

# Reinforcement Learning of Parameters in Complex Physical Systems

Nemo Fournier

September 4, 2019

# Outline

Introduction and Motivation

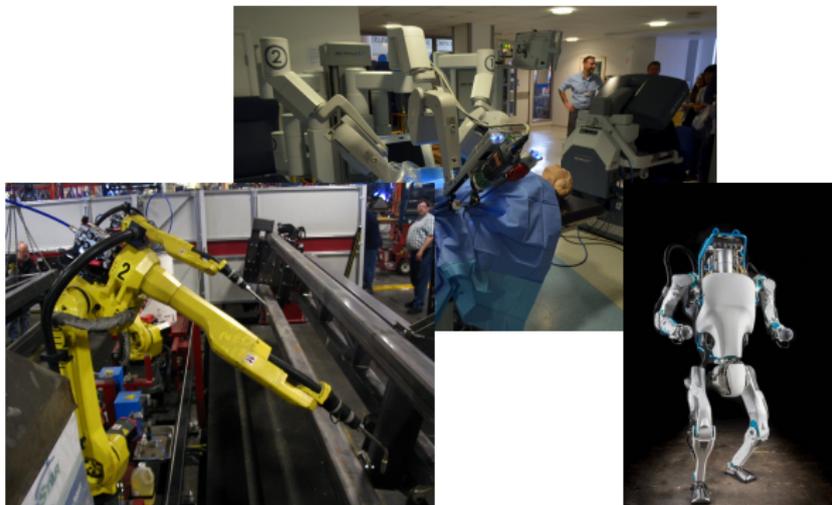
Reinforcement Learning

Robotics, RL and the Reality Gap

DR and UP

Embedding

# Progress in Robotic Hardware



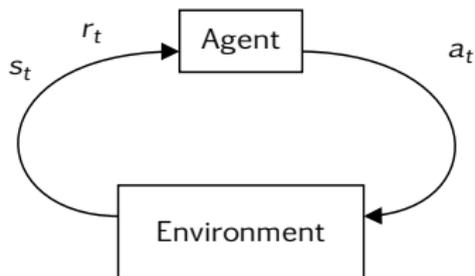
**Figure:** Up: Da Vinci surgical robot. Left: Fanuc welding robot. Right: Boston Dynamics' Atlas robot (*Images from Wikimedia*)

# Progress in Machine Learning



Figure: Up: AlphaGo match. Left: Dota2 AI. Right: Atari AI

# Main Principles of Reinforcement Learning



**Figure:** Reinforcement Learning (RL) feedback loop of the interactions between the agent and the environment.

# Some formalization

## Markov Decision Process

A MDP  $\mathcal{M}$  is a tuple  $(\mathcal{S}, \mathcal{A}, r, \gamma, p, p_0)$

- $\mathcal{S}$ : set of states
- $\mathcal{A}$ : set of actions
- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbf{R}$ : reward
- $\gamma$ : discount factor
- $p : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ : transition
- $p_0 \in \mathcal{P}(\mathcal{S})$ : initial state

## Some formalization

### Markov Decision Process

A MDP  $\mathcal{M}$  is a tuple  $(\mathcal{S}, \mathcal{A}, r, \gamma, p, p_0)$

- $\mathcal{S}$ : set of states
- $\mathcal{A}$ : set of actions
- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbf{R}$ : reward
- $\gamma$ : discount factor
- $p : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ : transition
- $p_0 \in \mathcal{P}(\mathcal{S})$ : initial state

### Policy

$\pi : \mathcal{S} \rightarrow \mathcal{A}$  or  $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$

# The RL problem

## Trajectory

$\tau = (s_0, a_0, s_1, a_1, s_2, \dots, s_a)$  is a trajectory over  $\mathcal{M}$  using a policy  $\pi$  if  $s_0 \sim p_0$ , and for  $t \geq 0$ ,  $a_t \sim \pi(\cdot | s_{t-1})$  and  $s_{t+1} \sim p(s_t, a_t)$ . We denote  $T_{\mathcal{M}, \pi}$  the distribution of such trajectories.

# The RL problem

## Trajectory

$\tau = (s_0, a_0, s_1, a_1, s_2, \dots, s_a)$  is a trajectory over  $\mathcal{M}$  using a policy  $\pi$  if  $s_0 \sim p_0$ , and for  $t \geq 0$ ,  $a_t \sim \pi(\cdot | s_{t-1})$  and  $s_{t+1} \sim p(s_t, a_t)$ . We denote  $T_{\mathcal{M}, \pi}$  the distribution of such trajectories.

## Performance of a policy

$$J(\pi) = \mathbb{E}_{\tau \sim T_{\mathcal{M}, \pi}} \left[ \sum_{t=0}^H \gamma^t r(s_t, a_t) \right]$$

# The RL problem

## Trajectory

$\tau = (s_0, a_0, s_1, a_1, s_2, \dots, s_a)$  is a trajectory over  $\mathcal{M}$  using a policy  $\pi$  if  $s_0 \sim p_0$ , and for  $t \geq 0$ ,  $a_t \sim \pi(\cdot | s_{t-1})$  and  $s_{t+1} \sim p(s_t, a_t)$ . We denote  $T_{\mathcal{M}, \pi}$  the distribution of such trajectories.

## Performance of a policy

$$J(\pi) = \mathbb{E}_{\tau \sim T_{\mathcal{M}, \pi}} \left[ \sum_{t=0}^H \gamma^t r(s_t, a_t) \right]$$

$$\pi^* = \arg \max_{\pi} J(\pi)$$

## Solving the RL problem: Policy Gradient Method

Idea: parametrize a policy  $\pi_\theta$  and perform gradient ascent:

$$\theta_{t+1} \leftarrow \theta_t + \alpha \nabla J(\theta_t)$$

## Solving the RL problem: Policy Gradient Method

Idea: parametrize a policy  $\pi_{\theta}$  and perform gradient ascent:

$$\theta_{t+1} \leftarrow \theta_t + \alpha \nabla J(\theta_t)$$

- REINFORCE
- Actor-Critic
- DPG
- TRPO
- PPO

# A Robotics Problem is a RL problem

- $\mathcal{S}$
- $\mathcal{A}$
- $p : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$

# A Robotics Problem is a RL problem

- $\mathcal{S}$
- $\mathcal{A}$
- $p : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$
- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbf{R}$

## Issues of RL when applied to robotics

- Sampling efficiency
- Random exploration
- Real-time rollouts

# Simulation

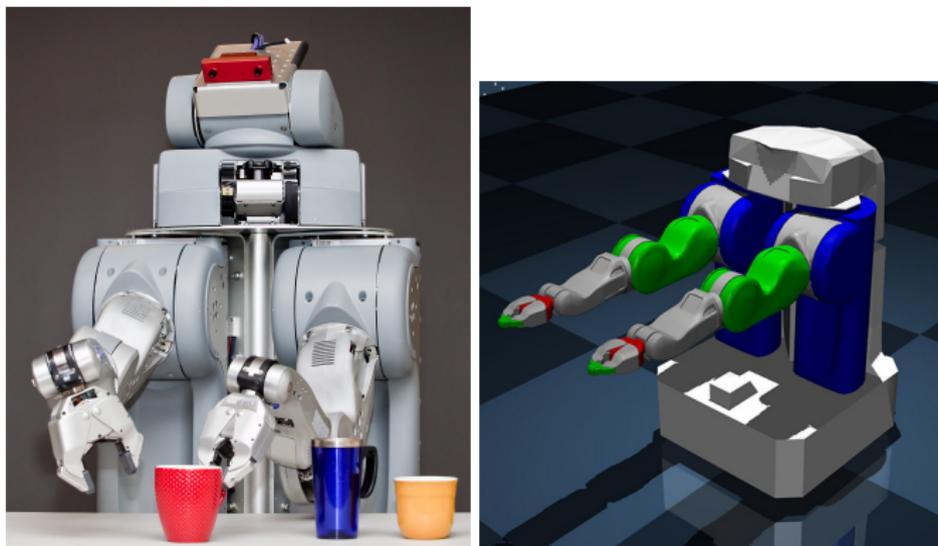


Figure: A real PR2 robot and its simulated equivalent.

# Strategies to Cross the Reality Gap

- Several learning phases
- Assess live discrepancies
- Dynamics randomization

# Dynamics Randomisation

Peng et al.<sup>1</sup> introduced parametrization of the environment using a vector  $\phi$ .

---

<sup>1</sup>X.B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel. *Sim-to-real transfer of robotic control with dynamics randomization*. 2018

# Dynamics Randomisation

Peng et al.<sup>1</sup> introduced parametrization of the environment using a vector  $\phi$ .

$$J_{\phi}(\pi) = \mathbb{E}_{\tau \sim \mathcal{T}_{\mathcal{M}, \pi, \phi}} \sum_{t=0}^H \gamma^t r(s_t, a_t)$$

---

<sup>1</sup>X.B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel. *Sim-to-real transfer of robotic control with dynamics randomization*. 2018

# Dynamics Randomisation

Peng et al.<sup>1</sup> introduced parametrization of the environment using a vector  $\phi$ .

$$J_{\phi}(\pi) = \mathbb{E}_{\tau \sim \Gamma_{\mathcal{M}, \pi, \phi}} \sum_{t=0}^H \gamma^t r(s_t, a_t)$$

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\phi \sim \rho} [J_{\phi}(\pi)]$$

---

<sup>1</sup>X.B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel. *Sim-to-real transfer of robotic control with dynamics randomization*. 2018

## Universal policy

Yu et al.<sup>2</sup> introduced the parametrized policy  $\pi_\phi = \pi(\cdot | \phi)$ .

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\phi \sim \rho} [J_\phi(\pi_\phi)]$$

---

<sup>2</sup>W. Yu, J. Tan, C. K. Liu, and G. Turk. *Preparing for the unknown: Learning a universal policy with online system identification*. 2017

## Universal policy

Yu et al.<sup>2</sup> introduced the parametrized policy  $\pi_\phi = \pi(\cdot | \phi)$ .

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\phi \sim \rho} [J_\phi(\pi_\phi)]$$

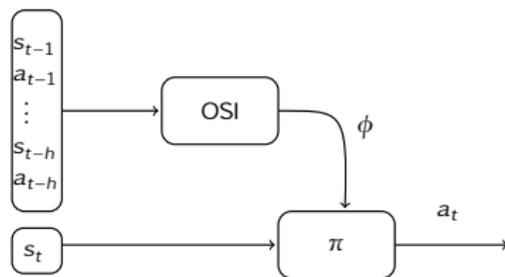


Figure: Universal Policy with Online System Identification

---

<sup>2</sup>W. Yu, J. Tan, C. K. Liu, and G. Turk. *Preparing for the unknown: Learning a universal policy with online system identification*. 2017

## New Issues of those two Methods (DR and UP-OSI)

- Sampling efficiency

## New Issues of those two Methods (DR and UP-OSI)

- Sampling efficiency
- Curse of dimensionality

## New Issues of those two Methods (DR and UP-OSI)

- Sampling efficiency
- Curse of dimensionality
- Choice of relevant parameters

# Dimensionality Reduction

We want a mapping

$$\Psi : \Phi \rightarrow \tilde{\Phi}$$

# Dimensionality Reduction

We want a mapping

$$\Psi : \Phi \rightarrow \tilde{\Phi}$$

$$\pi^*, \Psi^* = \arg \max_{\pi, \Psi} \mathbb{E}_{\phi \sim \rho} [J_{\phi}(\pi_{\Psi}(\phi))]$$

# Autoencoders

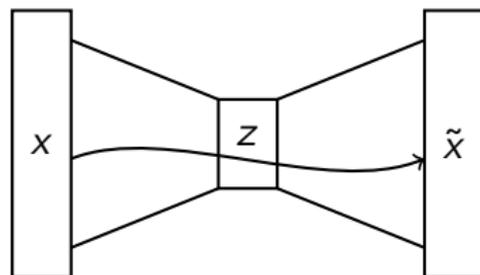


Figure: Standard autoencoder representation

## Our architecture

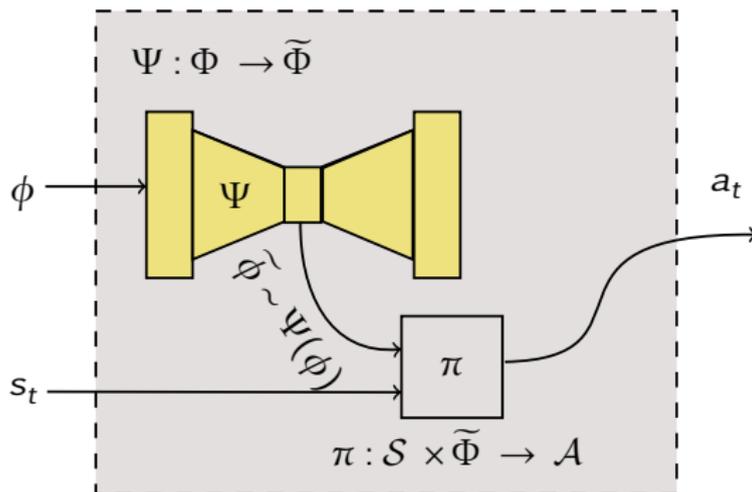


Figure: The new architecture we proposed.

# Analysing the embedding

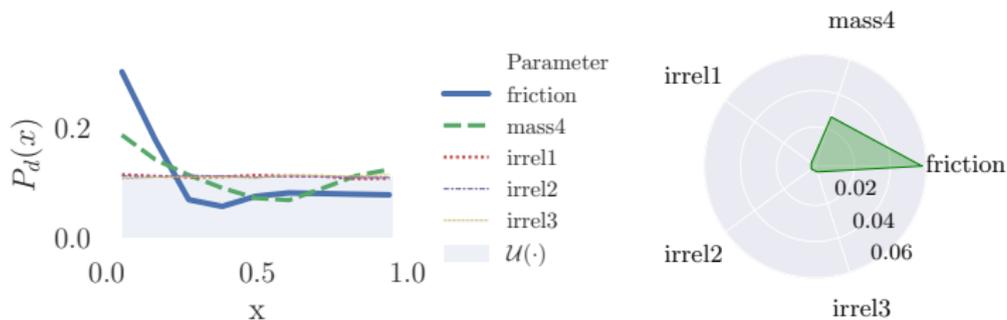
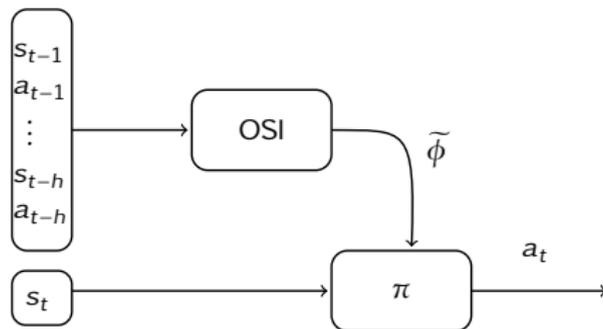


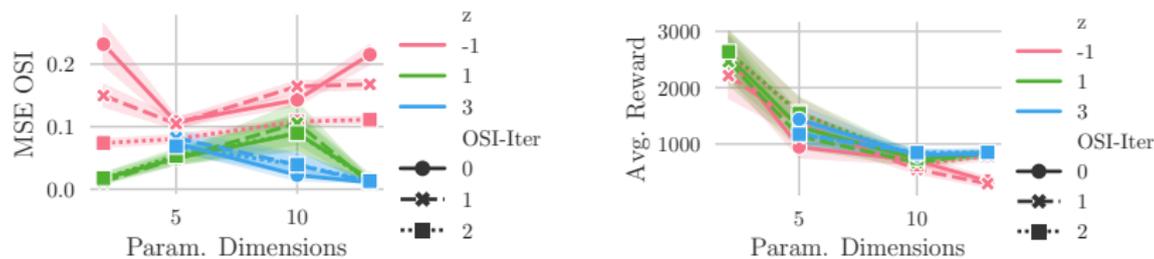
Figure: Toy problem on the Hopper environment.

# Training the embedded OSI



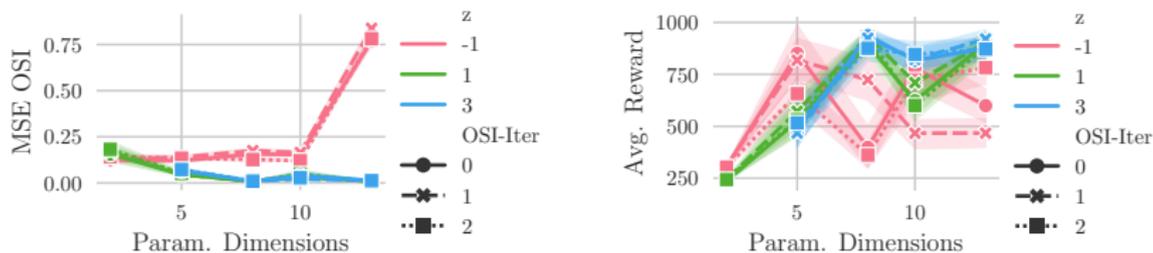
**Figure:** Embedded Universal Policy with Embedded Online System Identification

# Results



**Figure:** Effect of the embedding in terms of (Left) OSI prediction error and (Right)

# Transferability



**Figure:** Effect of the embedding for transfer in terms of (Left) OSI prediction error and (Right) Average reward on the task.

# Conclusion

- Promising direction and results
- Better evaluation needs to be done