



POLITECNICO
MILANO 1863

REINFORCEMENT LEARNING IN CONFIGURABLE CONTINUOUS ENVIRONMENTS

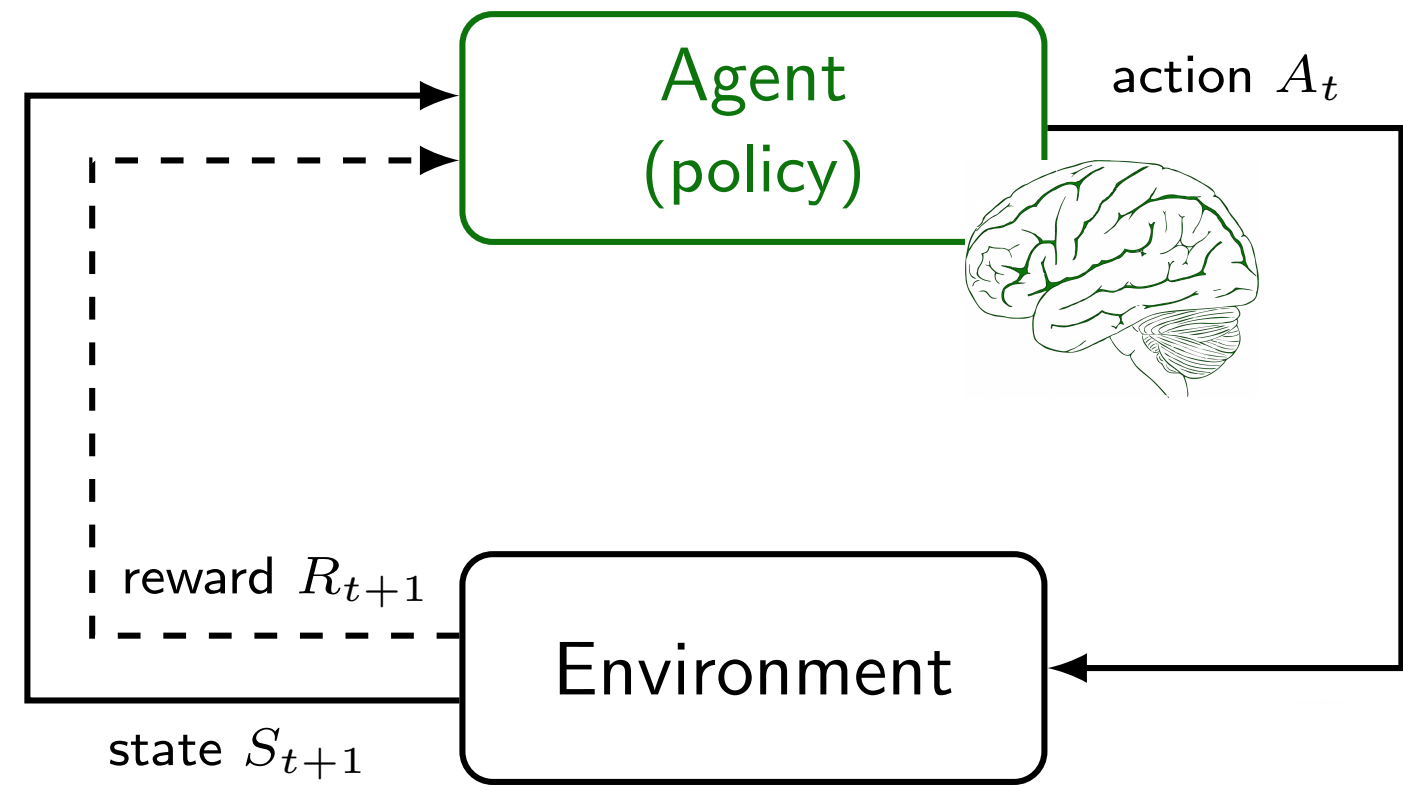
ALBERTO MARIA METELLI, EMANUELE GHELFI AND MARCELLO RESTELLI

{albertomaria.metelli, marcello.restelli}@polimi.it, emanuele.ghelfi@mail.polimi.it



PROBLEM

Markov Decision Process (MDP, Puterman, 2014)

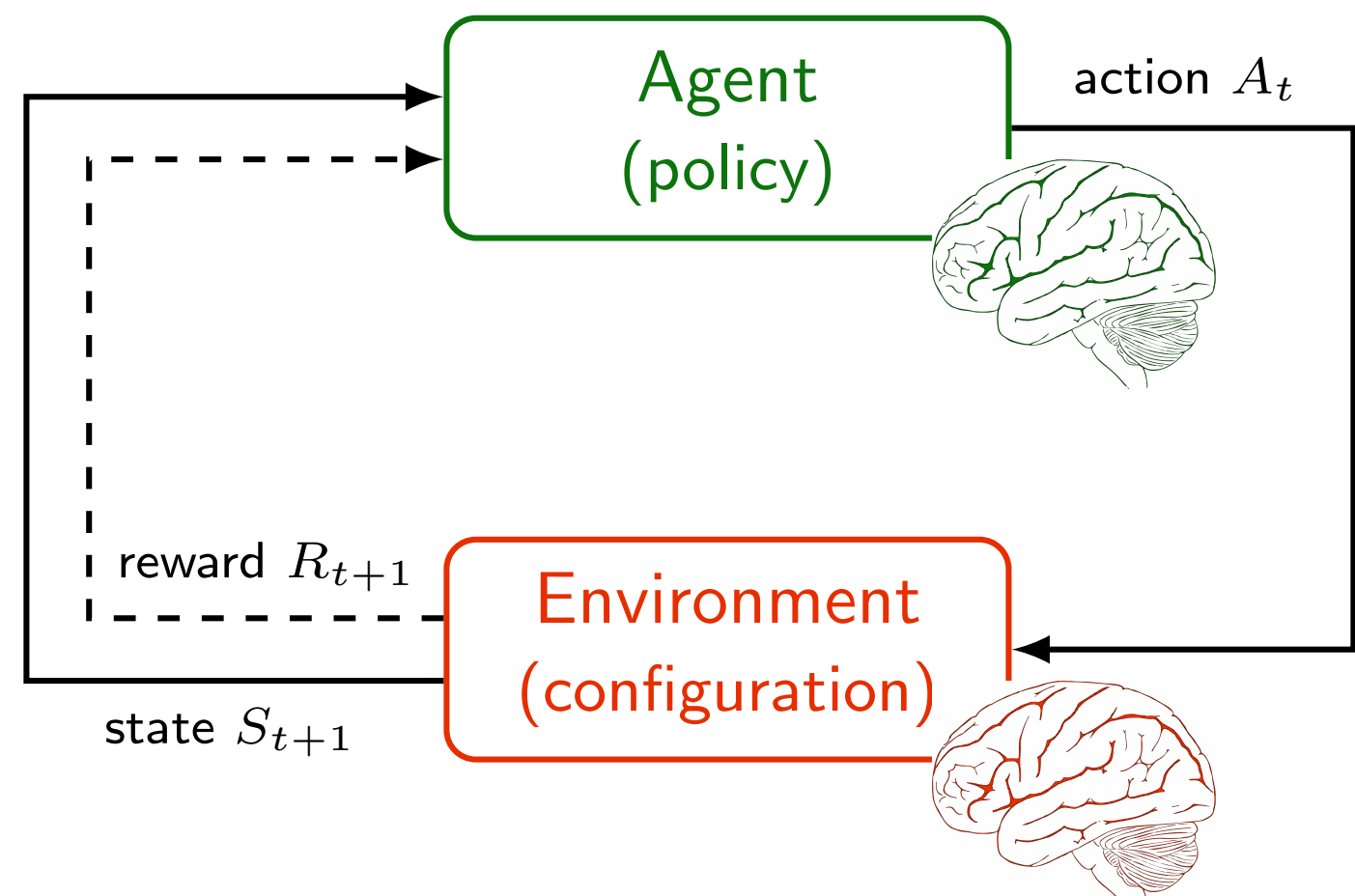


$$S_0 \sim \mu, A_t \sim \pi_\theta(\cdot|S_t), S_{t+1} \sim p(\cdot|S_t, A_t)$$

- Learn the policy parameters θ under the fixed environment p :

$$\theta^* = \arg \max_{\theta \in \Theta} J(\theta) = \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t R_{t+1} \right]$$

Configurable Markov Decision Process (Conf-MDP, Metelli et al., 2018)



$$S_0 \sim \mu, A_t \sim \pi_\theta(\cdot|S_t), S_{t+1} \sim p_\omega(\cdot|S_t, A_t)$$

- Learn the policy parameters θ together with the environment configuration ω :

$$\theta^*, \omega^* = \arg \max_{\theta \in \Theta, \omega \in \Omega} J(\theta, \omega) = \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t R_{t+1} \right]$$

MOTIVATION AND CONTRIBUTIONS

- Existing approaches (SPMI, Metelli et al., 2018):
 - work in **finite** state-actions spaces only
 - require the **full knowledge** of the environment dynamics

How to solve continuous Conf-MDPs without the knowledge of the exact dynamics?

- REMPS** is able to learn in **continuous Conf-MDPs** and can be equipped with an **approximate transition model**
 - REMPS alternates **optimization** and **projection** phases
 - Inspired to *trust region* methods (Peters et al., 2010)
 - We provide a *finite-sample analysis* for the single step
 - We *empirically* evaluate REMPS on Conf-MDPs

RELATIVE ENTROPY MODEL POLICY SEARCH (REMPS)

Optimization

- π_θ and p_ω induce a *stationary distribution* d_{π_θ, p_ω}
- Goal:** find a new stationary distribution d' in a *trust region* centered in d_{π_θ, p_ω} (PRIMAL $_\kappa$):

$$\begin{aligned} \max_{d'} J_{d'} &= \mathbb{E}_{S, A, S' \sim d'} [r(S, A, S')] \\ \text{s.t. } D_{\text{KL}}(d' \| d_{\pi_\theta, p_\omega}) &\leq \kappa \end{aligned}$$

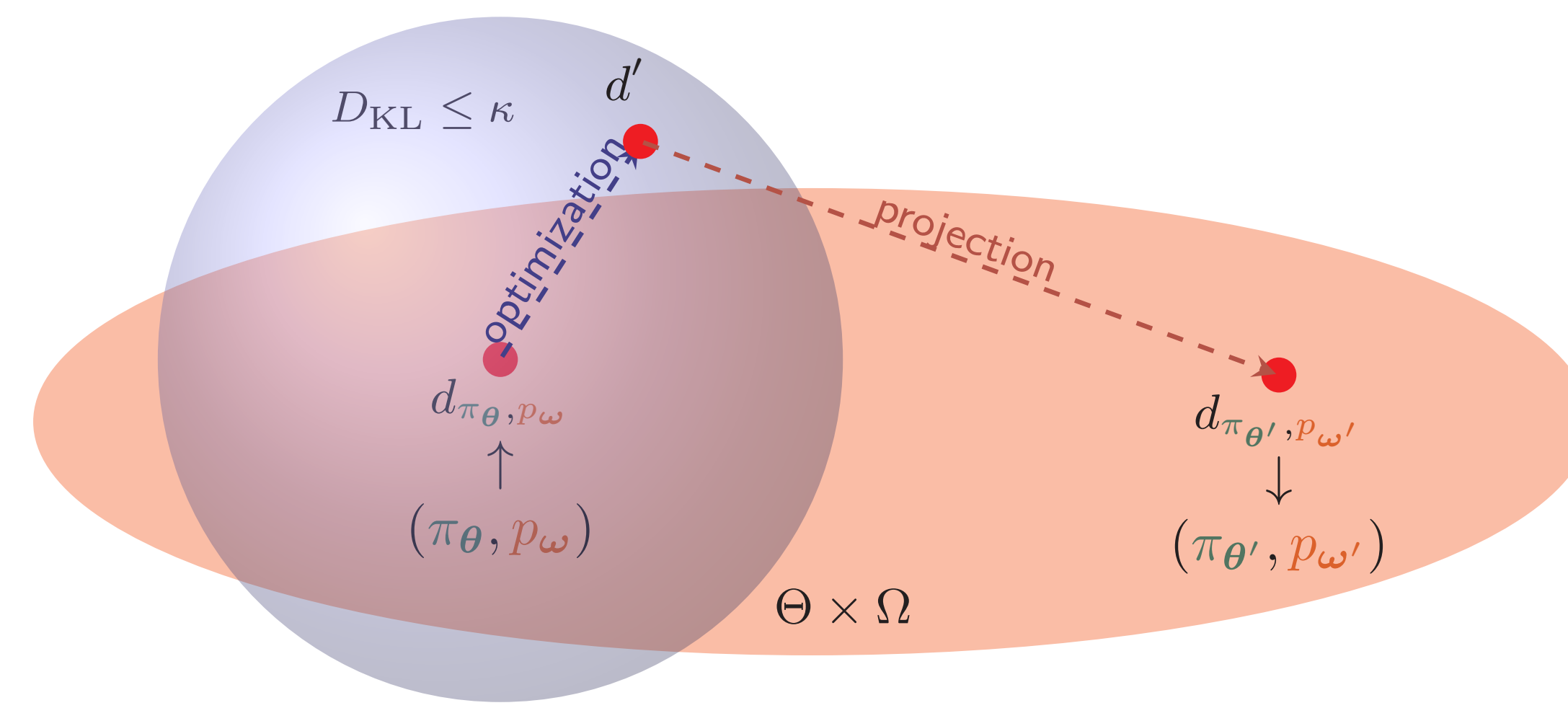
- $\kappa > 0$ defines the trust region in KL-divergence
- We solve the dual problem (DUAL $_\kappa$) in the Lagrange multiplier η :

$$\min_{\eta \in [0, +\infty)} \eta \log \mathbb{E}_{S, A, S' \sim d_{\pi_\theta, p_\omega}} \left[\exp \left(\frac{1}{\eta} r(S, A, S') + \kappa \right) \right]$$

- d' *exponentially reweighs* the probability of each (s, a, s') by the reward $r(s, a, s')$

$$d'(s, a, s') \propto d_{\pi_\theta, p_\omega}(s, a, s') \exp \left(\frac{1}{\eta} r(s, a, s') \right)$$

- When using *samples* $\{(s_i, a_i, s'_i, r_i)\}_{i=1}^N$ we minimize the empirical version of DUAL $_\kappa$



Projection

- d' might fall outside the space of representable stationary distribution, given $\Theta \times \Omega$
- Goal:** perform a *moment projection* onto $\Theta \times \Omega$:
 - project the stationary distribution (PROJ $_{d'}$)

$$\min_{\theta' \in \Theta, \omega' \in \Omega} D_{\text{KL}}(d' \| d_{\pi_{\theta'}, p_{\omega'}})$$

- project the state kernel (PROJ $_{p_\pi}$)

$$\min_{\theta' \in \Theta, \omega' \in \Omega} \mathbb{E}_{S \sim d'} \left[D_{\text{KL}}(p^{\pi'}(\cdot|S) \| p_{\omega'}^{\pi_{\theta'}}(\cdot|S)) \right]$$

- project the policy and model separately (PROJ $_{\pi, p}$)

$$\begin{aligned} \min_{\theta \in \Theta} \mathbb{E}_{S \sim d'} [D_{\text{KL}}(\pi'(\cdot|S) \| \pi_{\theta'}(\cdot|S))] \\ \min_{\omega \in \Omega} \mathbb{E}_{S, A \sim d'} [D_{\text{KL}}(p'(\cdot|S, A) \| p_{\omega'}(\cdot|S, A))] \end{aligned}$$

- When using *samples* $\{(s_i, a_i, s'_i, r_i)\}_{i=1}^N$ we minimize the log-likelihood employing *importance sampling* (Owen, 2013)

Approximate Transition Model

- In the projection phase, we can use an *approximate transition model* \hat{p} learned from samples $\{(s_i, a_i, s'_i, \omega_i)\}_{i=1}^N$ instead of the real one

$$\max_{\hat{p} \in \hat{\mathcal{P}}_\Omega} \frac{1}{N} \sum_{i=1}^N \log \hat{p}(s'_i | s_i, a_i, \omega_i),$$

- \hat{p} can be any model approximator (e.g., neural network, Gaussian process)

THEORETICAL ANALYSIS

- Goal:** bound the performance of the stationary distribution d' obtained with infinite samples and the performance of $(\pi_{\tilde{\theta}'}, p_{\tilde{\omega}'})$ obtained after projection using N i.i.d. samples

$$\begin{aligned} J_{d'} - J(\tilde{\theta}', \tilde{\omega}') &\leq \sqrt{2} r_{\max} \sup_{d: D_{\text{KL}}(d' \| d_{\pi_\theta, p_\omega}) \leq \kappa} \inf_{\bar{\theta} \in \Theta, \bar{\omega} \in \Omega} \sqrt{D_{\text{KL}}(d \| d_{\pi_{\bar{\theta}}, p_{\bar{\omega}}})} \\ &\quad + \tilde{O} \left(\sqrt{\frac{v \log \frac{2eN}{v} + \log \frac{8}{\delta}}{N}} \right) \end{aligned}$$

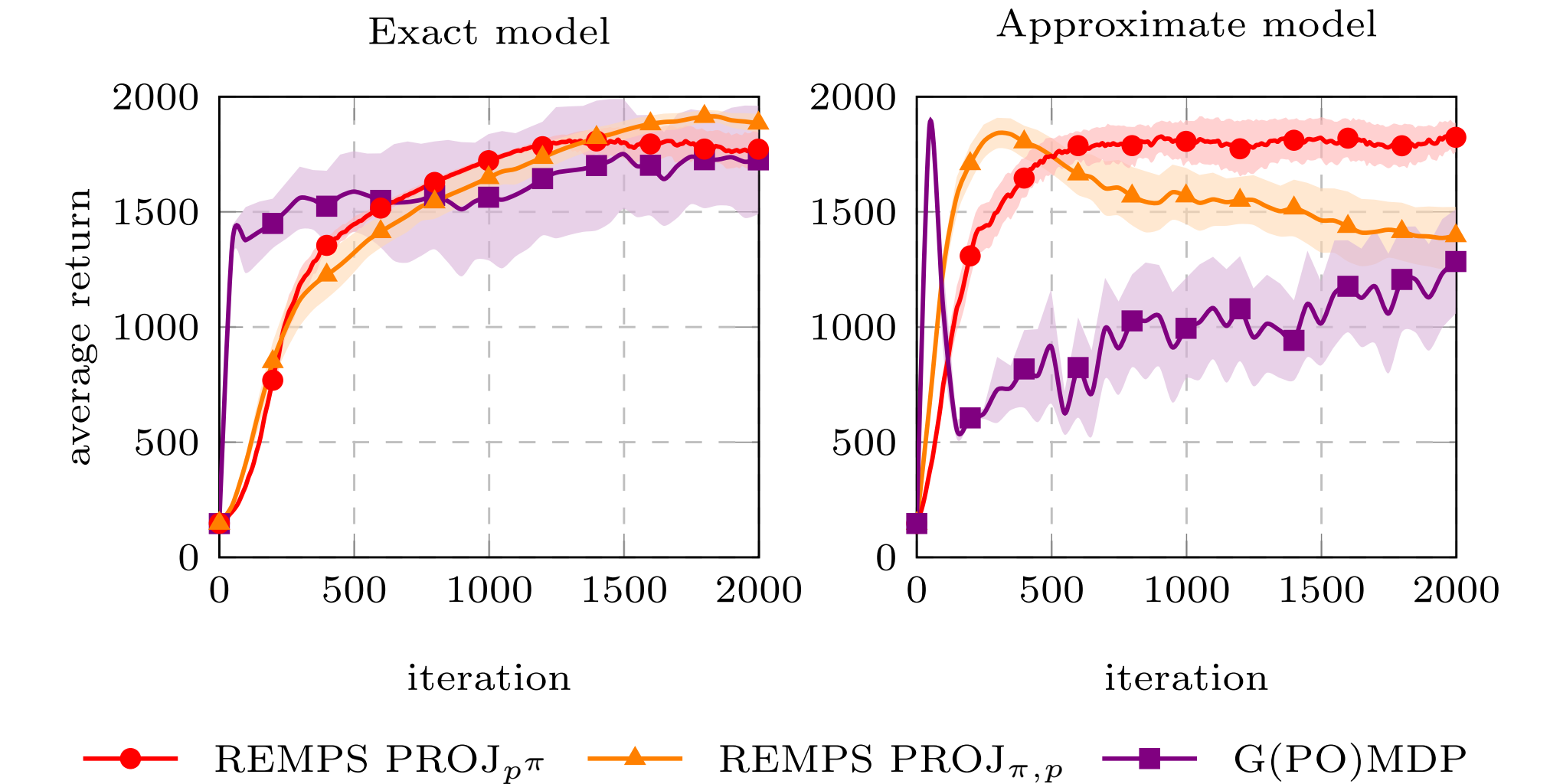
Approximation Error: due to the finite capacity of the parametric space $\Theta \times \Omega$, depends also on the threshold κ

Estimation Error: due to the finite samples N , depends also on the pseudo-dimension v (Cortes et al., 2013)

EXPERIMENTS

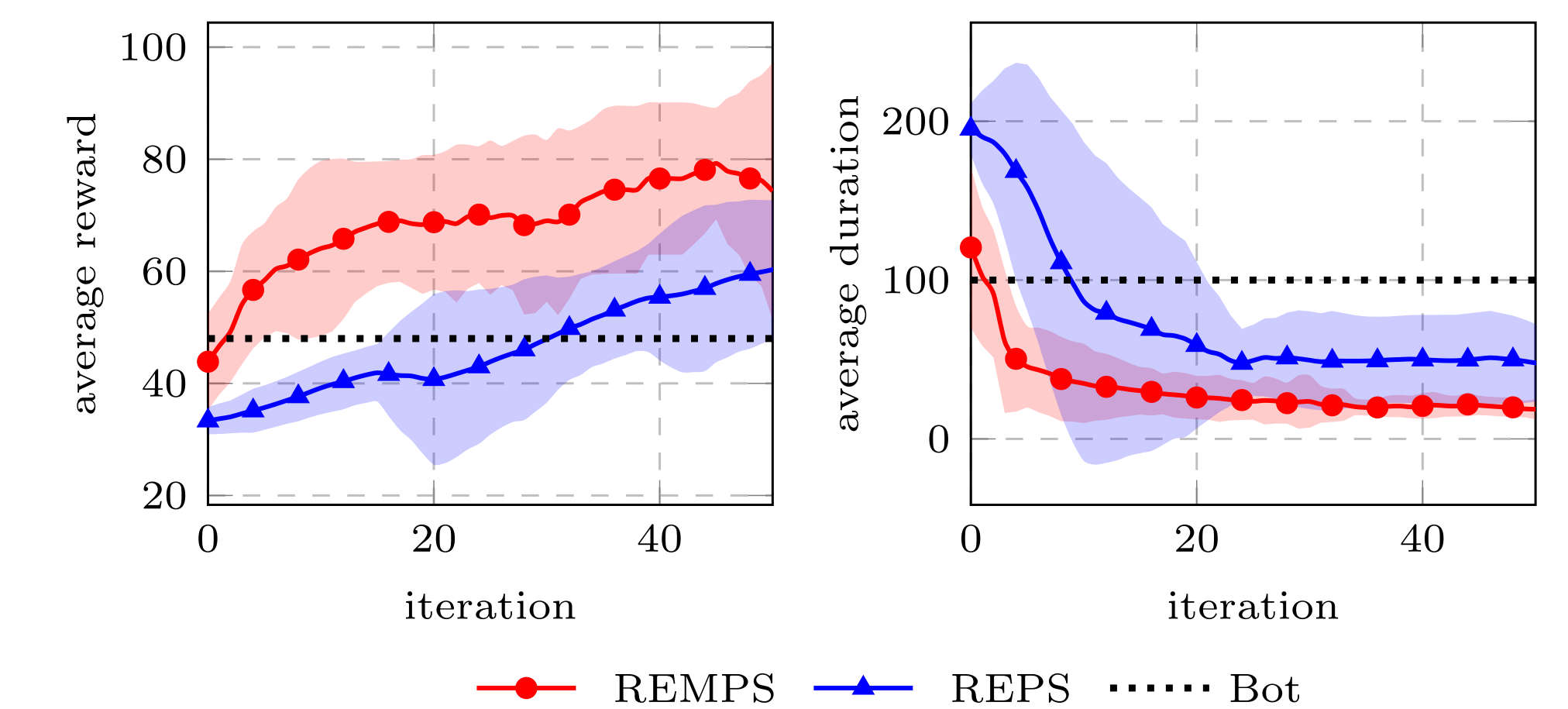
CARTPOLE

- Configure the *cart force*

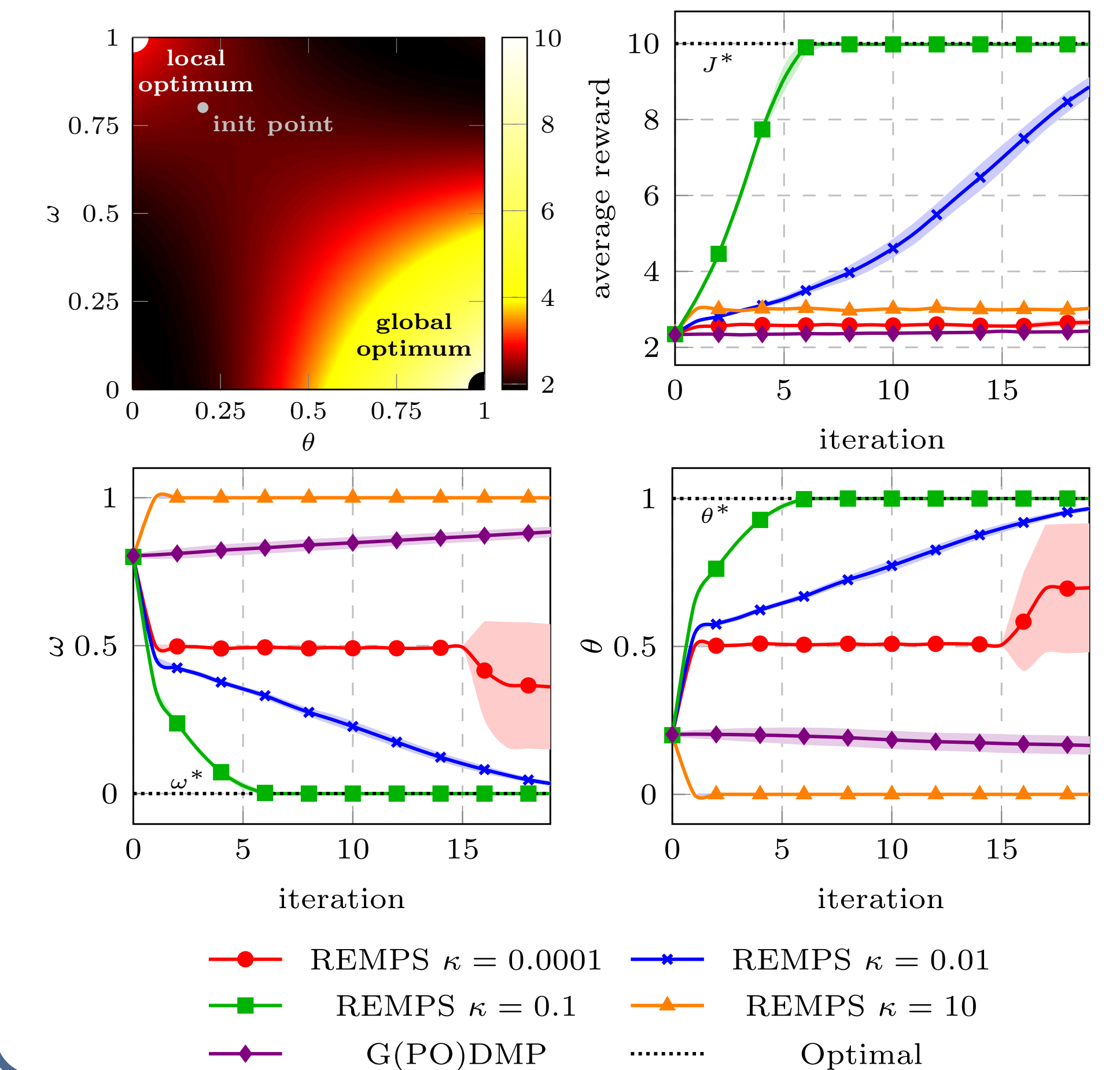
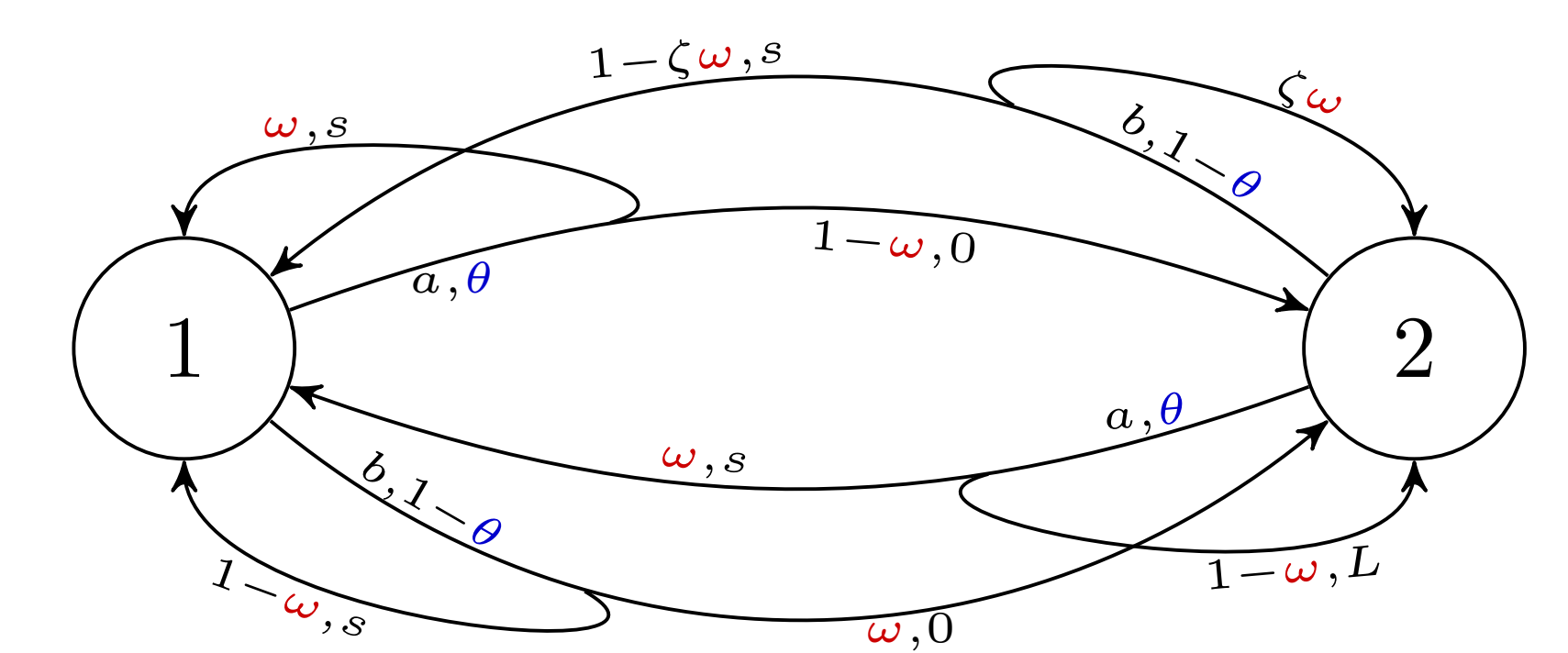


TORCS

- Configure the *front-rear wing orientation* and *brake repartition*



CHAIN DOMAIN



REFERENCES

- C. Cortes, S. Greenberg, and M. Mohri. Relative deviation learning bounds and generalization with unbounded loss functions. *arXiv preprint arXiv:1310.5796*, 2013.
- A. M. Metelli, M. Mutti, and M. Restelli. Configurable markov decision processes. In *35th International Conference on Machine Learning*, volume 80, pages 3491–3500. PMLR, 2018.
- A. B. Owen. *Monte Carlo theory, methods and examples*. 2013.
- J. Peters, K. Mülling, and Y. Altun. Relative entropy policy search. In *AAAI*, pages 1607–1612. Atlanta, 2010.
- M. L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.