

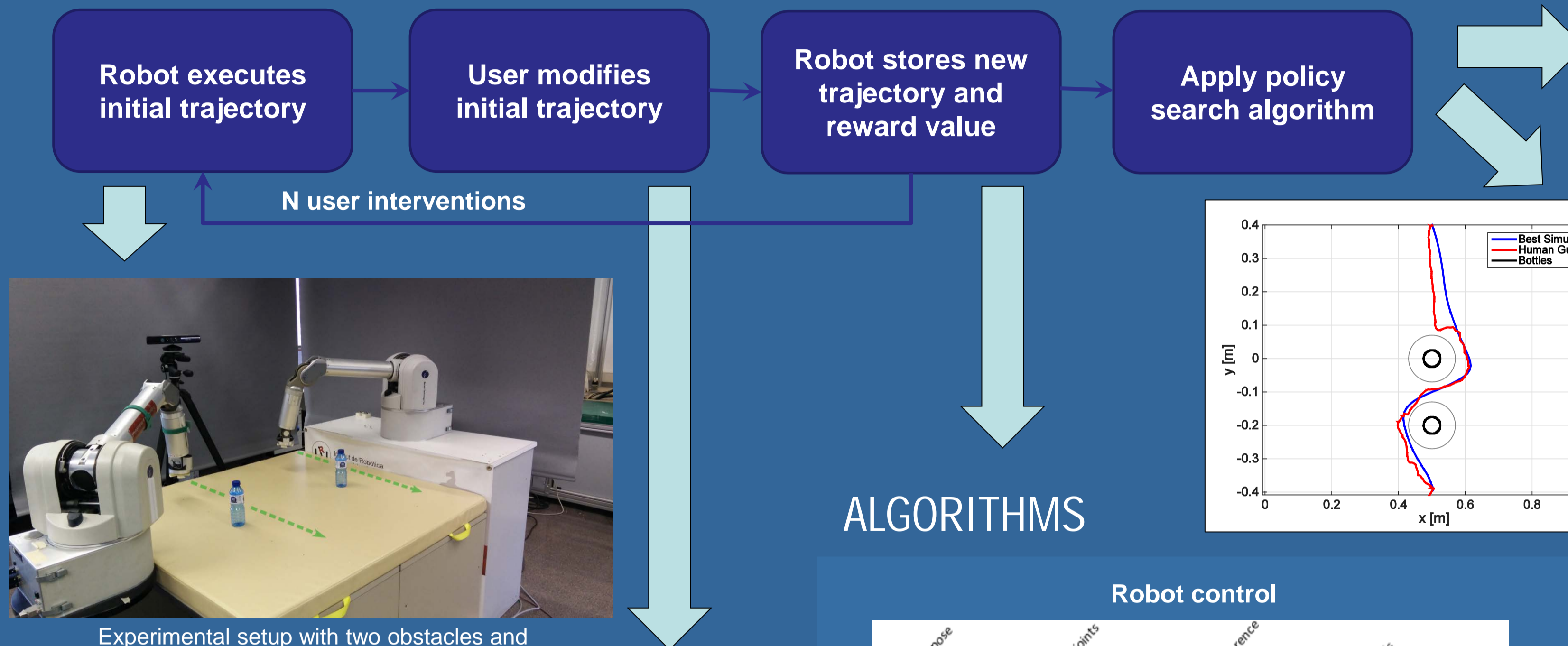
Learning Robot Motion through User Intervention and Policy Search

Aleksandar Jevtić, Adrià Colomé, Guillem Alenyà and Carme Torras
Institut de Robòtica i Informàtica industrial, CSIC-UPC, Barcelona, Spain

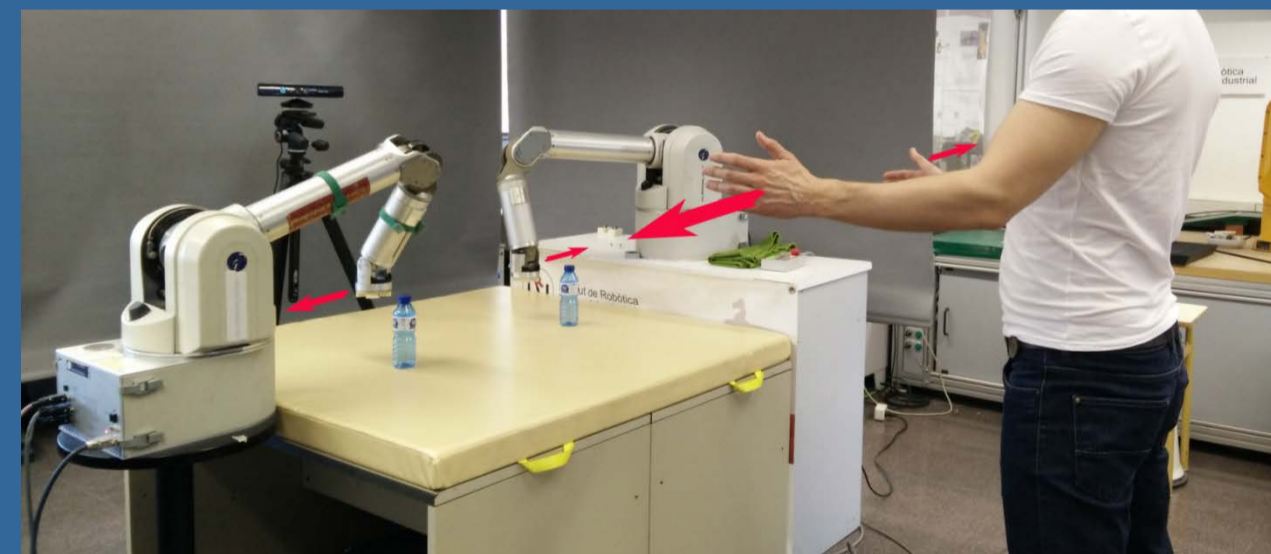
ABSTRACT

Reinforcement Learning (RL) through policy search is well suited for robot learning of manipulation tasks but requires a high number of policy updates to produce a satisfactory robot behavior. We propose a method that exploits **user intervention as a strategy for policy search updates**, which significantly **reduces the number of necessary updates** and allows **users to tailor robot behavior** according to their needs. We developed an interface for **gesture recognition and hand motion following** that allows a user to select and modify a desired manipulation task segment and speed up the learning process through a series of interventions. We implemented the method on a **single-arm and dual-arm robot** and applied it to a manipulation task on the laboratory tabletop. We compared our method with the conventional policy search strategy in simulation; the results show that the number of necessary policy search updates was significantly reduced by user intervention.

PROPOSED INTERACTIVE LEARNING FRAMEWORK

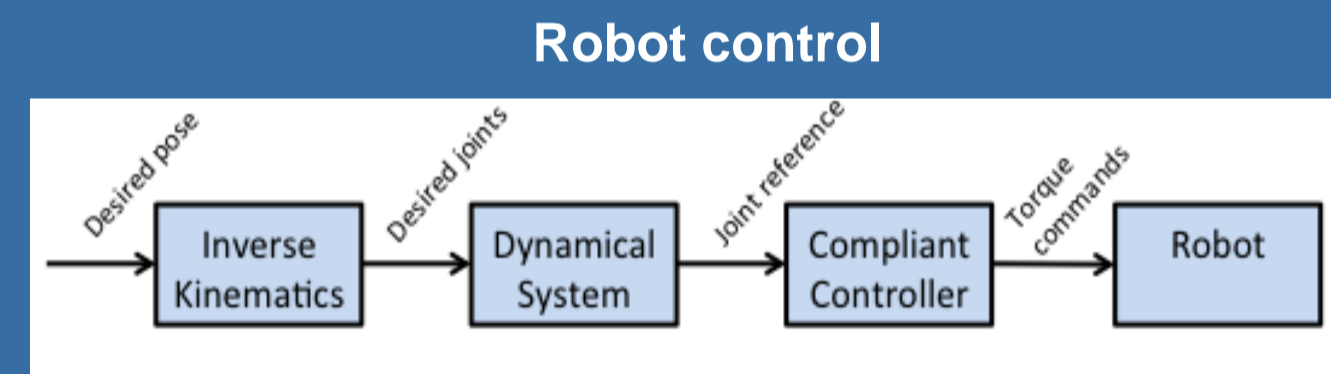


Experimental setup with two obstacles and presentation of the initial robots' trajectories



User intervention using hand gestures to guide the robots around the obstacles

ALGORITHMS



User intervention algorithm and reinforcement learning

Algorithm 2 Reinforcement learning of new trajectories

- 1: Input: User-generated trajectories $\tau_k = \{x_k^1, \dots, x_k^{N_t}\}$, for $\forall k$; cut and connect points indexes I_k^{cut}, I_k^{conn} , for $\forall k$; rewards of user-generated trajectories R_k , for $\forall k$
- 2: Find common unmodified part of the initial trajectory for all rollouts: $I^{cut} = \min_{k=1..N_{rollouts}} (I_k^{cut}), I^{conn} = \max_{k=1..N_{rollouts}} (I_k^{conn})$
- 3: Extract segment of trajectory that will be modified by RL, for all rollouts: $T_k = \{x_k^{I^{cut}}, \dots, x_k^{I^{conn}}\}$, for $\forall k$
- 4: Apply time alignment to T_k , for $\forall k$, using trajectory from the first rollout as reference; resample to have same number of points in all trajectories
- 5: Apply REPS to compute relative importance weight for each rollout: $w_k = \text{REPS}(R)$
- 6: Obtain new trajectory by a weighted sum of the time-aligned, resampled rollouts:

$$x_{new}^t = \sum_{k=1}^{N_{rollouts}} w_k x_k^t, \forall t = I^{cut} .. I^{conn}$$

- 7: Merge unmodified trajectory segments from step 2, $\{x^1, \dots, x^{I^{cut}-1}\}, \{x^{I^{conn}+1}, \dots, x^{N_t}\}$ same for $\forall k$, with modified trajectory segment from step 6 to obtain the RL-modified trajectory:

$$\tau_{new} = \{x^1, \dots, x^{I^{cut}-1}, x_{new}^{I^{cut}}, \dots, x_{new}^{I^{conn}}, x_{new}^{I^{conn}+1}, \dots, x^{N_t}\}$$

Algorithm 1 User intervention

- 1: set state to "EXECUTE";
- 2: go to first trajectory point;
- 3: **while** not end of trajectory **do**
- 4: go to next trajectory point;
- 5: **if** hand raised **then**
- 6: set state to "FOLLOW";
- 7: store trajectory cut point;
- 8: **while** in "FOLLOW" state **do**
- 9: follow user's hands;
- 10: **if** hand raised **then**
- 11: set state to "EXECUTE";
- 12: **end if**
- 13: **end while**
- 14: find closest trajectory connect point;
- 15: **end if**
- 16: **end while**

Static gesture recognition



Hand motion following

New robot position is obtained as follows:

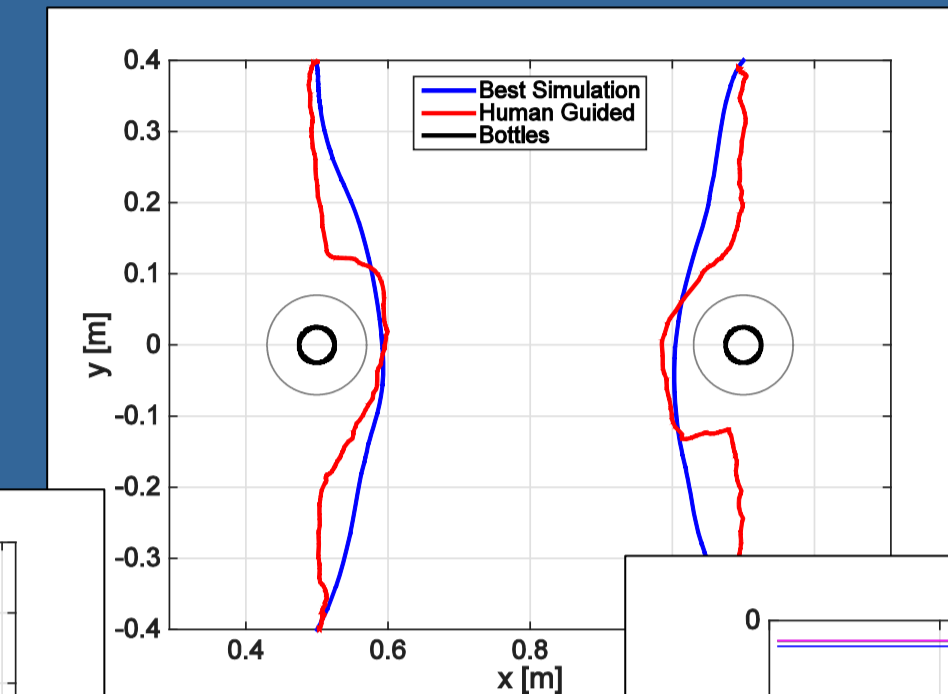
$$p_{goal} = p_{cut} + (p_{hand} - p_0), p \in \mathbb{R}^2$$

p_{cut} – starting point of a robot's end-effector during the user's intervention.

CONTACT

Aleksandar Jevtić
Institut de Robòtica i Informàtica industrial, CSIC-UPC
08028 Barcelona, Spain
Email: ajevtic@iri.upc.edu
Website: www.iri.upc.edu

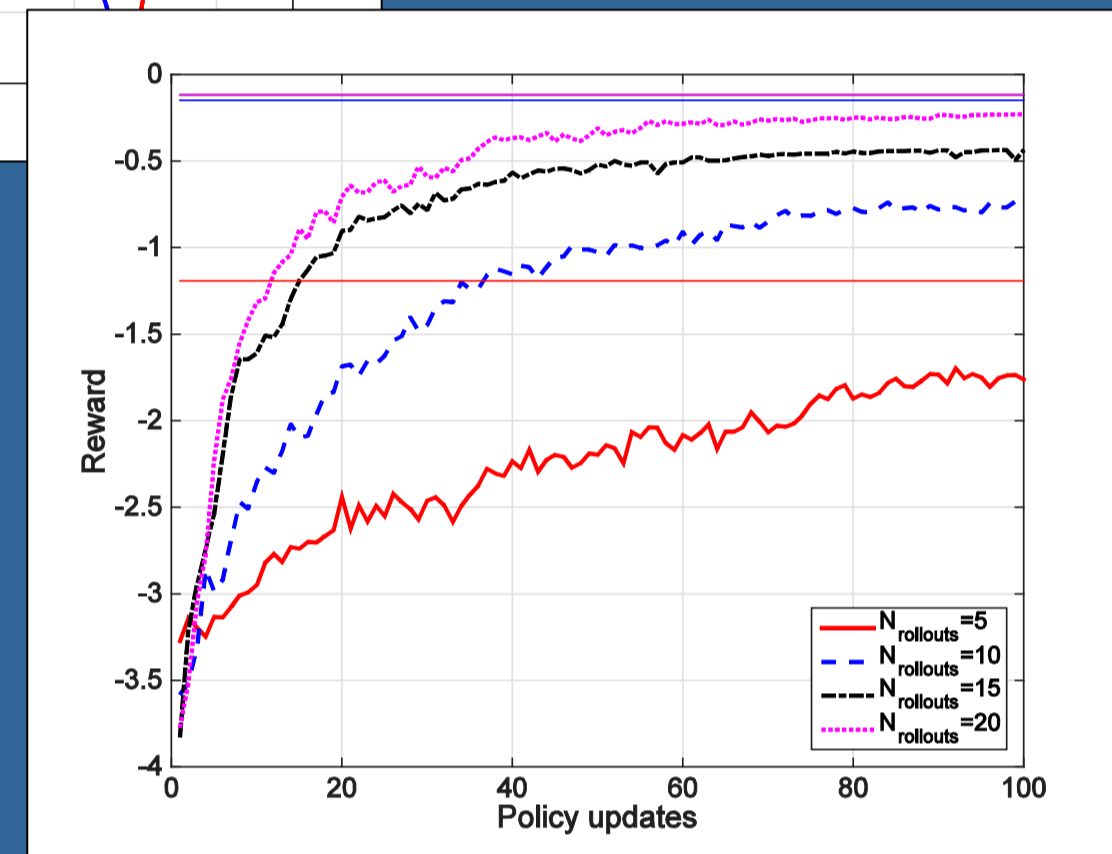
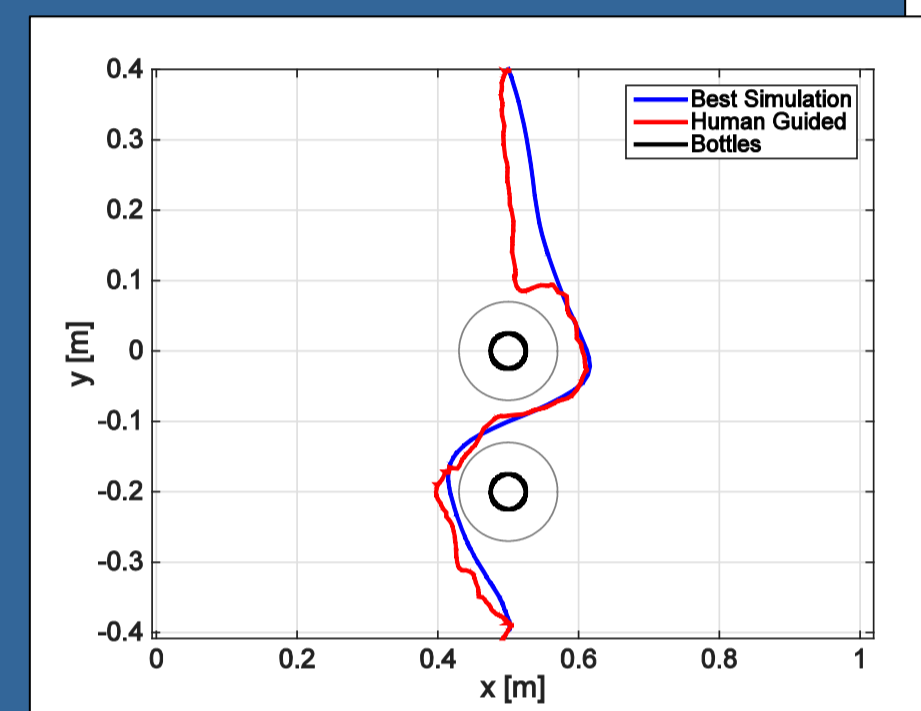
RESULTS



Reward function

$$R = -\alpha \cdot D_{total} - B$$

B is the number of knocked-down obstacles, D_{total} trajectory length, and α is a relative weight.



REPS algorithm uses Kullback-Liebler (KL) divergence as an indicator of difference between two probability distributions p, q over a random variable x :

$$KL(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

Given the previous policy $q(\theta)$, the new policy $\pi(\theta)$ is obtained by adding a KL-Divergence bound ϵ between the new and old policies to the optimization of the expected reward:

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \int \pi(\theta) R(\theta) d\theta \text{ s.t. } \epsilon \geq KL(\pi(\theta)||q(\theta))$$

$$1 = \int \pi(\theta) d\theta$$

where Θ are the trajectory parameters, $R(\theta)$ is their associated reward, and $\pi(\theta)$ is the probability, according to the policy π , of having such parameters.

REFERENCES

1. M. Deisenroth, G. Neumann and J. Peters, A survey on policy search for robotics, Foundations and Trends in Robotics, vol. 2, no. 1-2, pp. 1-142, 2013.
2. A. Jevtić, G. Doisy, Y. Parmet and Y. Edan, Comparison of interaction modalities for mobile indoor robot guidance: Direct physical interaction, person following, and pointing control, IEEE Transactions on Human-Machine Systems, vol. 45, no. 6, pp. 653-663, 2015.
3. A. Colomé, A. Planells and C. Torras, A friction-model-based framework for reinforcement learning of robotic tasks in non-rigid environments, 2015 IEEE International Conference on Robotics and Automation, pp. 5649-5654, 2015.
4. J. Peters, K. Mülling and Y. Altün, Relative entropy policy search, Twenty-Fourth National Conference on Artificial Intelligence, pp. 182-189, 2011.