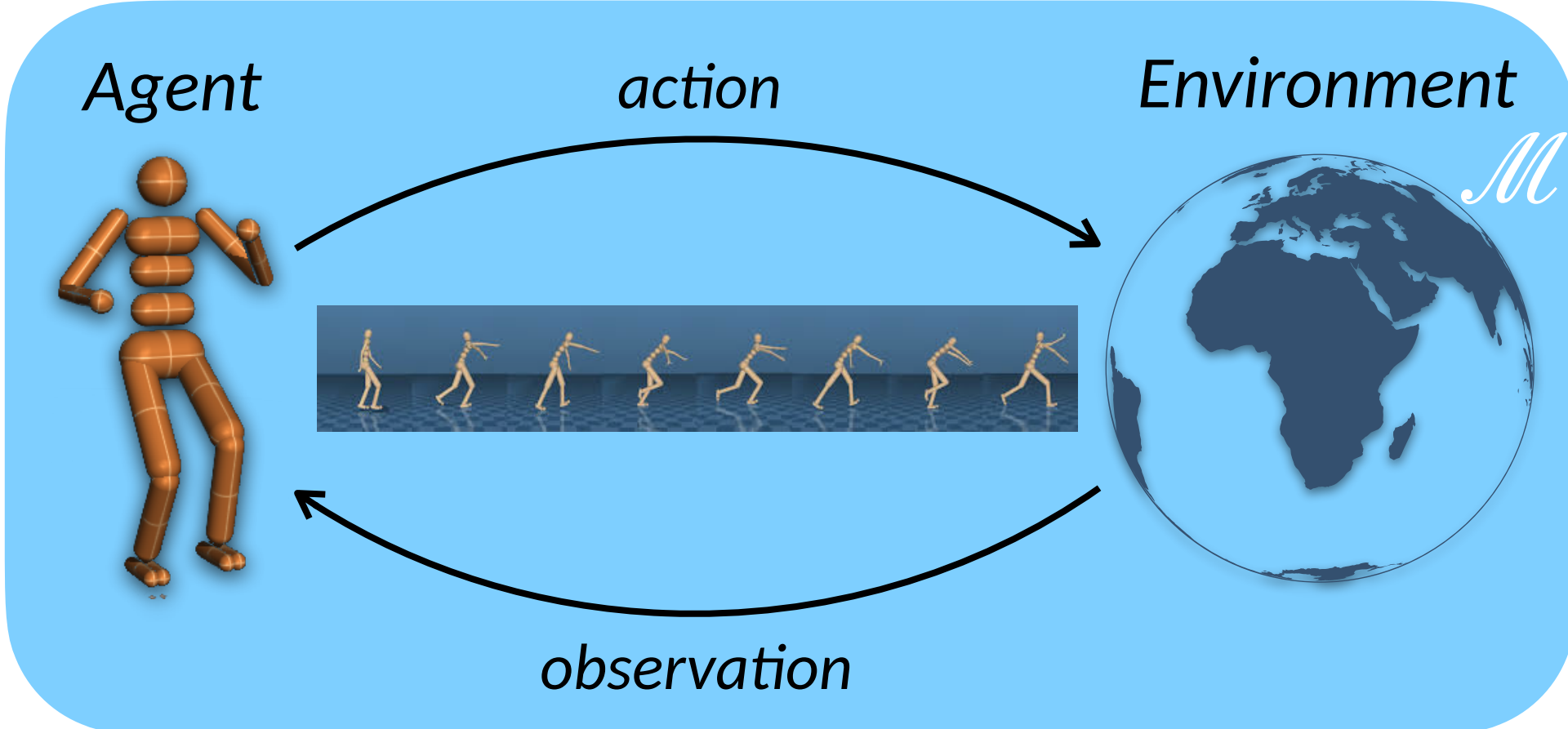


Reinforcement Learning



- Unknown environment
- Continuous state/action spaces

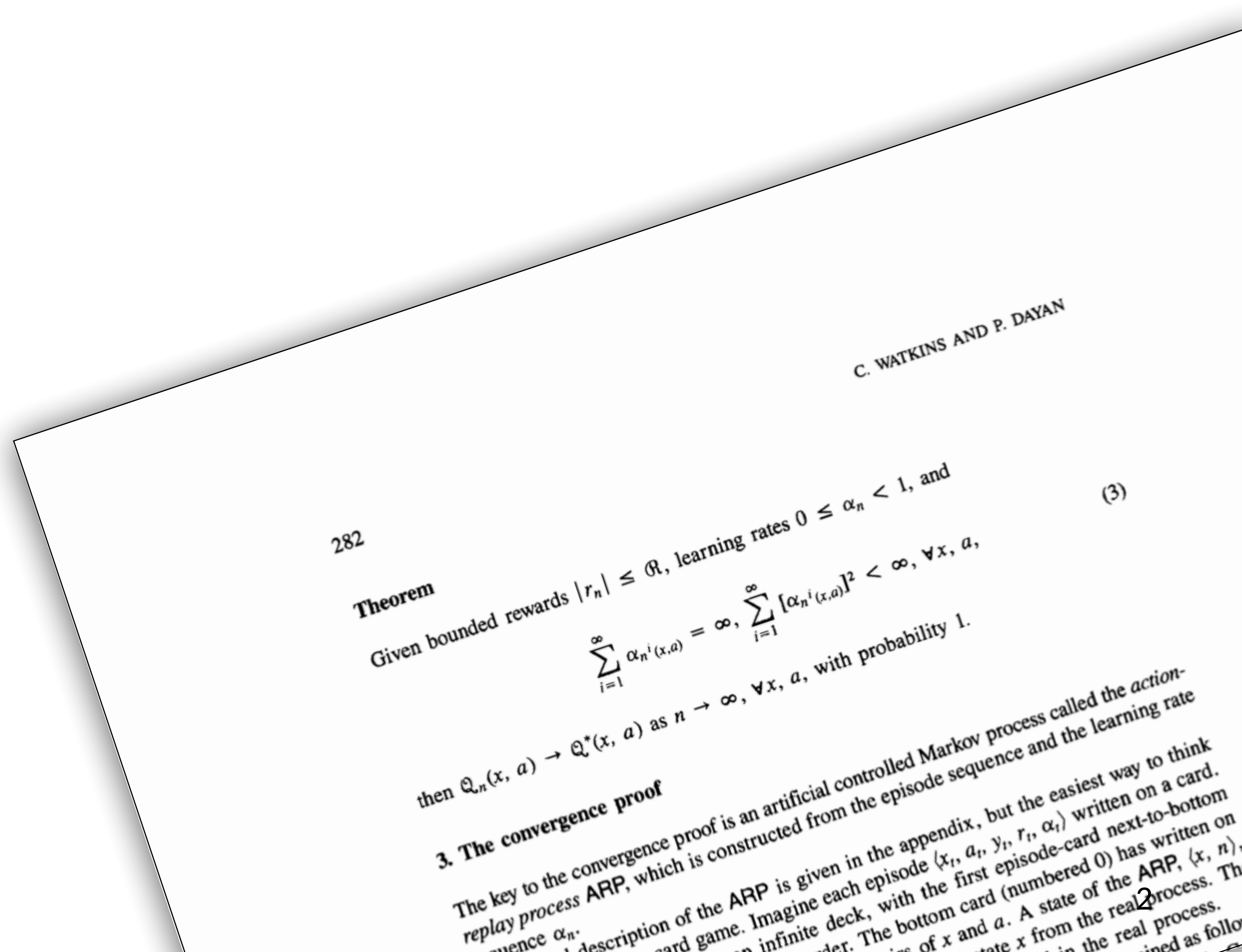
Reinforcement Learning



π

control policy

- Unknown environment
- Continuous state/action spaces



282

Theorem

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

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

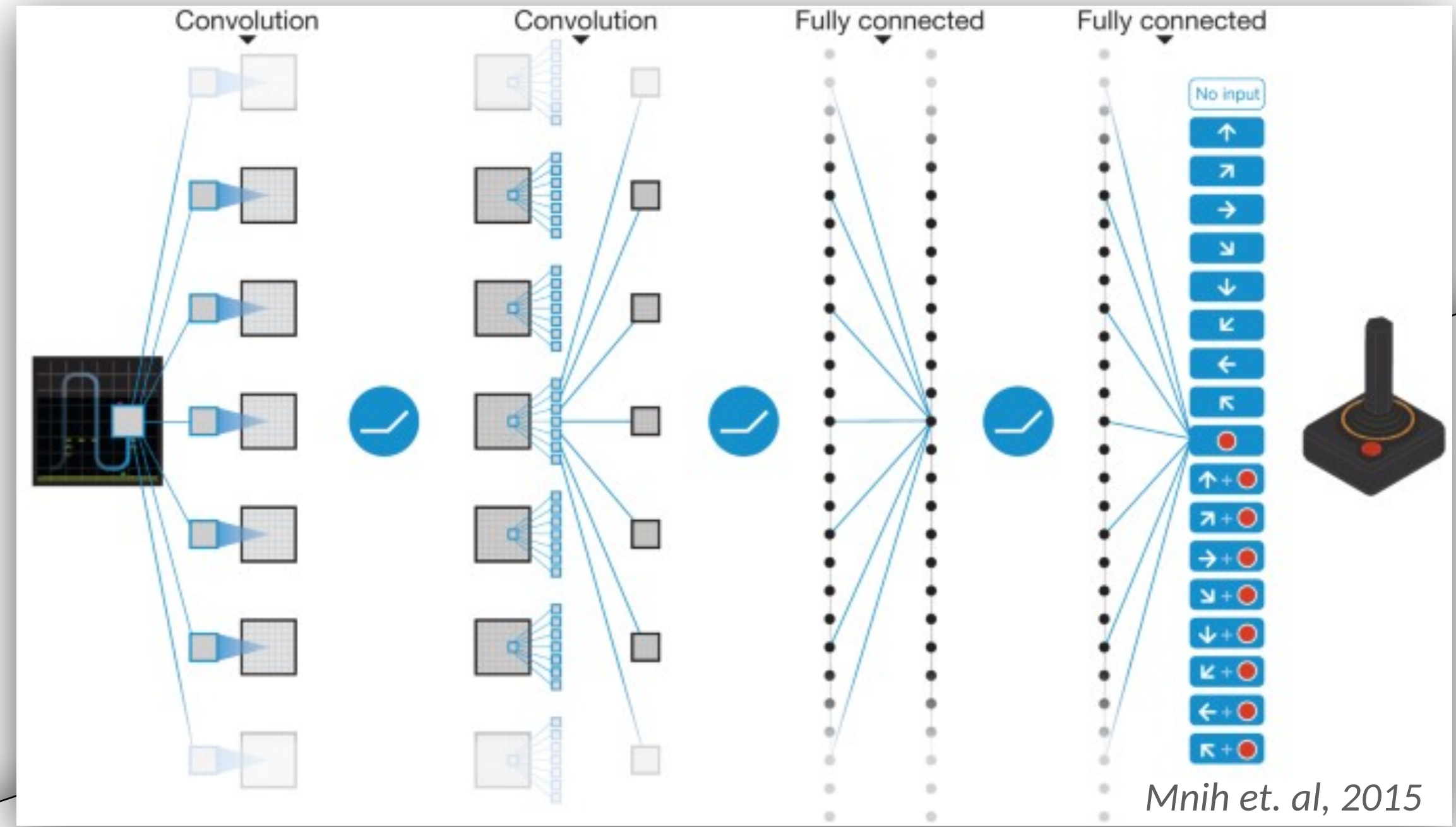
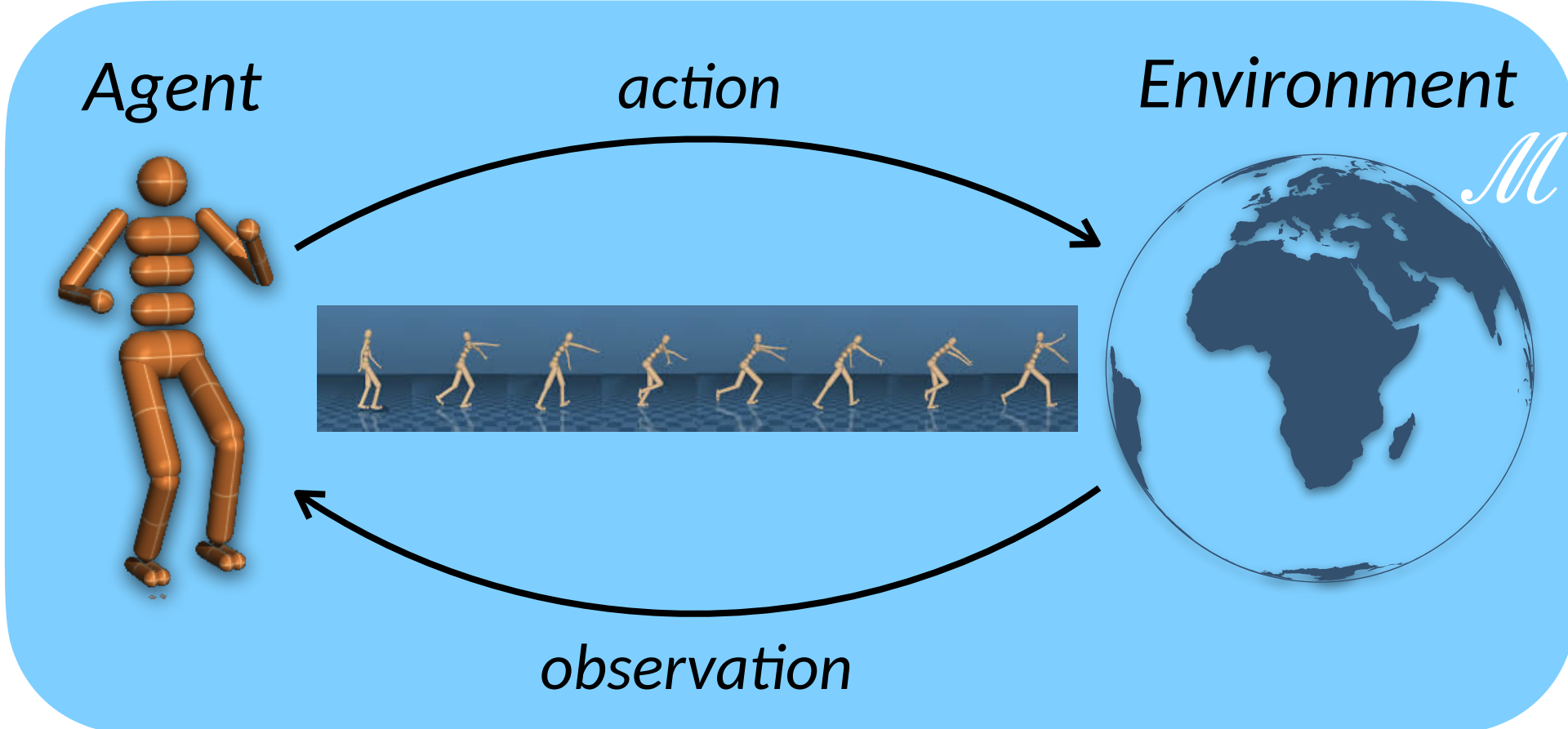
then $Q_n(x, a) \rightarrow Q^*(x, a)$ as $n \rightarrow \infty$, $\forall x, a$, with probability 1.

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 α_n . The 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 $\langle x_i, a_i, y_i, r_i, \alpha_i \rangle$ written on a card. The bottom card (numbered 0) has written on it the state x and a . A state of the ARP, $\langle x, n \rangle$, is defined as follows: the card numbered n is the top card, and the state x is the state written on the bottom card. The bottom card (numbered 0) has written on it the state x and a . A state of the ARP, $\langle x, n \rangle$, is defined as follows: the card numbered n is the top card, and the state x is the state written on the bottom card.

C. WATKINS AND P. DAYAN

Reinforcement Learning



control policy

- Unknown environment
- Continuous state/action spaces

282

Theorem

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

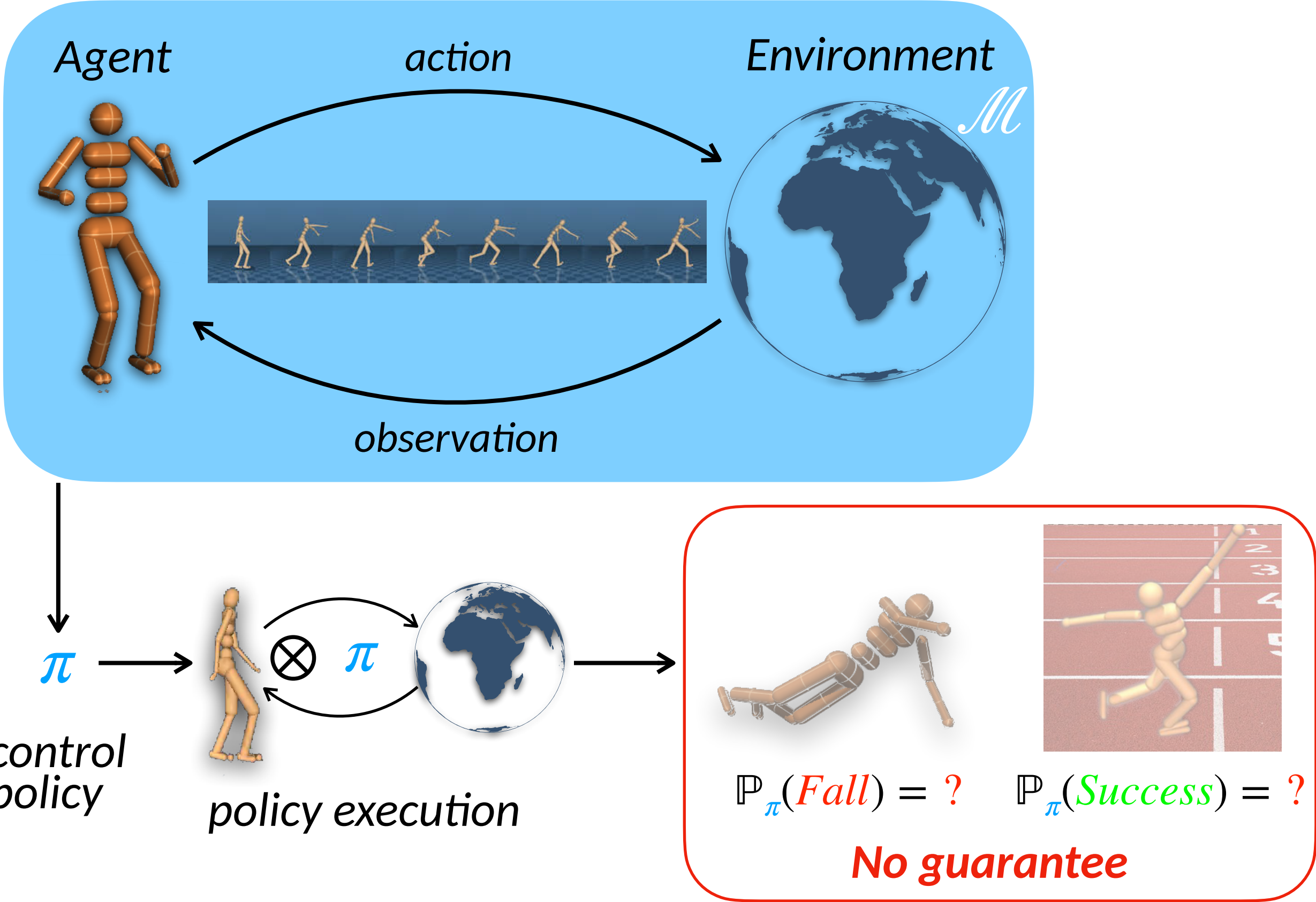
$$\sum_{i=1}^{\infty} \alpha_n^{i(x,a)} = \infty, \sum_{i=1}^{\infty} [\alpha_n^{i(x,a)}]^2 < \infty, \forall x, a,$$

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

3. The convergence proof

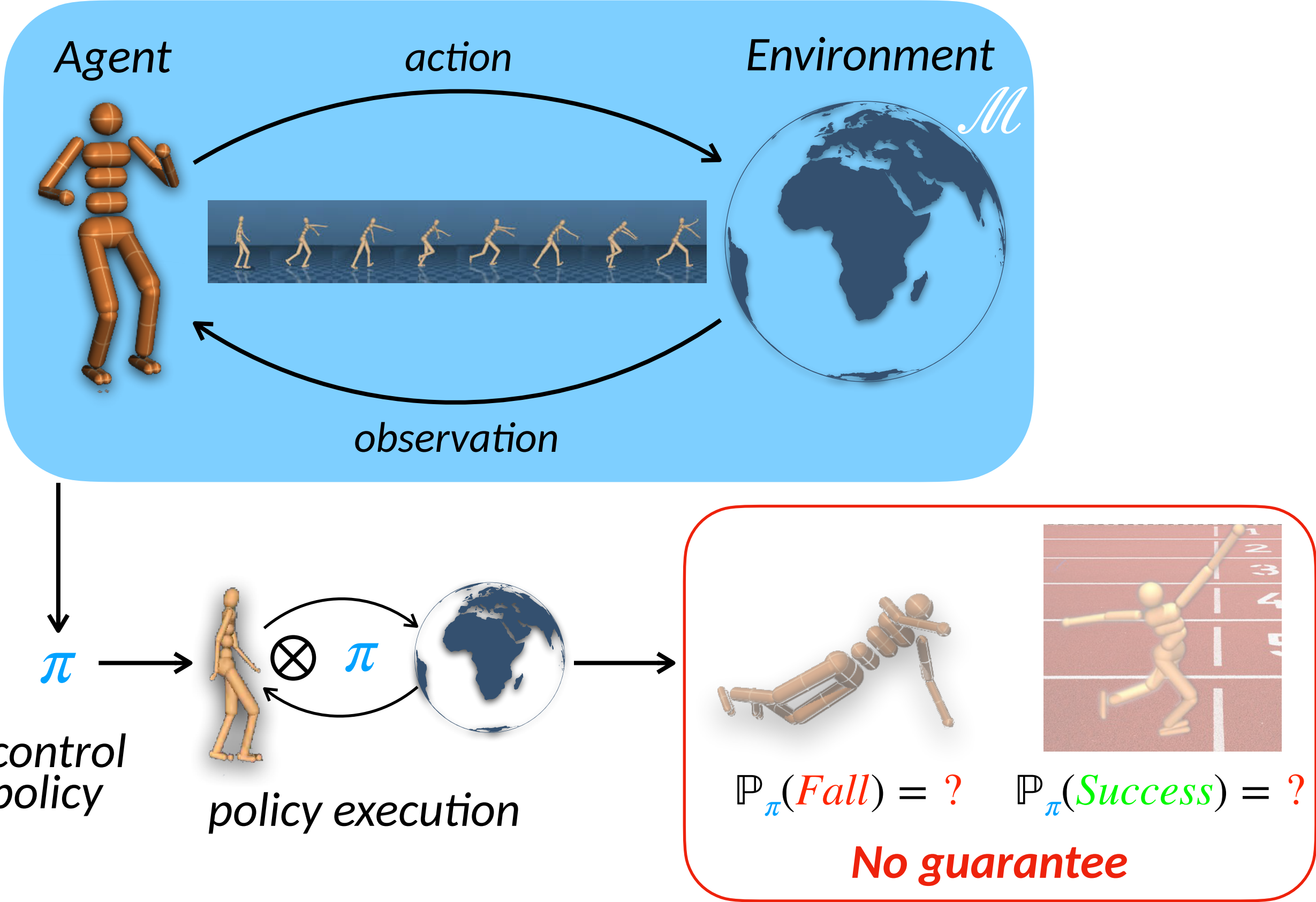
The key to the convergence proof is an artificial controlled Markov process called the action-replay process ARP, which is constructed from the episode sequence and the learning rate α_n . The description of the ARP is given in the appendix, but the easiest way to think of the ARP is as a card game. Imagine each episode $\langle x_i, a_i, y_i, r_i, \alpha_i \rangle$ written on a card. The bottom card (numbered 0) has written on it the state x and a . A state of the ARP, $\langle x, n \rangle$, is a sequence of n cards. The top card is the state x from the real process. The cards are drawn from the real process. The cards are drawn from the real process. The cards are drawn from the real process.

Reinforcement Learning



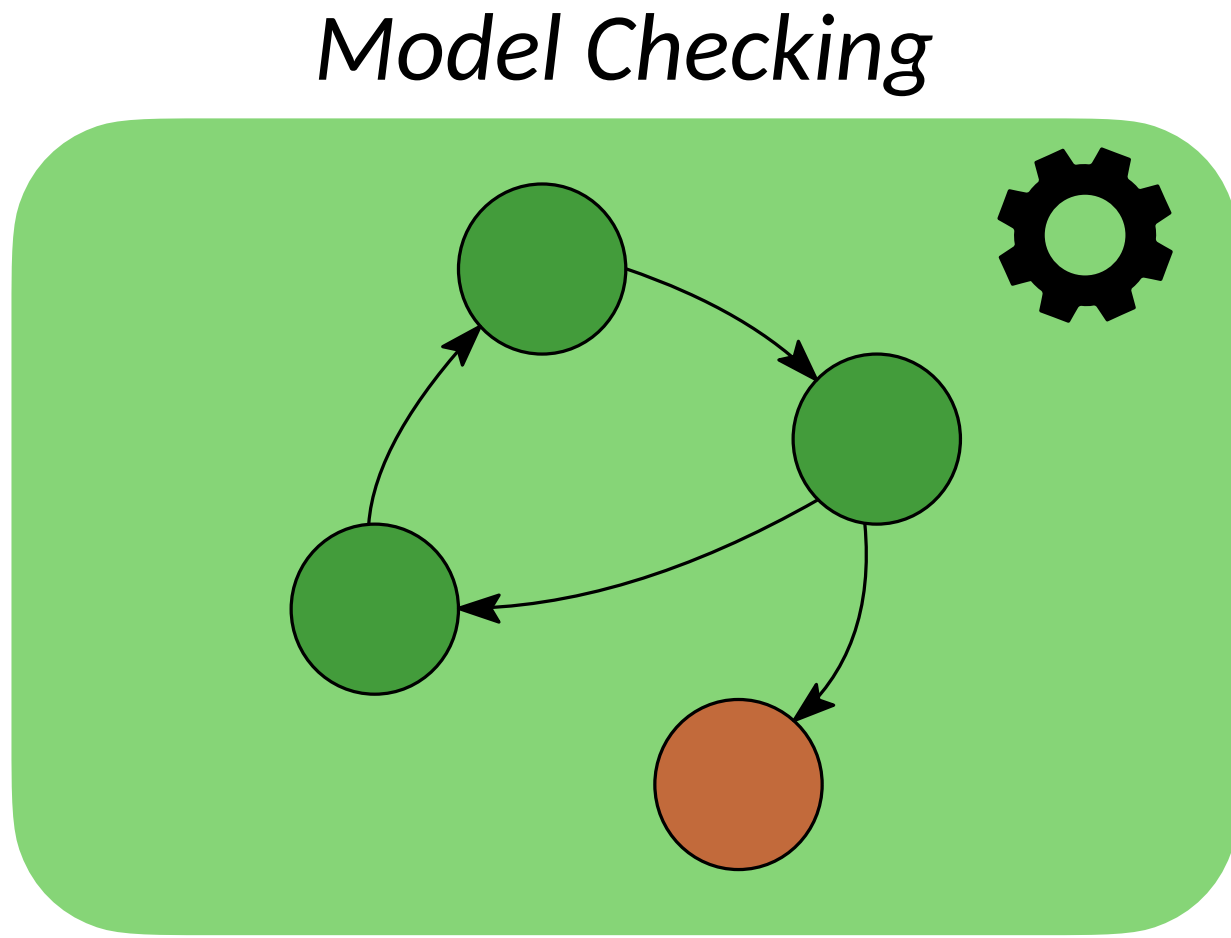
- Unknown environment
- Continuous state/action spaces

Reinforcement Learning



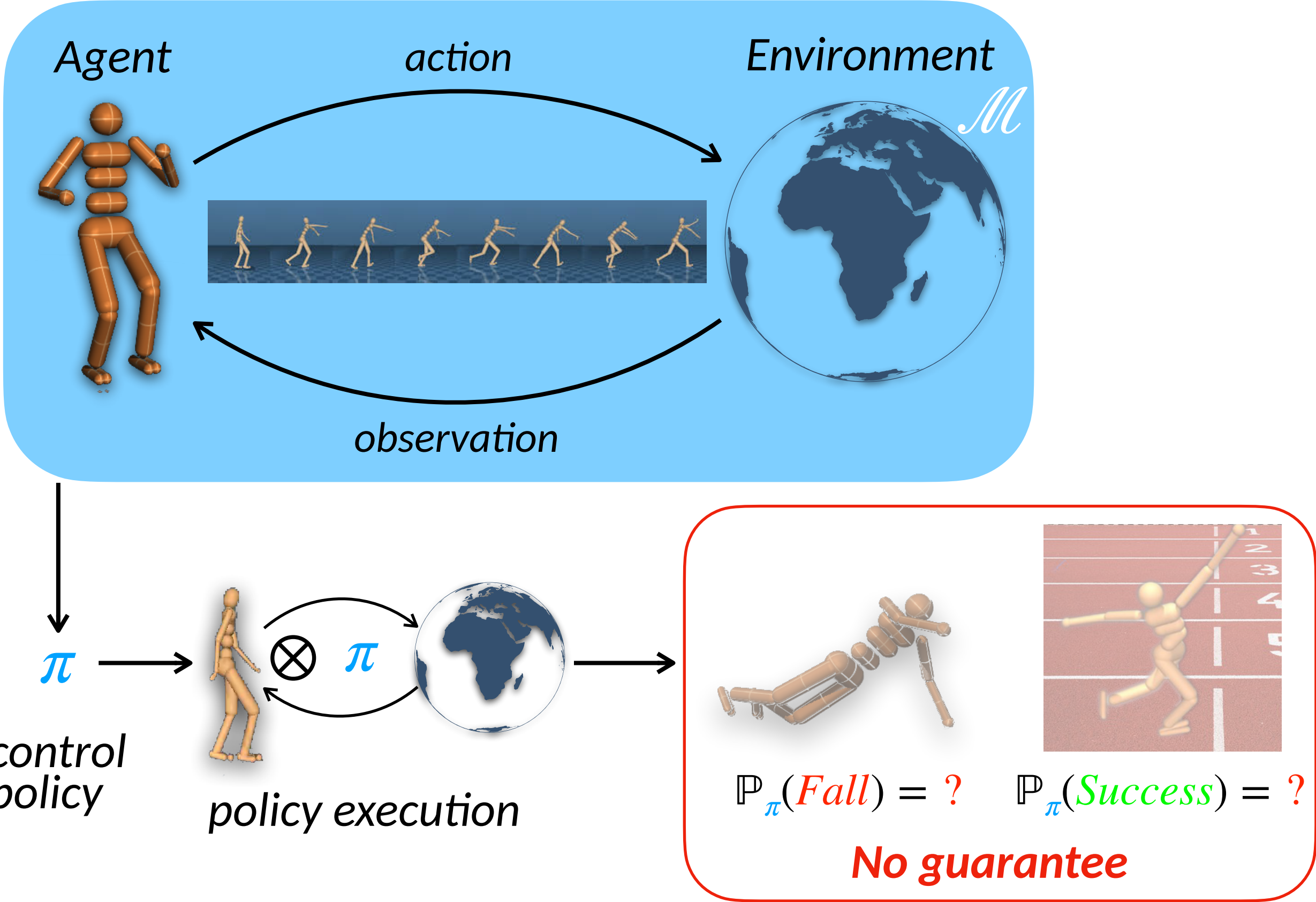
- Unknown environment
- Continuous state/action spaces

Formal Guarantees



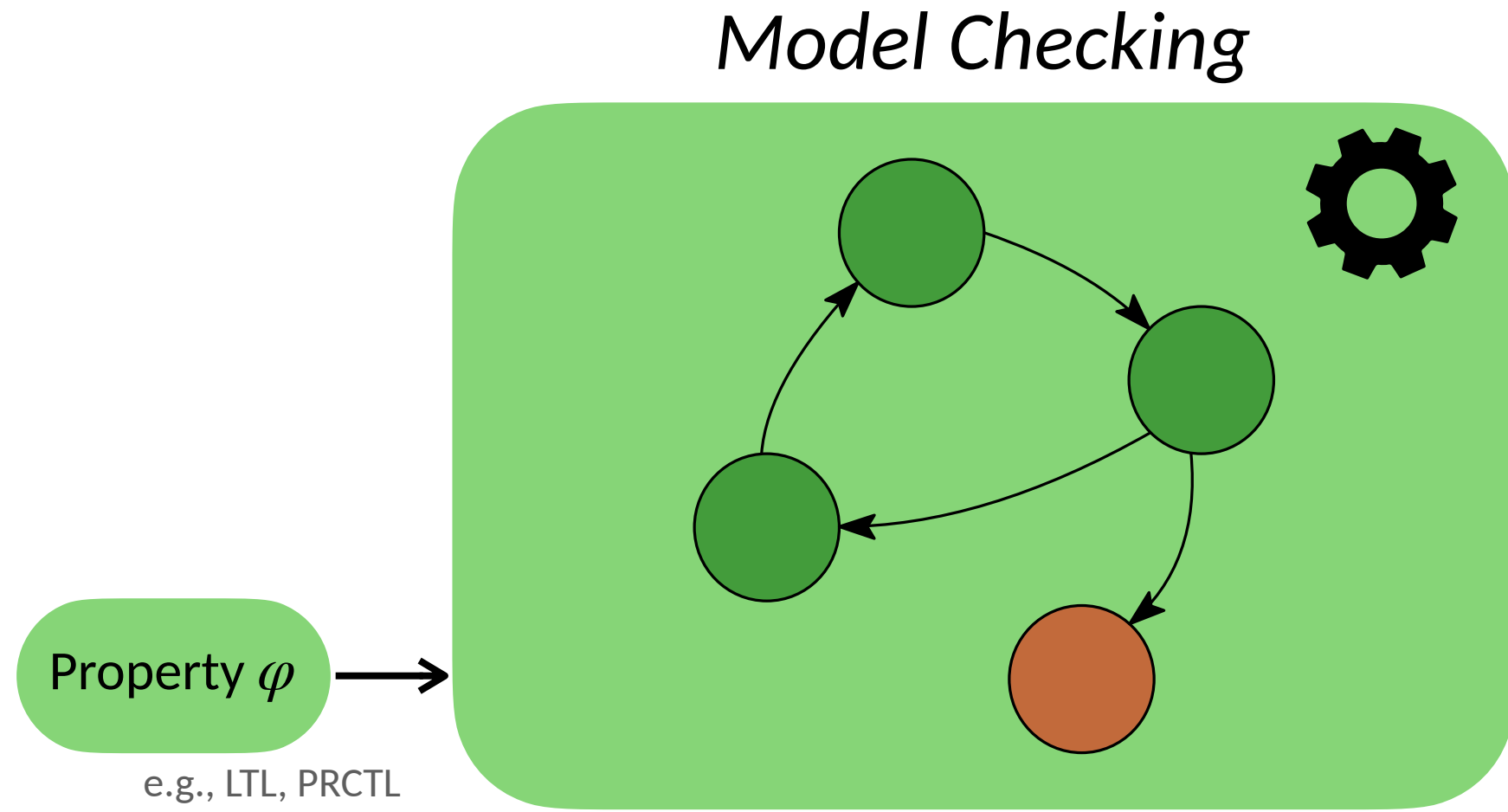
- Full knowledge of the model of the interaction

Reinforcement Learning



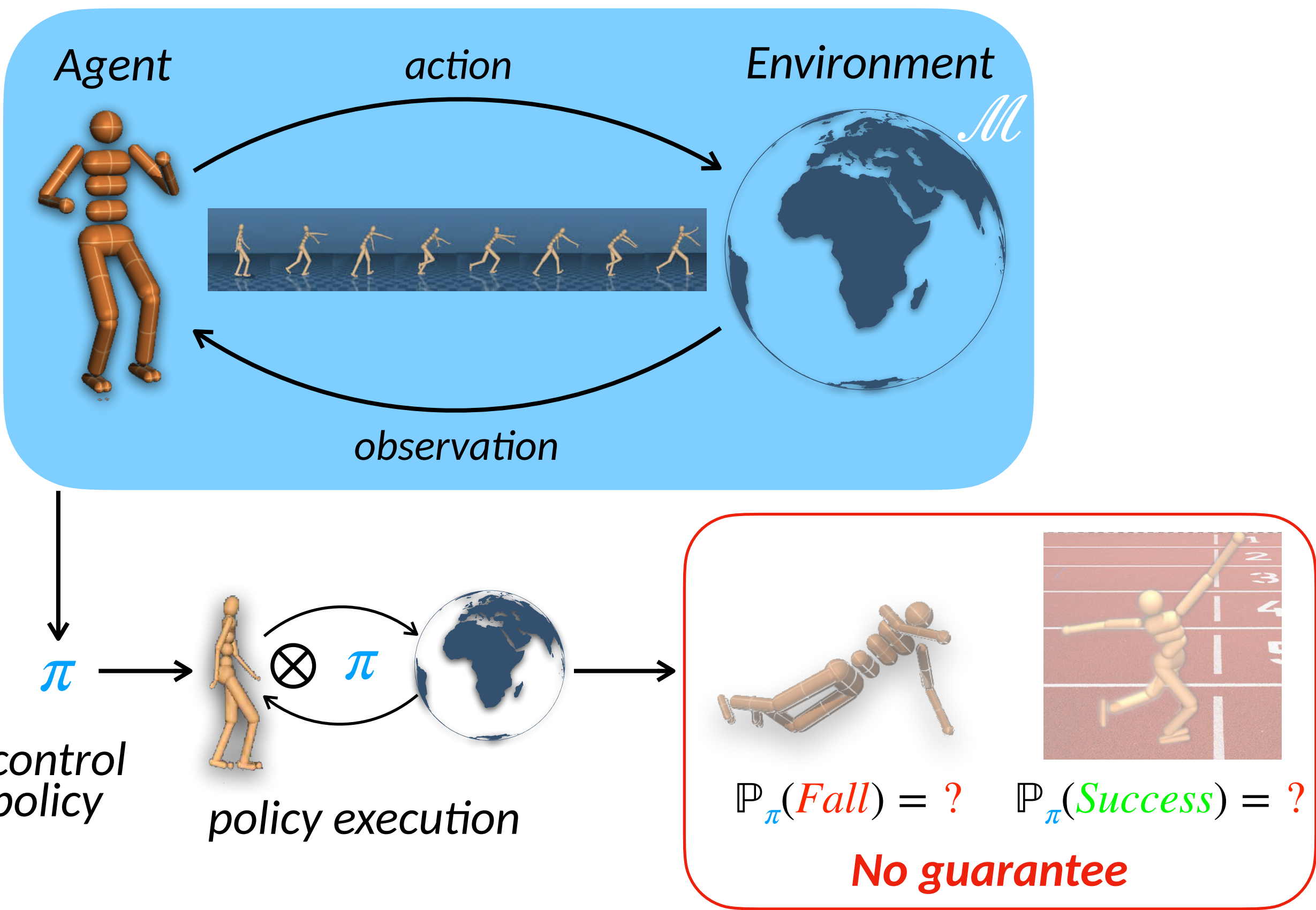
- Unknown environment
- Continuous state/action spaces

Formal Guarantees



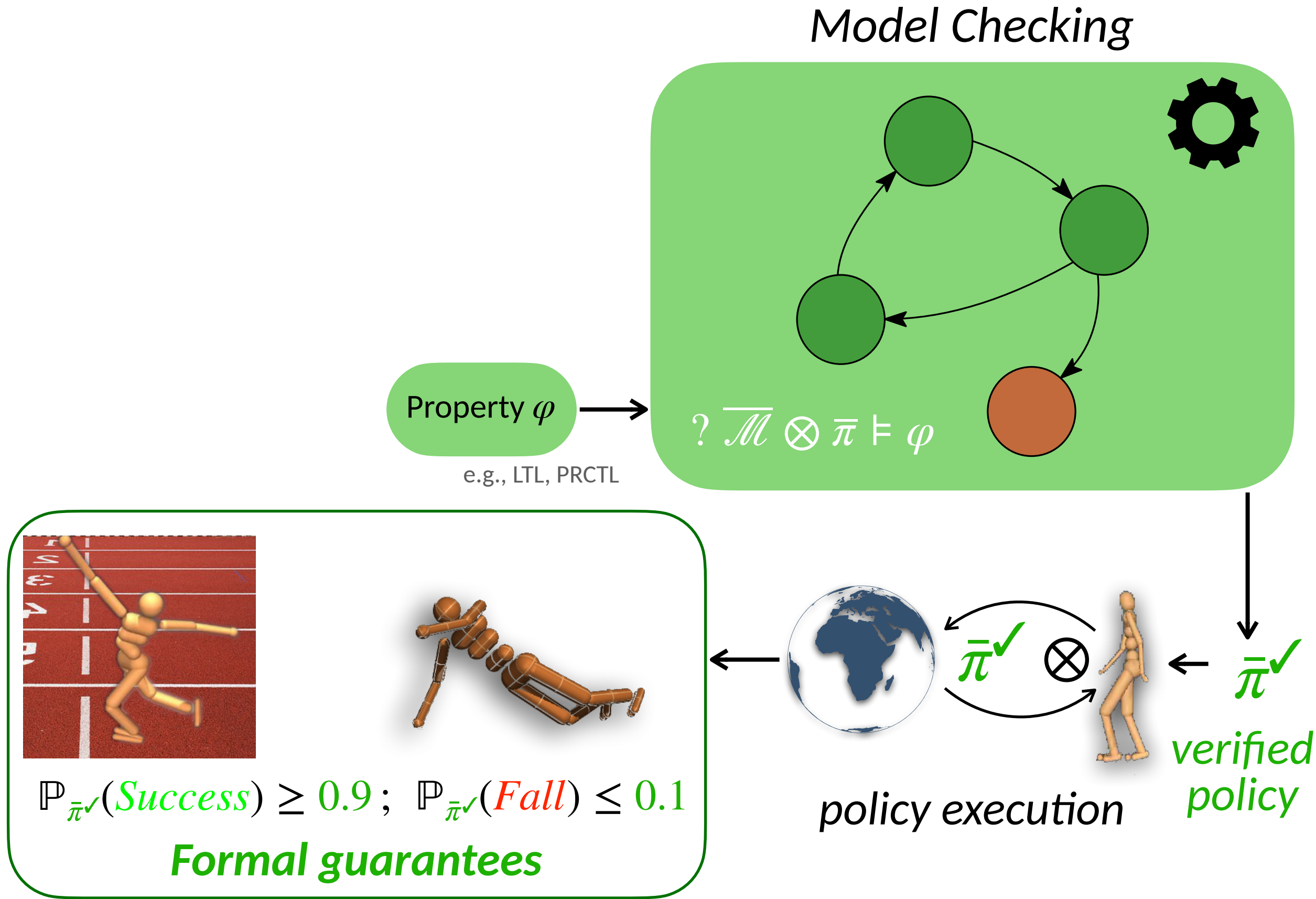
- Full knowledge of the model of the interaction

Reinforcement Learning



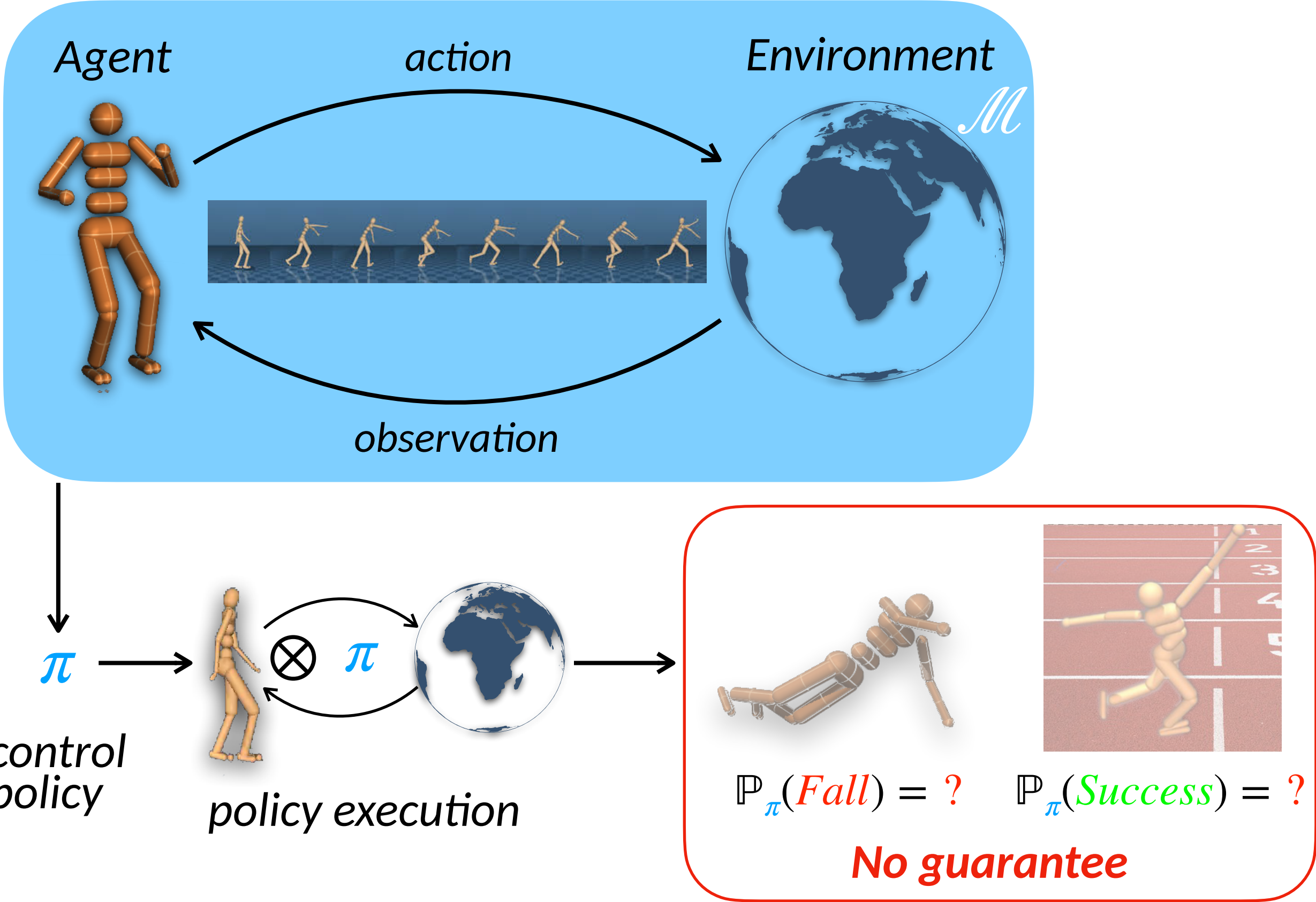
- Unknown environment
- Continuous state/action spaces

Formal Guarantees



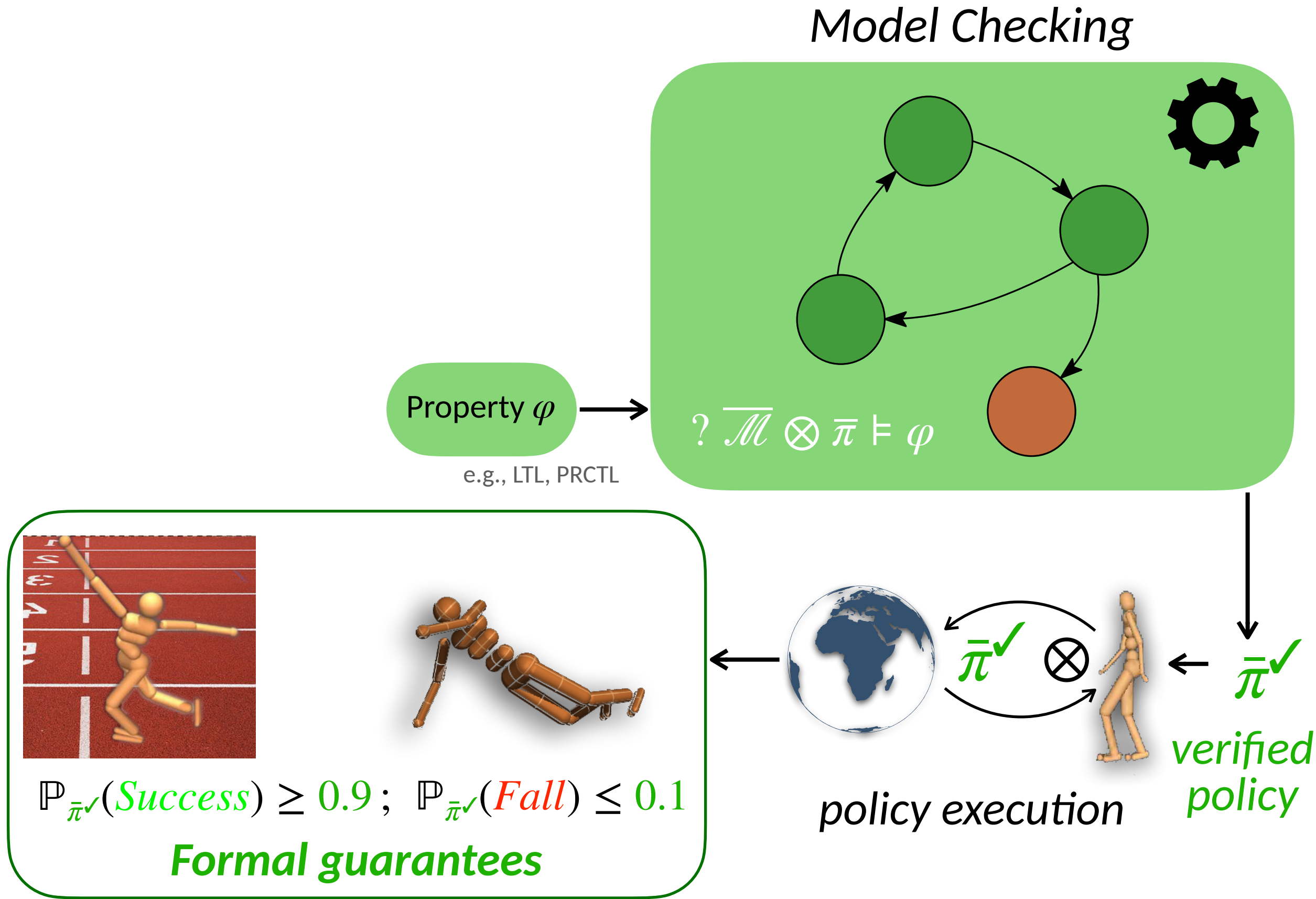
- Full knowledge of the model of the interaction

Reinforcement Learning



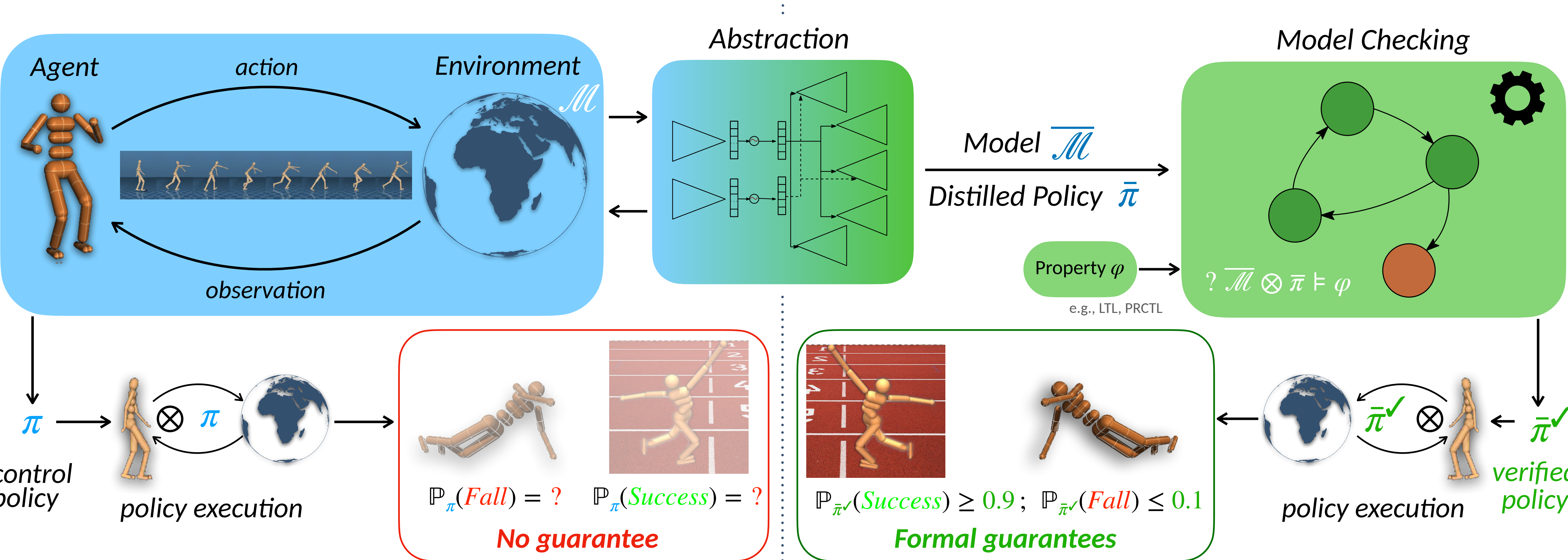
- Unknown environment
- Continuous state/action spaces

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



- Unknown environment
- Continuous state/action spaces

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

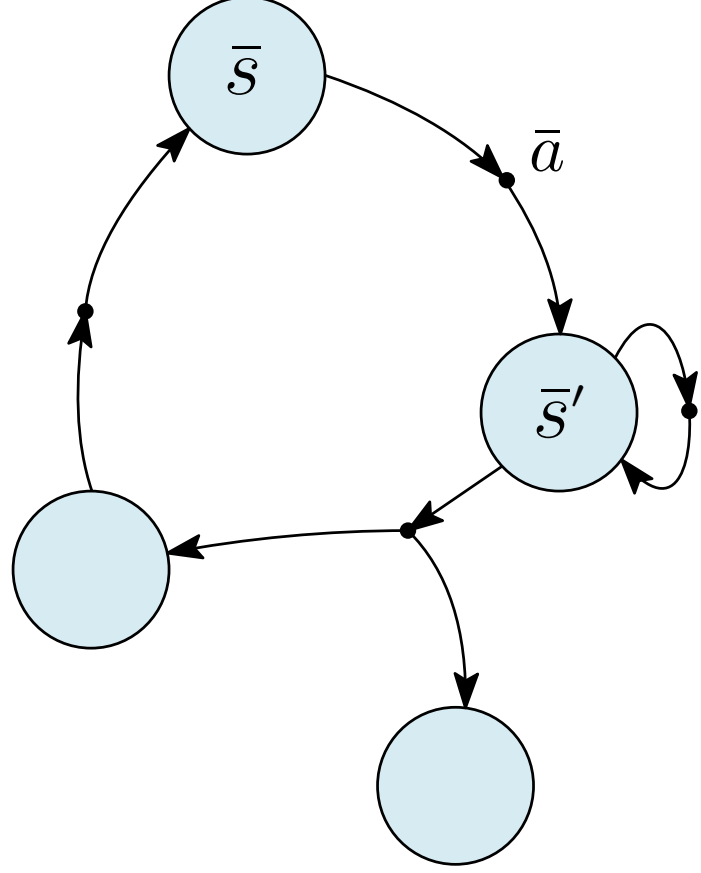
Bisimulation distance

Continuous-spaces MDP

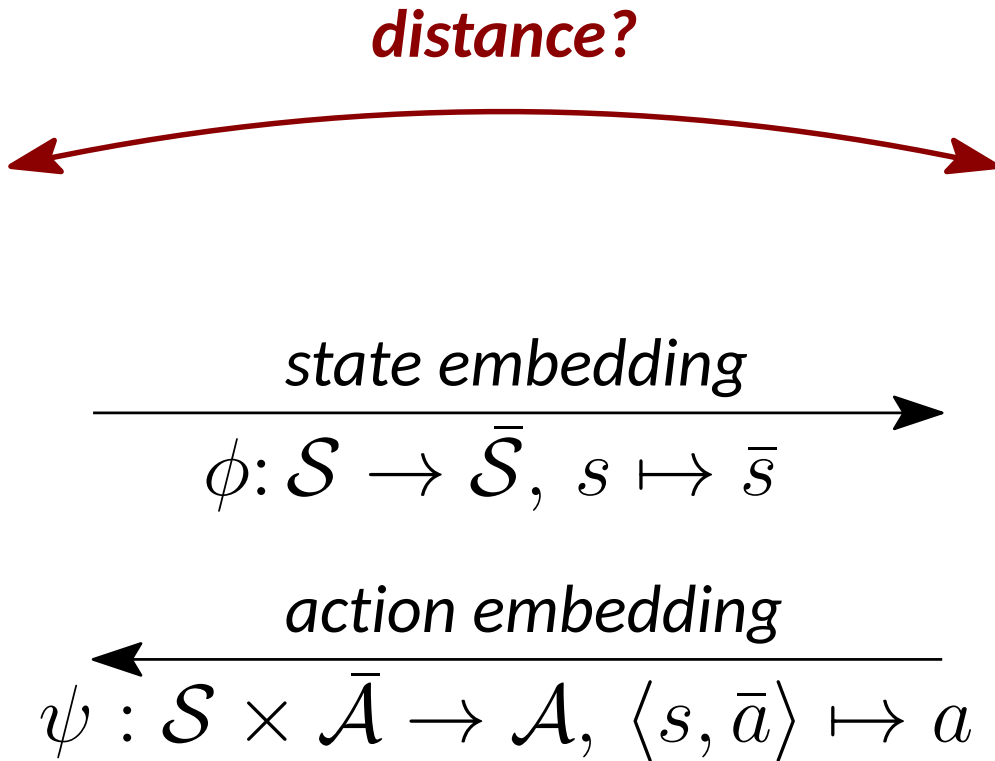


$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbf{P}, \ell \rangle$$

Discrete latent MDP



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



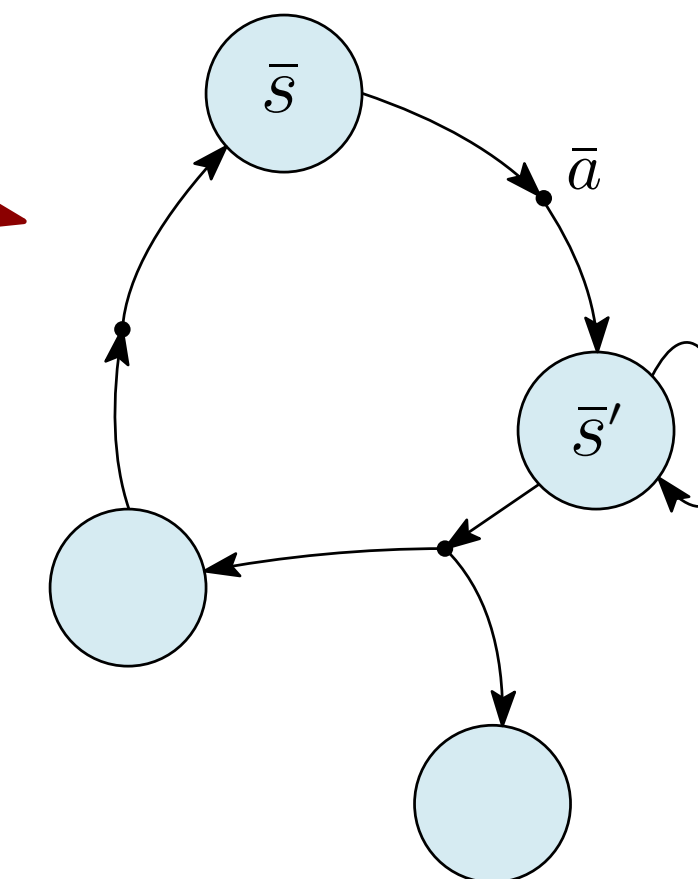
Bisimulation distance

Continuous-spaces MDP

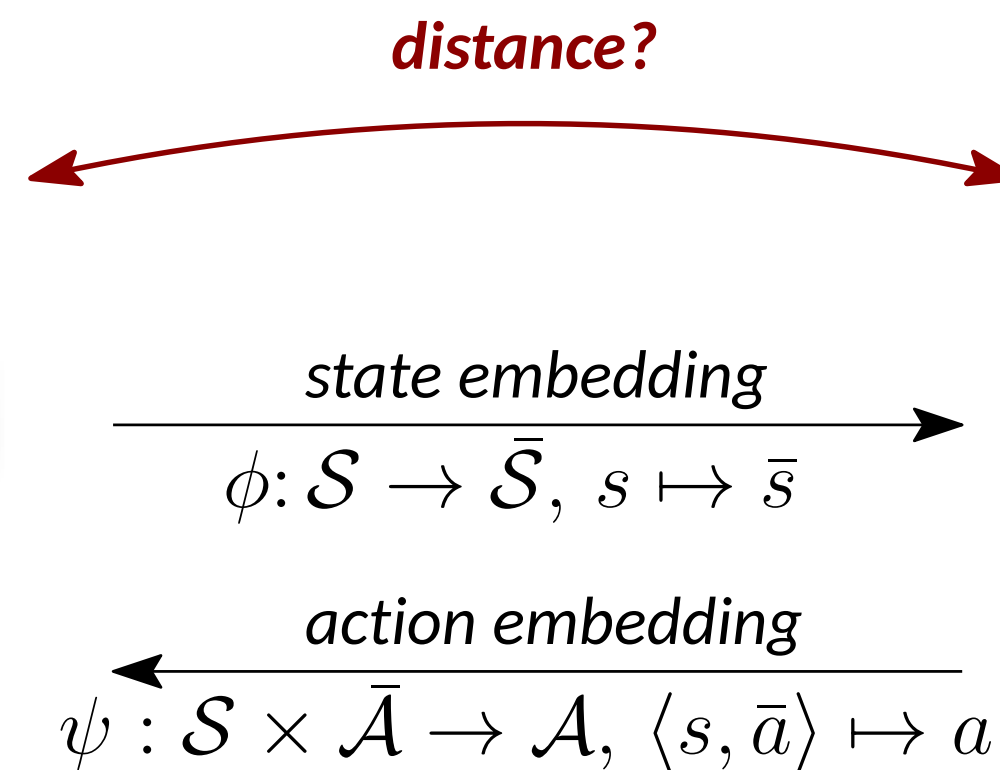


$$\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbf{P}, \ell \rangle$$

Discrete latent MDP



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



- For policy π , $\gamma \in [0, 1[$, 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}

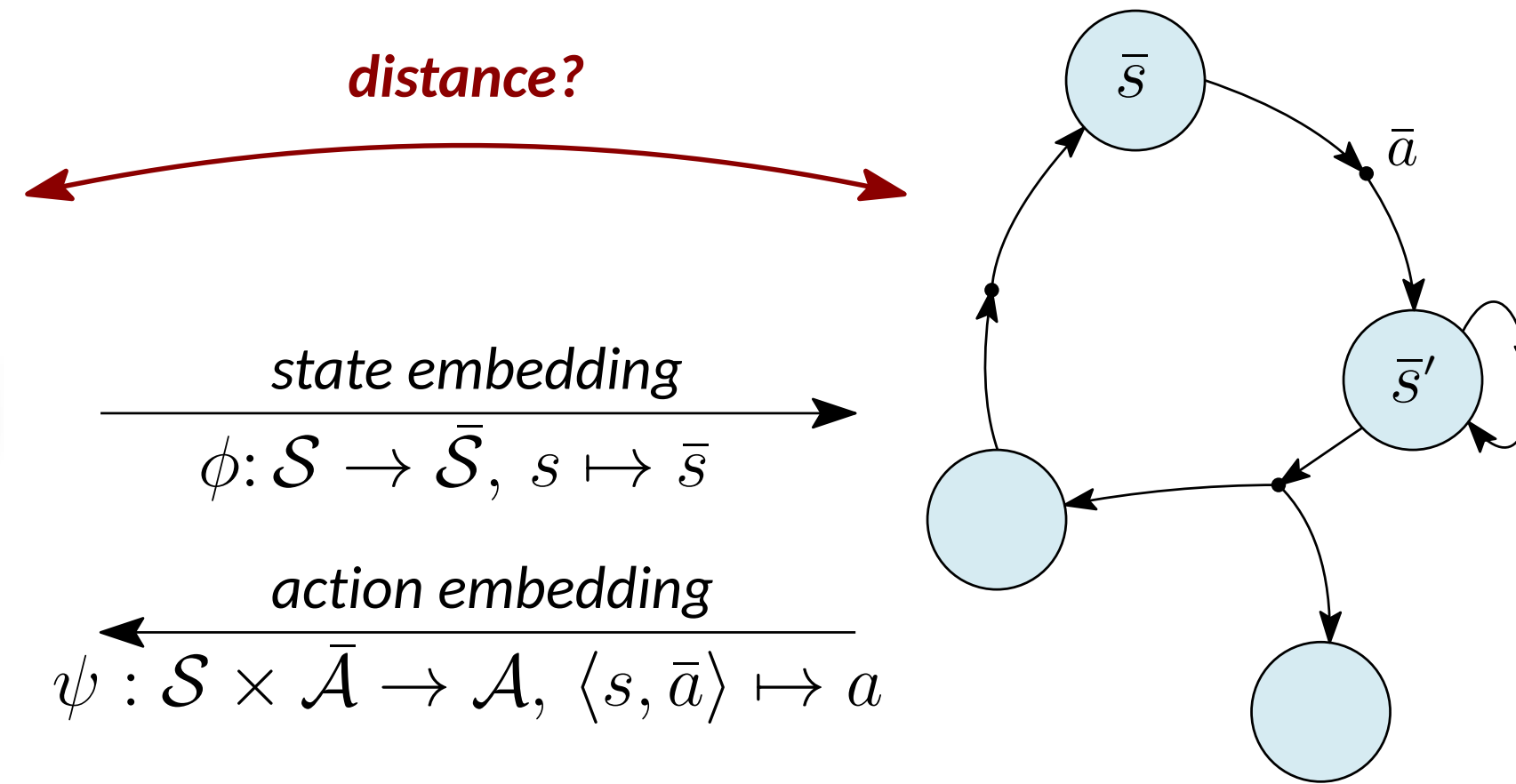
Bisimulation distance

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 π , $\gamma \in [0, 1[$, and formal logic \mathcal{L} :

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

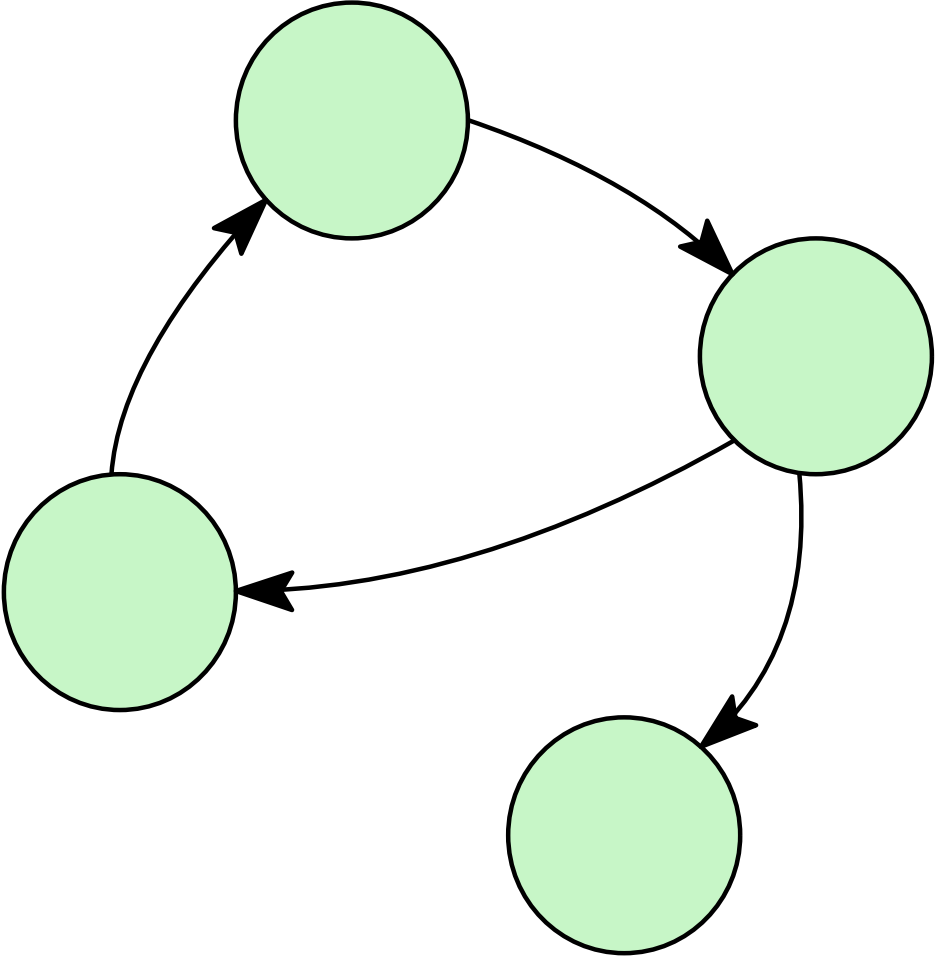
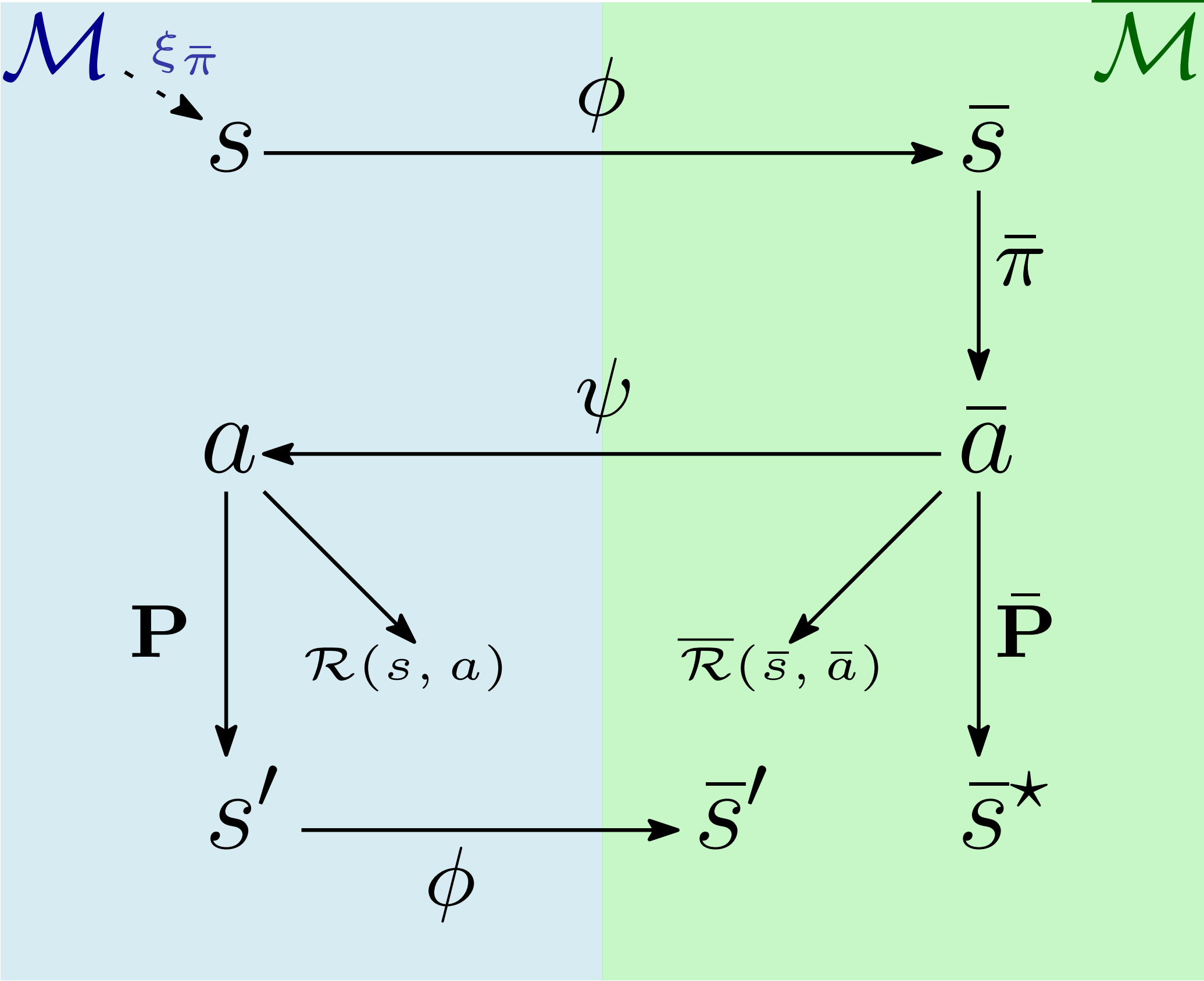
$$\tilde{d}_{\pi}(s_1, s_2) = \sup_{V \in \mathcal{F}_{\gamma}^{\mathcal{L}}(\pi)} \left| V_{\pi}(s_1) - V_{\pi}(s_2) \right| \quad \forall s_1, s_2 \in \mathcal{S}$$

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

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

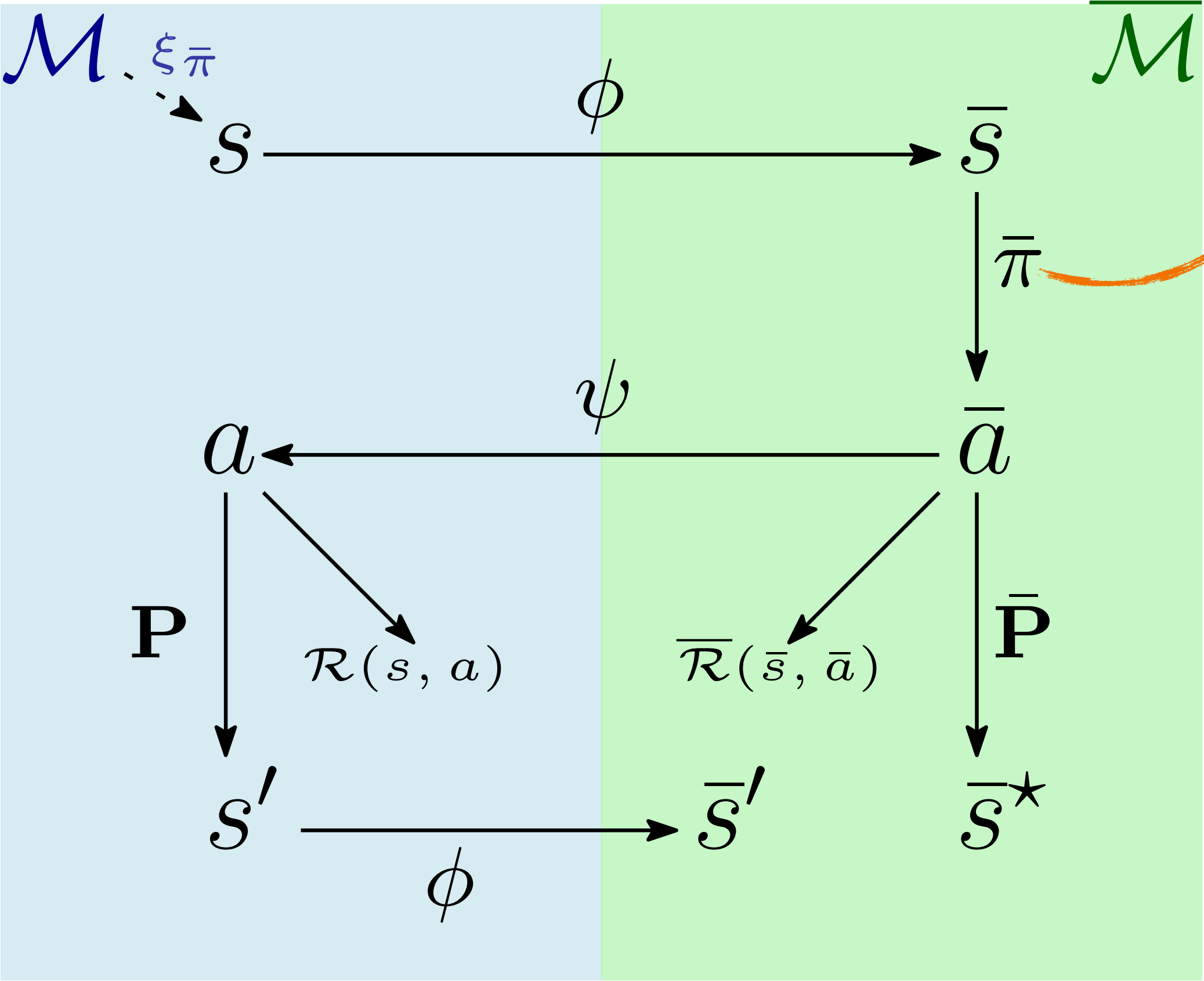
Latent Flow

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

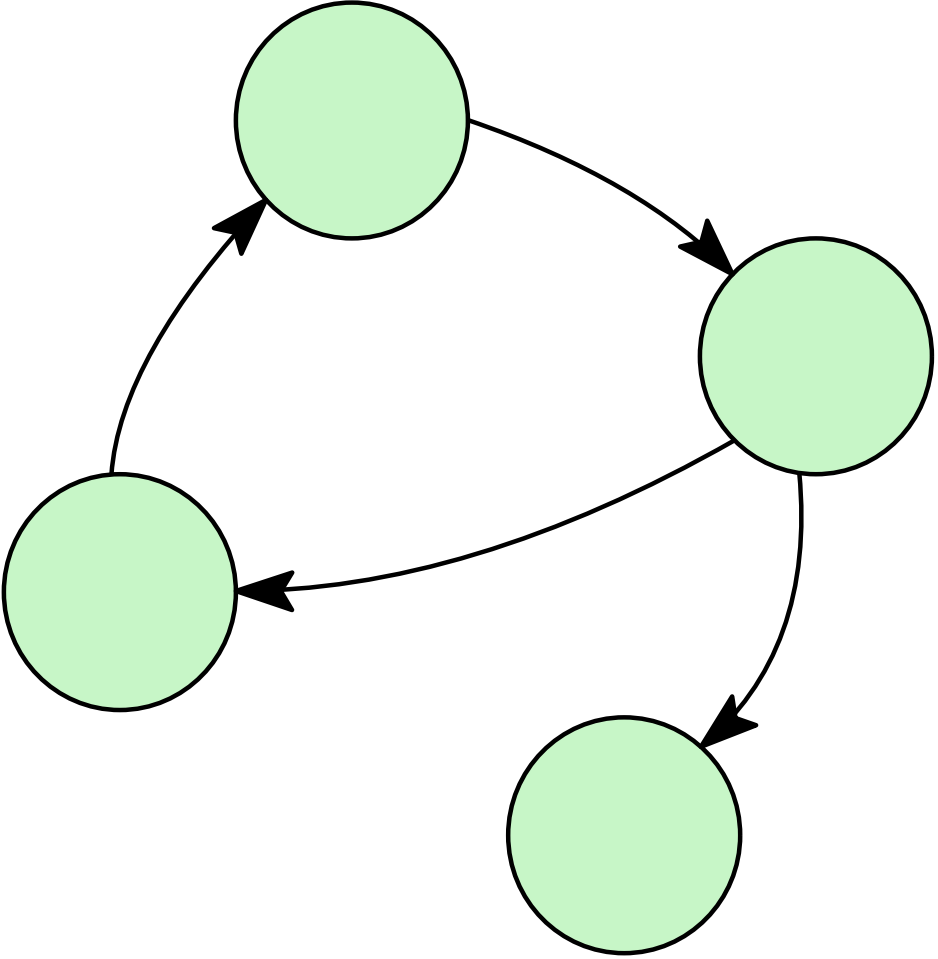


Latent Flow

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



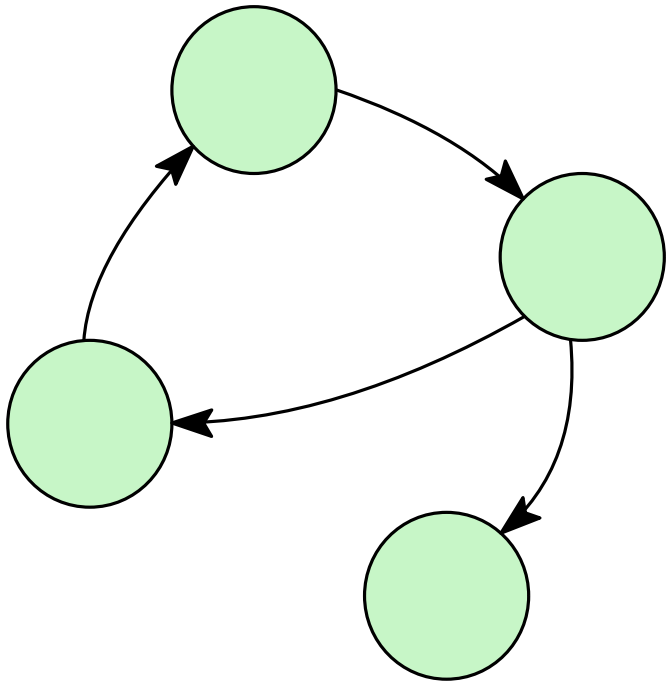
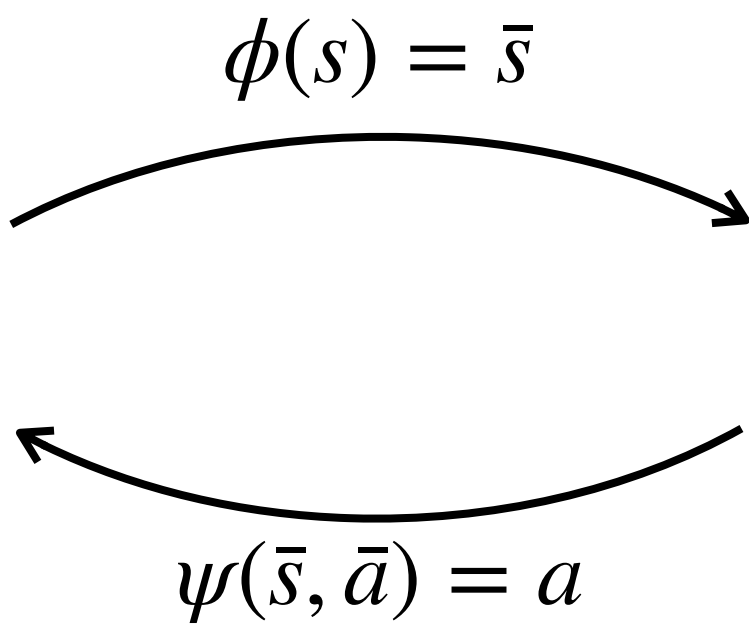
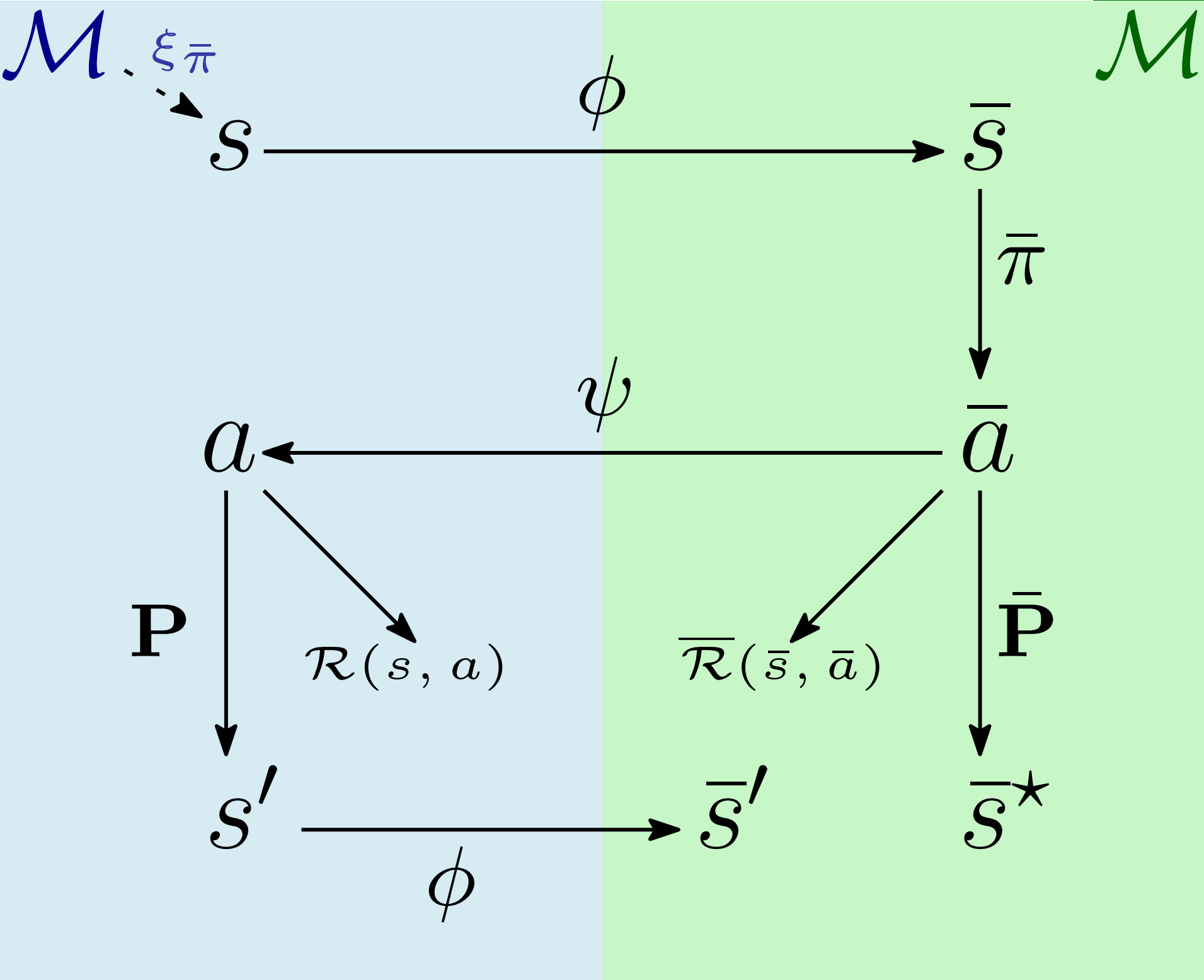
Policy distilled from the RL policy



Latent Flow

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

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



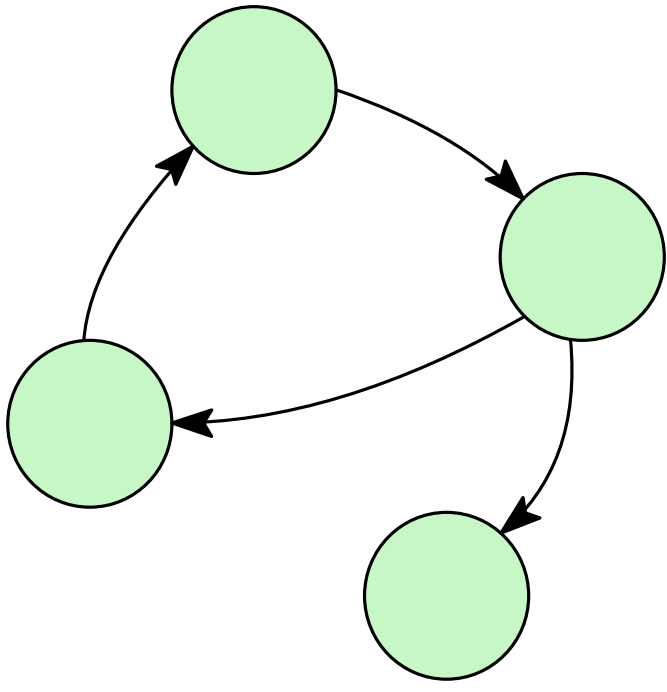
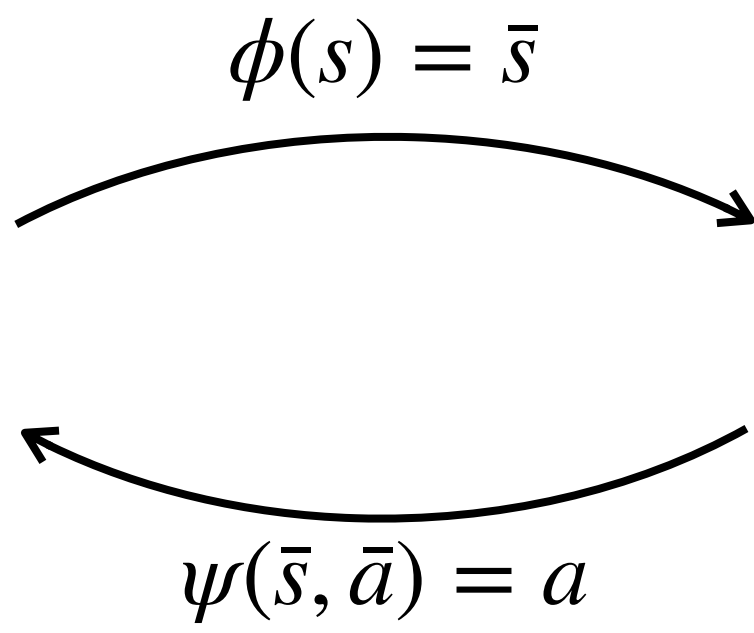
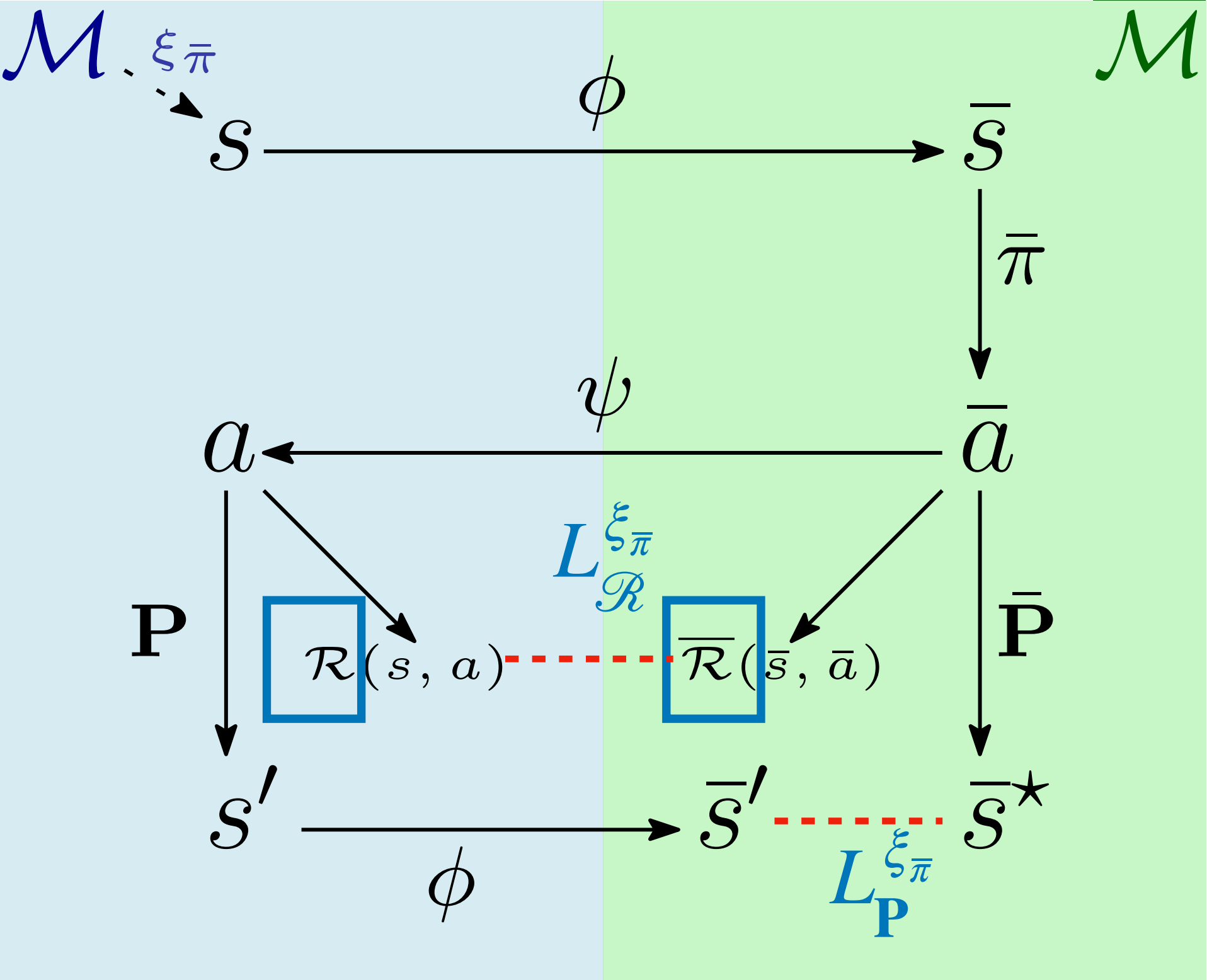
Latent Flow

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

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

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

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



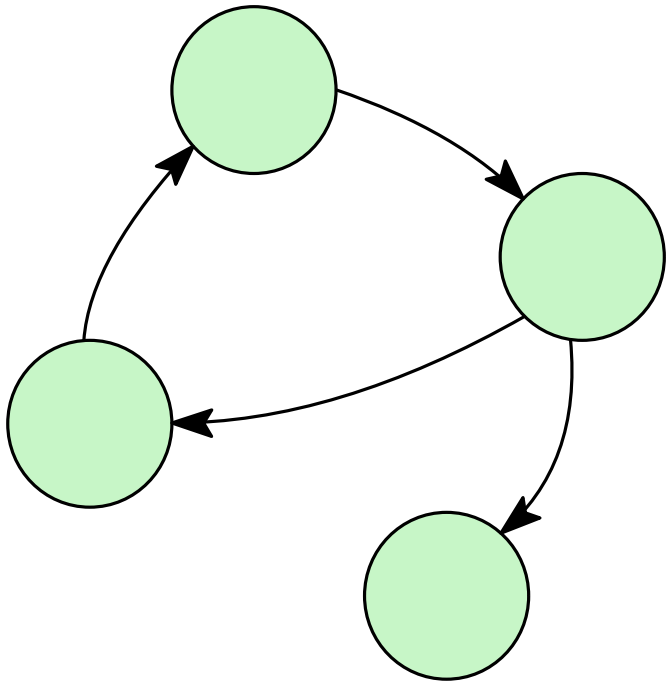
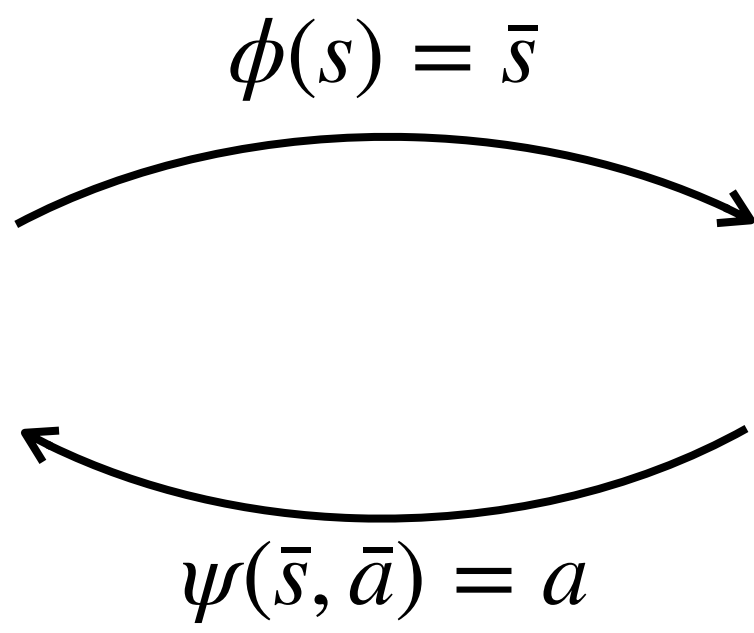
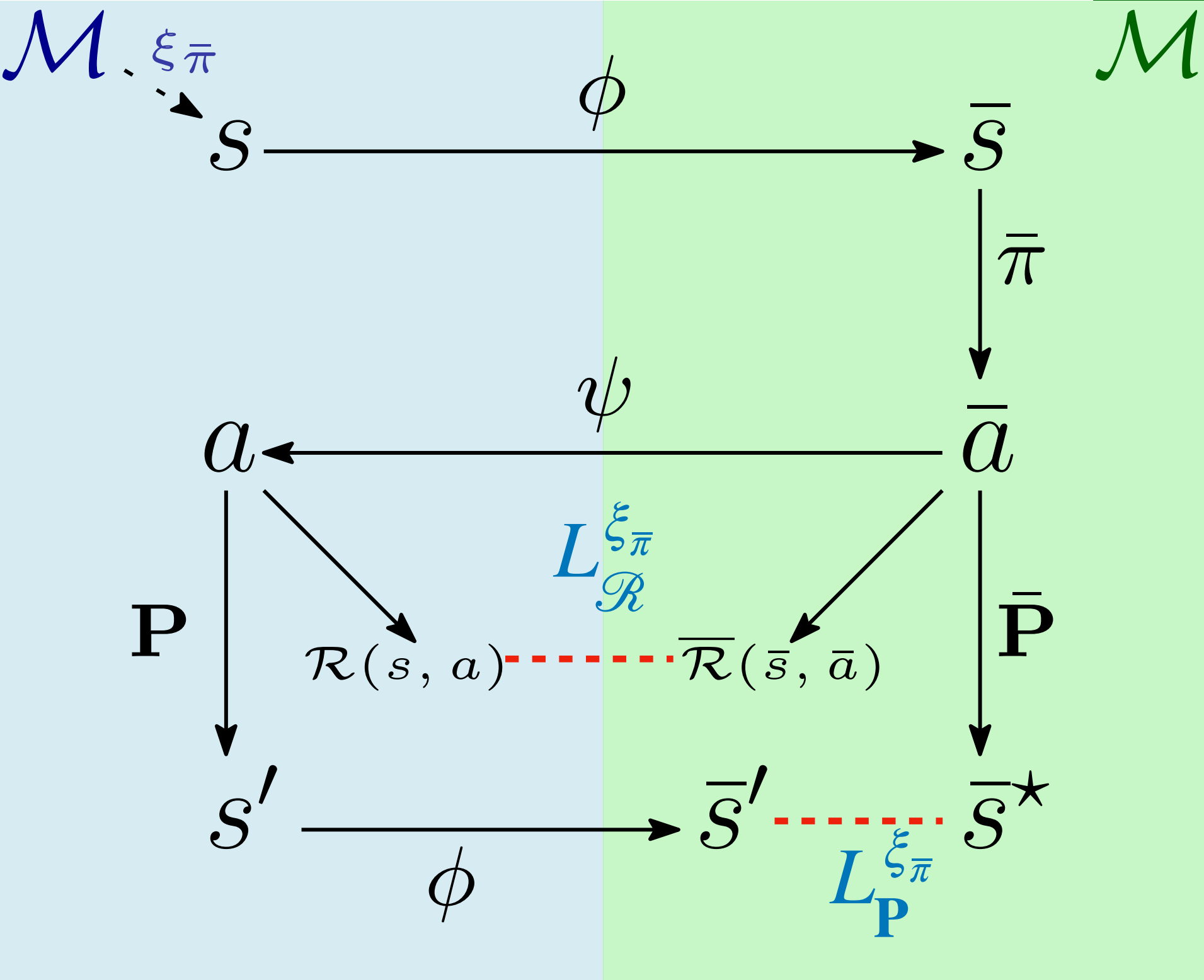
Latent Flow

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

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

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

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



Latent Flow

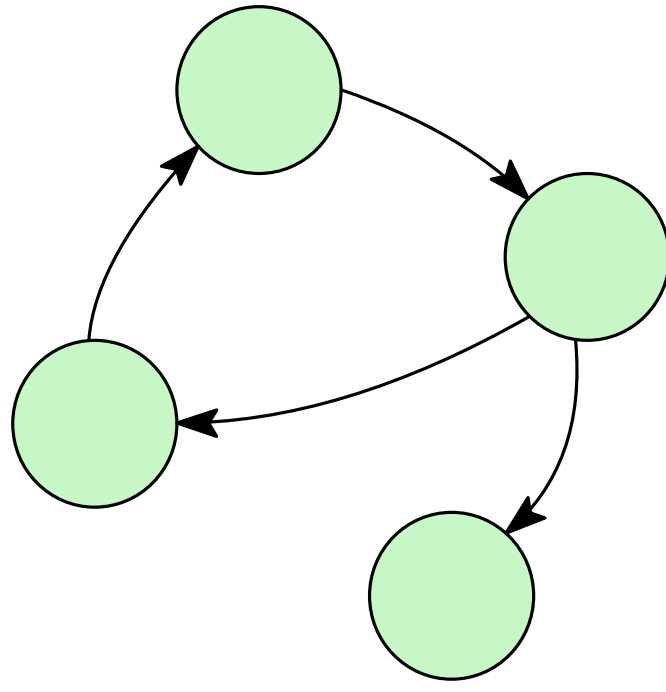
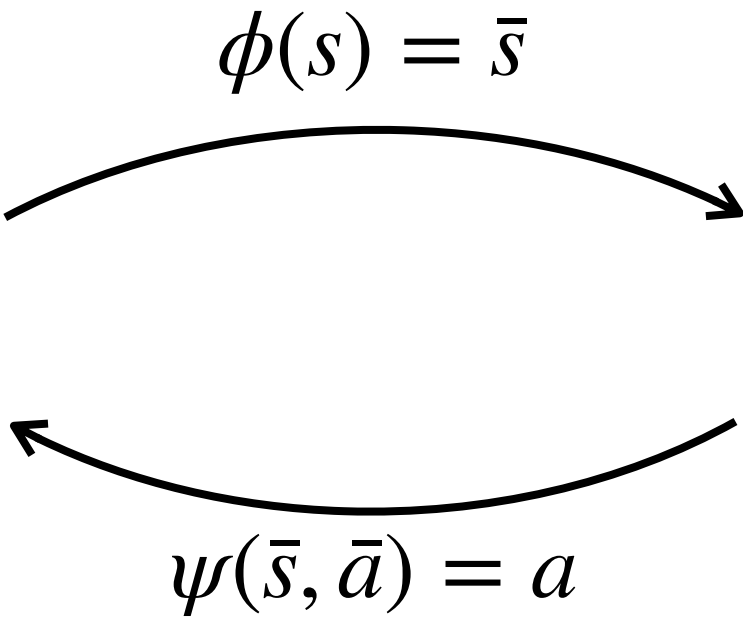
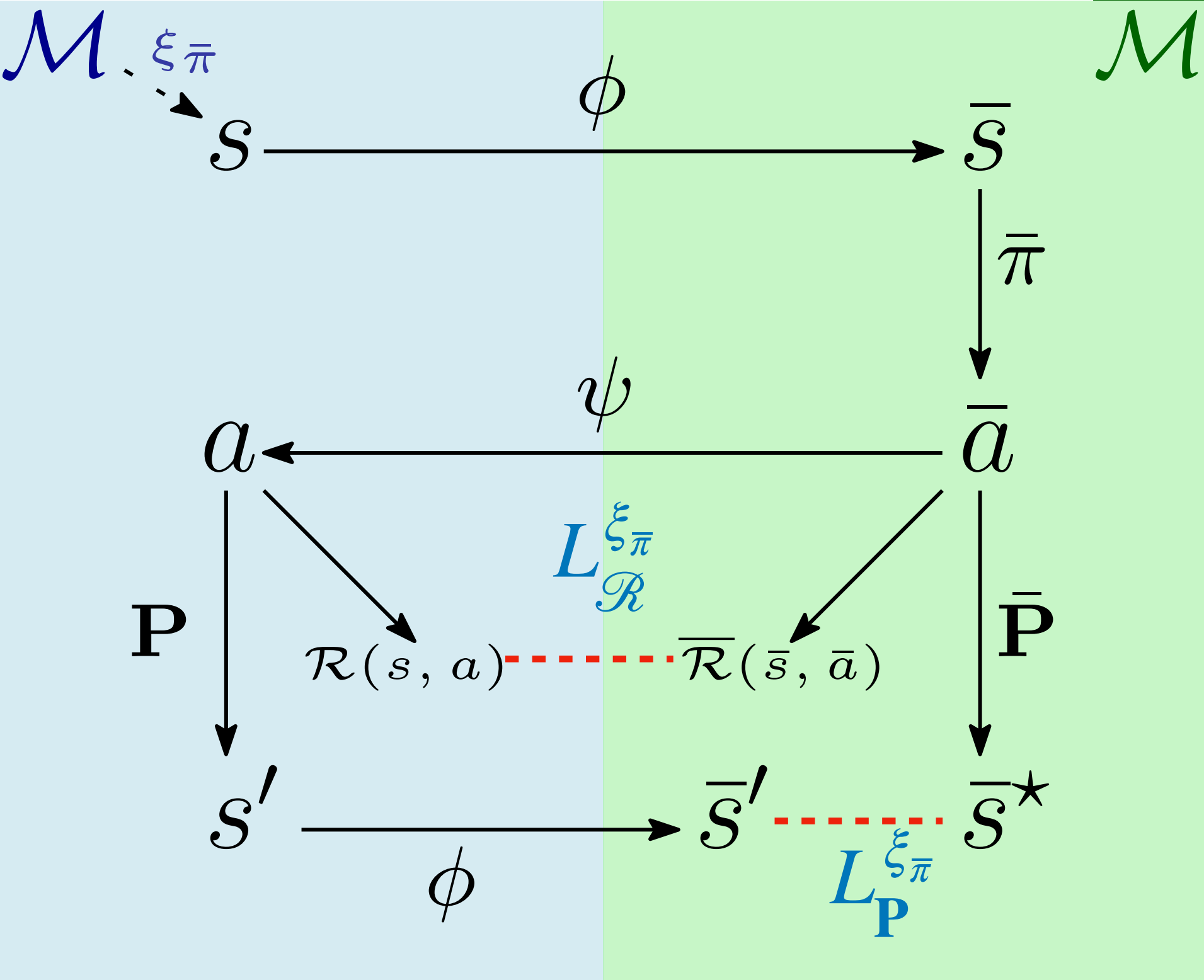
Execution of a latent policy $\bar{\pi}$ in the original model: **Local Losses**

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

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

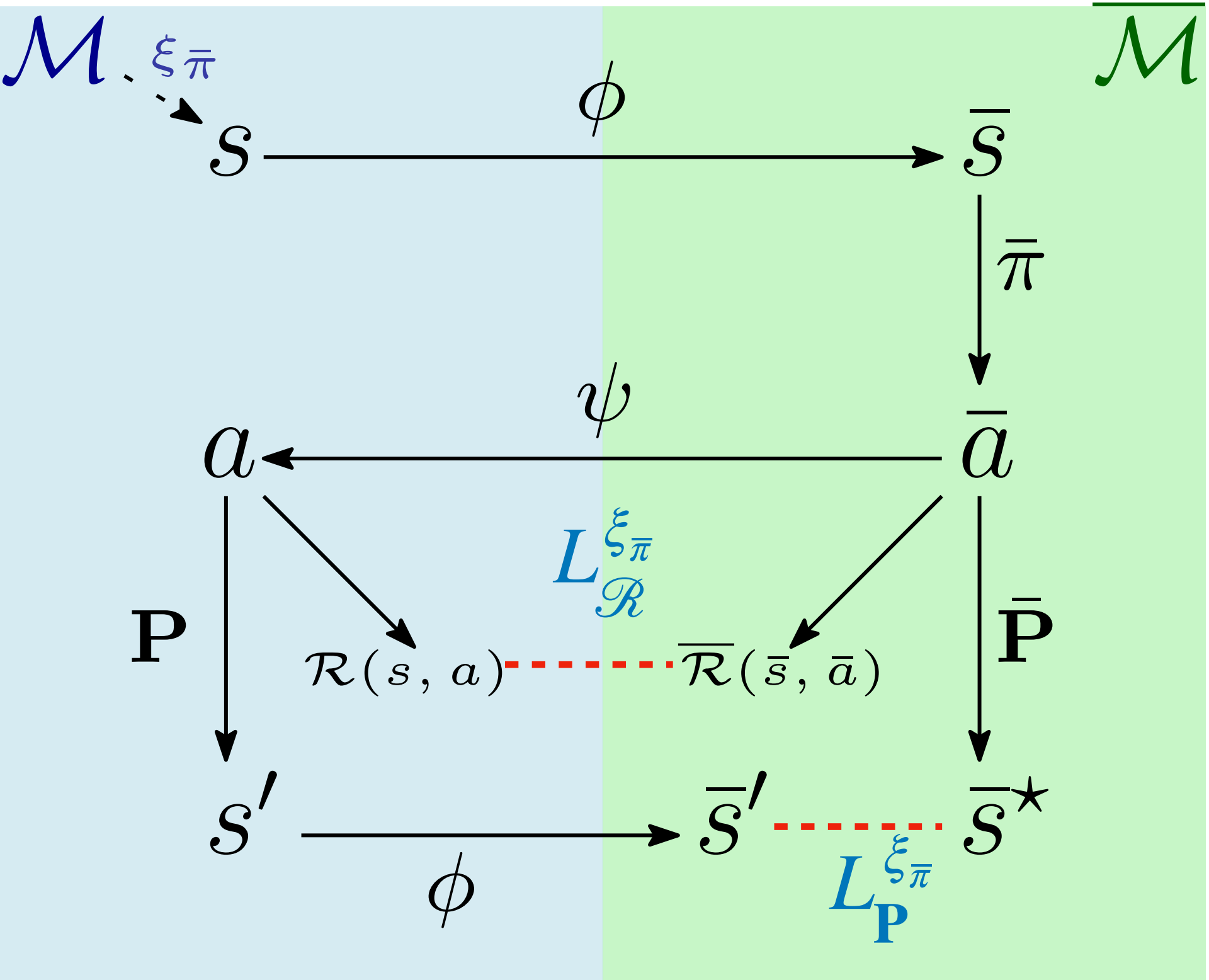
$$L_{\mathcal{R}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} \left| \mathcal{R}(s, \bar{a}) - \bar{\mathcal{R}}(\phi(s), \bar{a}) \right|$$

- Abstraction quality:** $\mathbb{E}_{s \sim \xi_{\bar{\pi}}} \tilde{d}_{\bar{\pi}}(s, \phi(s)) \leq \frac{L_{\mathcal{R}}^{\xi_{\bar{\pi}}} + \gamma L_{\mathbf{P}}^{\xi_{\bar{\pi}}}}{1 - \gamma}$



Latent Flow

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



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

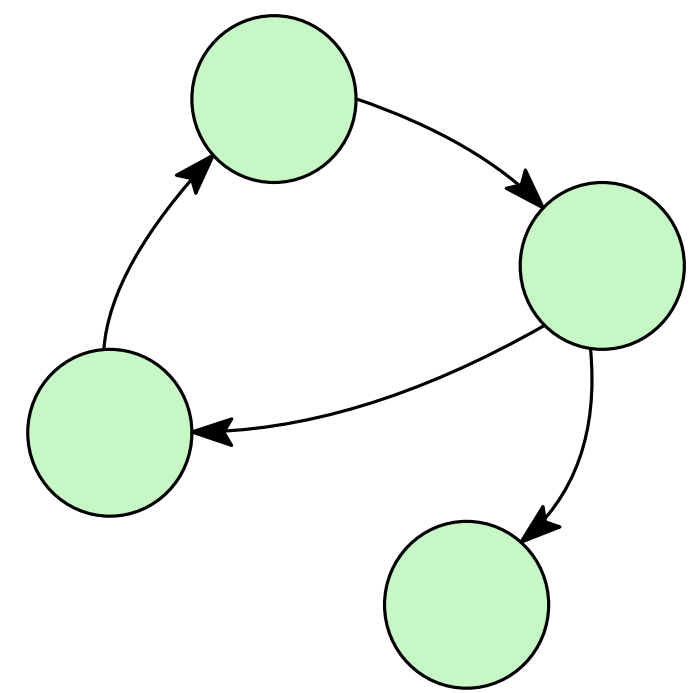
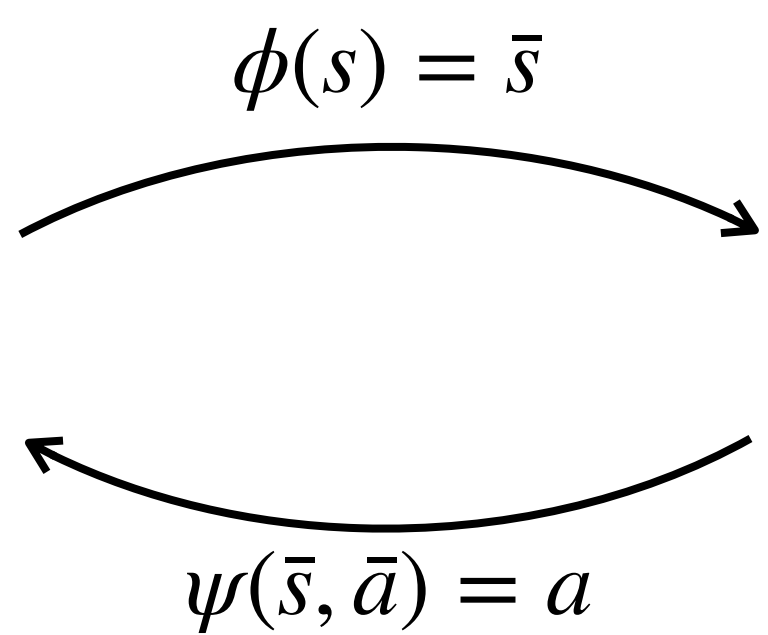
$$L_{\mathbf{P}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} W_{d_{\bar{\mathcal{S}}}} \left(\phi \mathbf{P}(\cdot | s, \bar{a}), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a}) \right)$$

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

- **Abstraction quality:** $\mathbb{E}_{s \sim \xi_{\bar{\pi}}} \tilde{d}_{\bar{\pi}}(s, \phi(s)) \leq \frac{L_{\mathcal{R}}^{\xi_{\bar{\pi}}} + \gamma L_{\mathbf{P}}^{\xi_{\bar{\pi}}}}{1 - \gamma}$

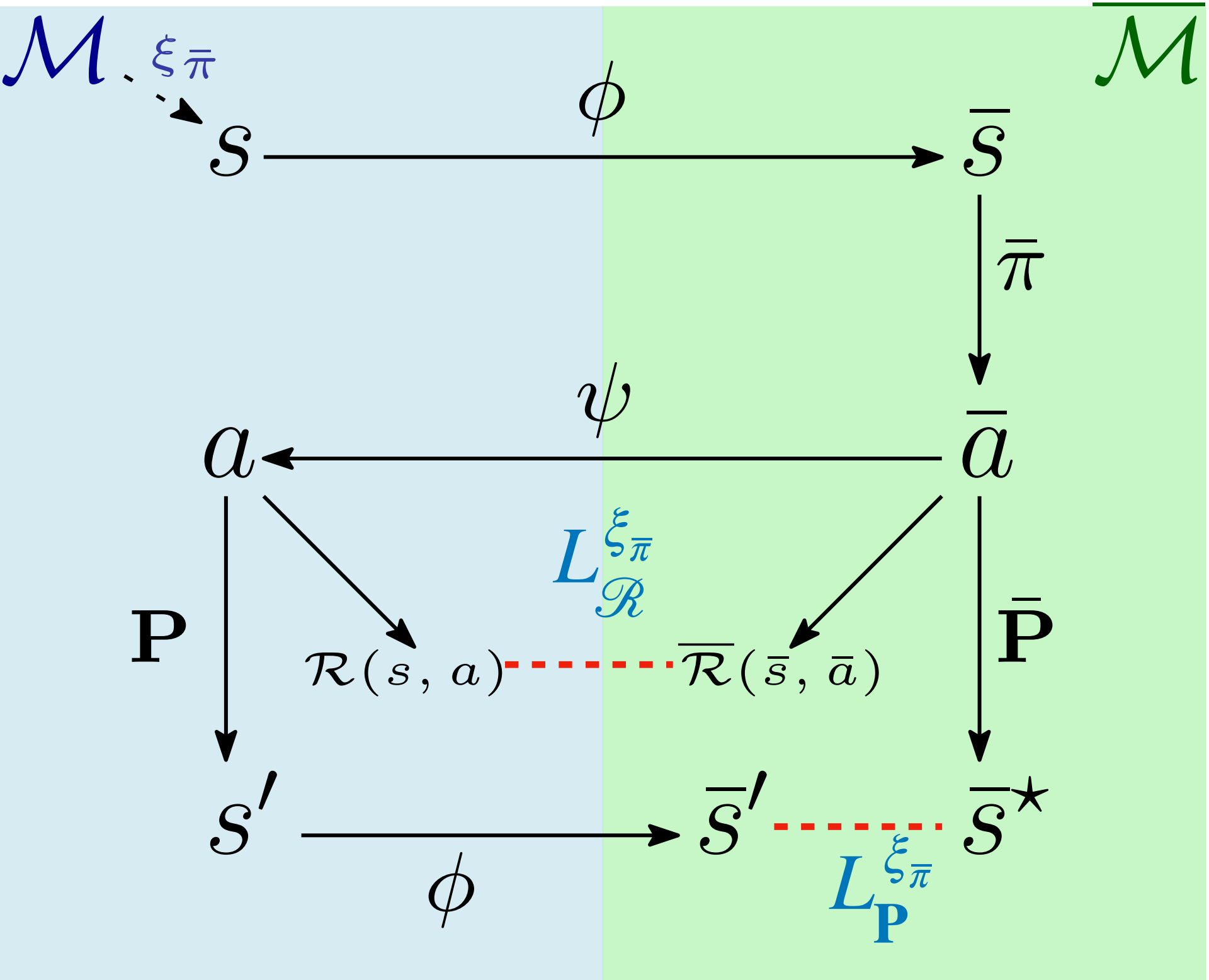
- **Representation quality:** for all $s_1, s_2 \in \mathcal{S}$ such that $\phi(s_1) = \phi(s_2)$

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



Latent Flow

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



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

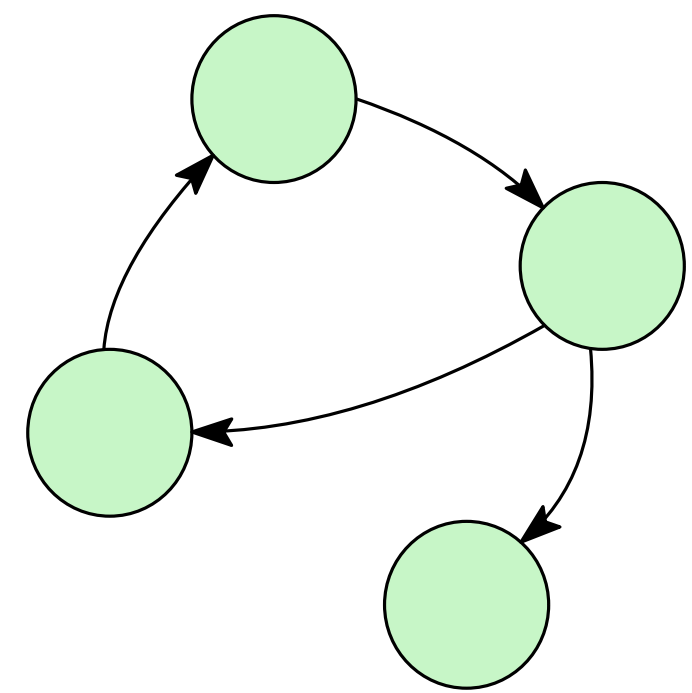
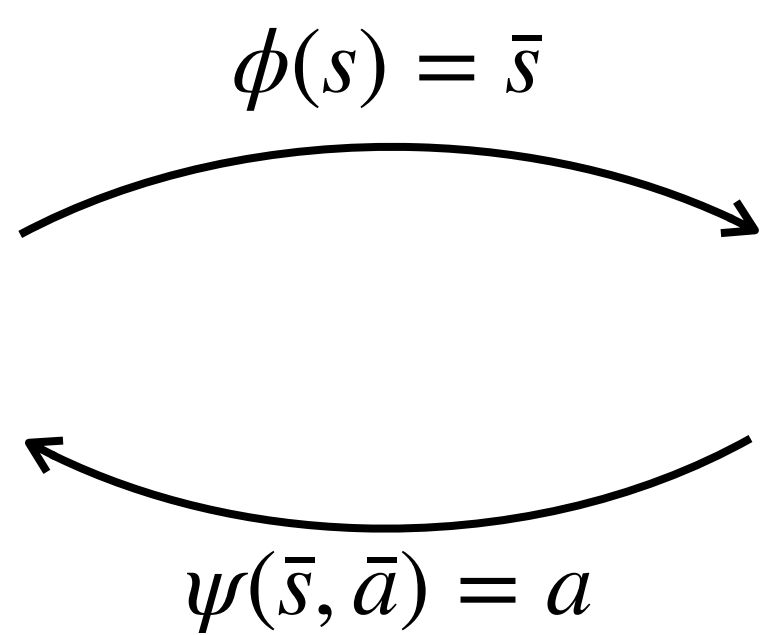
$$L_{\mathbf{P}}^{\xi_{\bar{\pi}}} = \mathbb{E}_{s, \bar{a} \sim \xi_{\bar{\pi}}} W_{d_{\bar{S}}} \left(\phi \mathbf{P}(\cdot | s, \bar{a}), \bar{\mathbf{P}}(\cdot | \phi(s), \bar{a}) \right)$$

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

- **Abstraction quality:** $\mathbb{E}_{s \sim \xi_{\bar{\pi}}} \tilde{d}_{\bar{\pi}}(s, \phi(s)) \leq \frac{L_{\mathcal{R}}^{\xi_{\bar{\pi}}} + \gamma L_{\mathbf{P}}^{\xi_{\bar{\pi}}}}{1 - \gamma}$

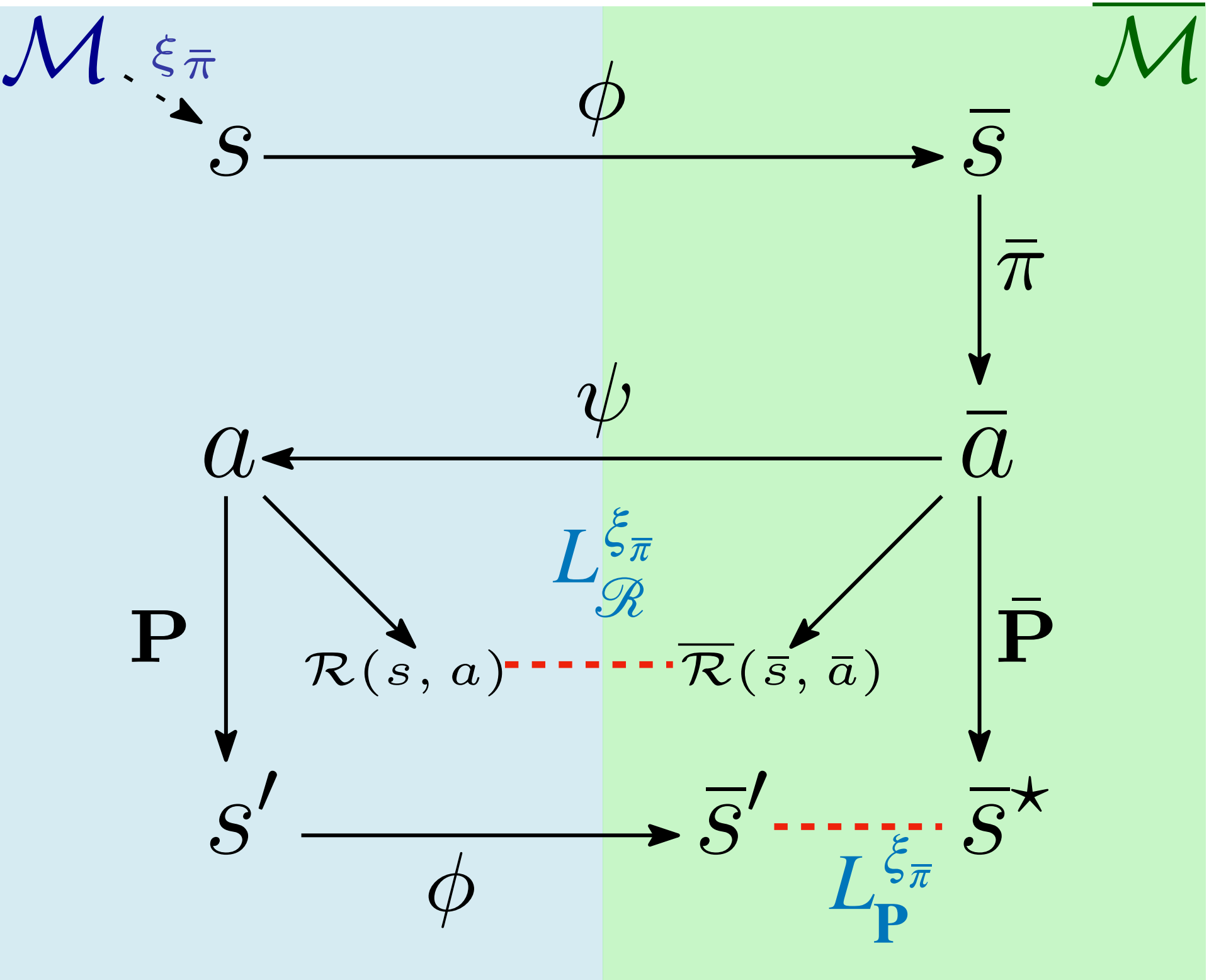
- **Representation quality:** for all $s_1, s_2 \in \mathcal{S}$ such that $\phi(s_1) = \phi(s_2)$

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



Latent Flow

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



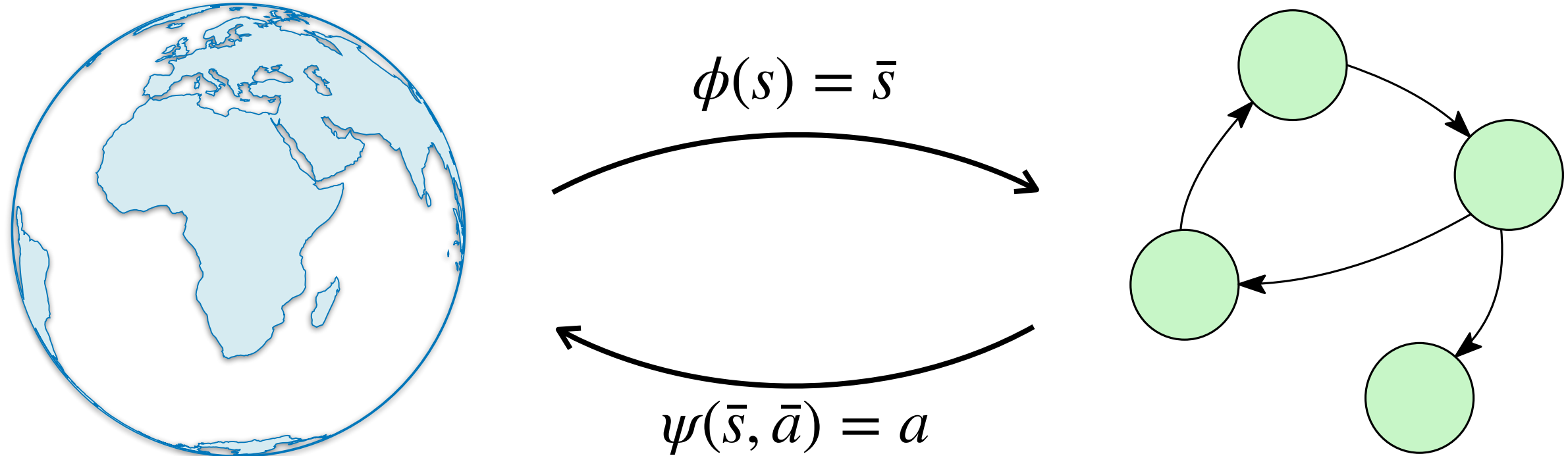
- Latent policy $\bar{\pi}$, stationary distribution $\xi_{\bar{\pi}}$

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

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

- **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 $\phi(s_1) = \phi(s_2)$

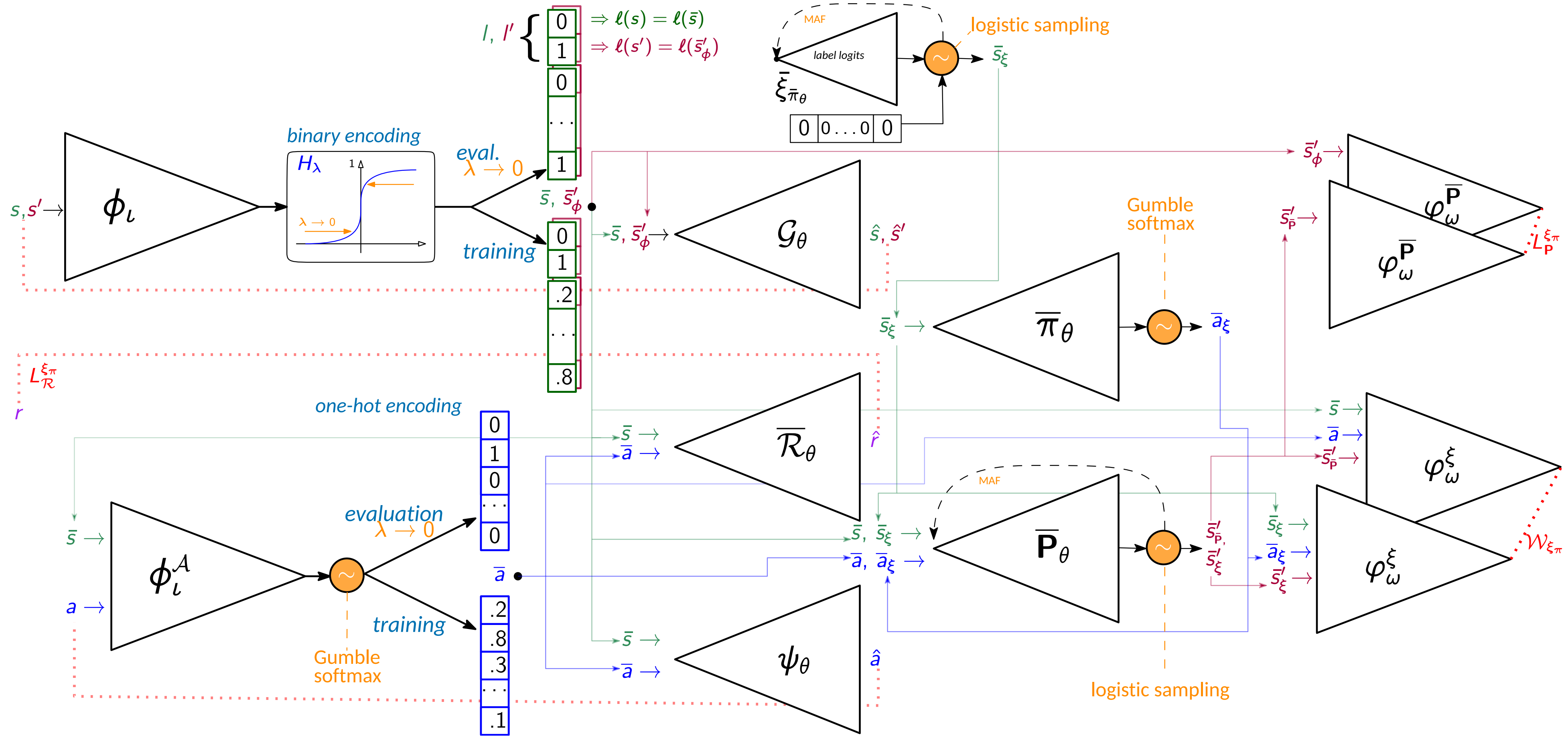
$$\tilde{d}_{\bar{\pi}}(s_1, s_2) \leq \left(\frac{L_{\mathcal{R}}^{\xi_{\bar{\pi}}} + \gamma L_{\mathbf{P}}^{\xi_{\bar{\pi}}}}{1 - \gamma} \right) \cdot \left(\xi_{\bar{\pi}}^{-1}(s_1) + \xi_{\bar{\pi}}^{-1}(s_2) \right)$$



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

Wasserstein Auto-encoded Markov Decision Process

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

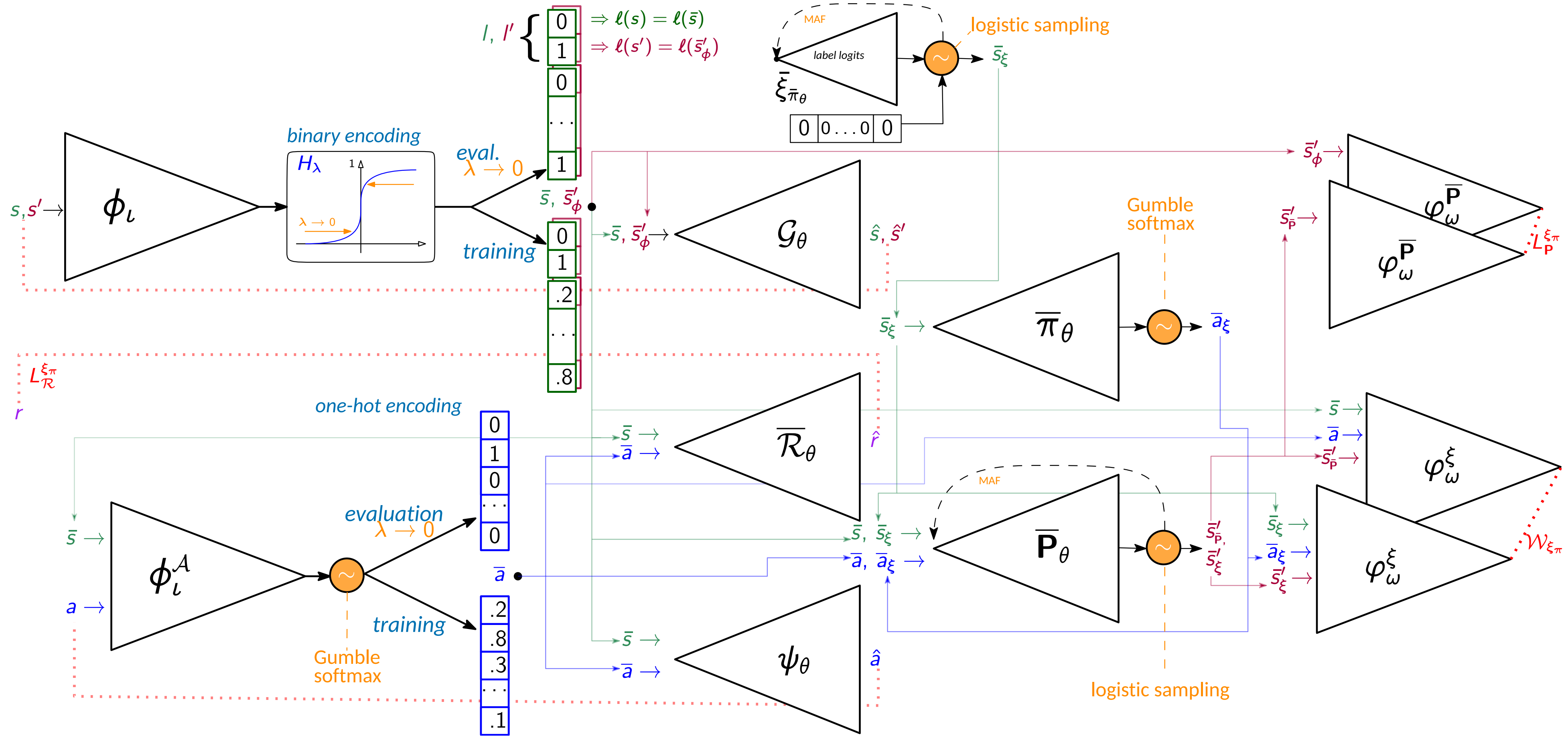


- $\mathcal{W}_{\xi_\pi} = \max_{\omega: \|\varphi_\omega^\xi\| \leq 1} \mathbb{E}_{s, a \sim \xi_\pi} \mathbb{E}_{\bar{a} \sim \phi_l^A(\cdot | \phi_l(s), a)} \mathbb{E}_{\bar{s}' \sim \bar{\mathbf{P}}_\theta(\cdot | \bar{s}, \bar{a})} \varphi_\omega^\xi(\phi_l(s), \bar{a}, \bar{s}') - \mathbb{E}_{\bar{s},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

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

Wasserstein Auto-encoded Markov Decision Process

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

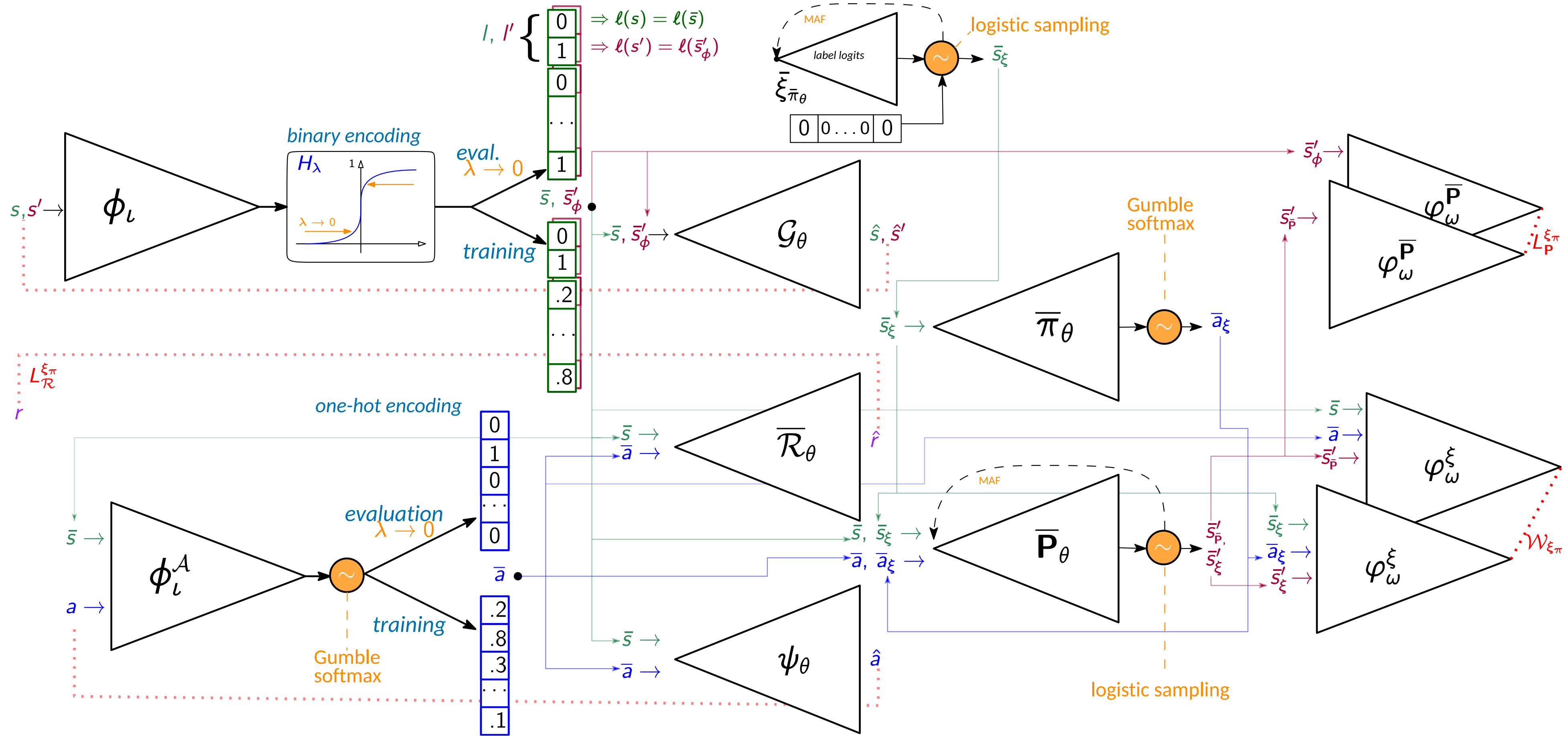


- $\mathcal{W}_{\xi_\pi} = \max_{\omega: \|\varphi_\omega^\xi\| \leq 1} \mathbb{E}_{s, a \sim \xi_\pi} \mathbb{E}_{\bar{a} \sim \phi_l^A(\cdot | \phi_l(s), a)} \mathbb{E}_{\bar{s}' \sim \bar{\mathbf{P}}_\theta(\cdot | \bar{s}, \bar{a})} \varphi_\omega^\xi(\phi_l(s), \bar{a}, \bar{s}') - \mathbb{E}_{\bar{s},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

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

Wasserstein Auto-encoded Markov Decision Process

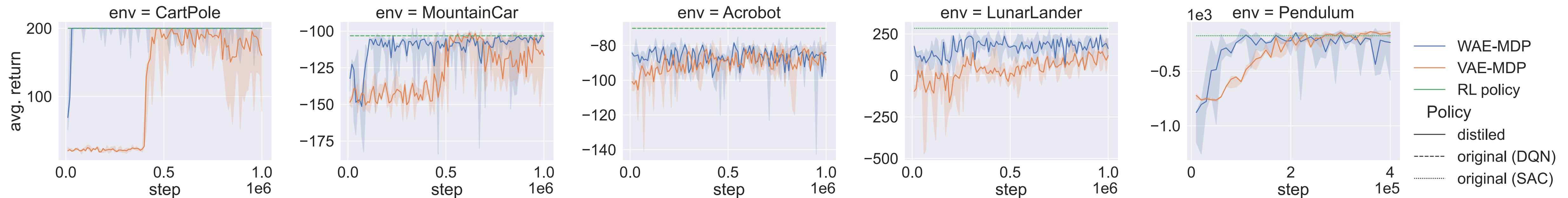
$$\min_{l, \theta} \mathbb{E}_{s, a, s' \sim \xi_\pi} \mathbb{E}_{\bar{s}, \bar{a}, \bar{s}' \sim \phi_l(\cdot | s, a, s')} \|\langle s, a, s' \rangle - \langle \mathcal{G}_\theta(\bar{s}), \psi_\theta(\bar{s}, \bar{a}), \mathcal{G}_\theta(\bar{s}') \rangle\| + L_{\mathcal{R}}^{\xi_\pi} + \beta \left(\mathcal{W}_{\xi_\pi}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



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

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

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*

Distillation of RL Policies with Formal Guarantees via Variational Abstraction of Markov Decision Processes.
Florent Delgrange, Ann Nowé, Guillermo A. Pérez (2022). AAI 2022.