

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier 10.1109/ACCESS.2017.DOI

# Crossing the Reality Gap: a Survey on Sim-to-Real Transferability of Robot Controllers in Reinforcement Learning

ERICA SALVATO<sup>1</sup>, GIANFRANCO FENU<sup>1</sup>, ERIC MEDVET<sup>1</sup> and FELICE ANDREA PELLEGRINO<sup>1</sup>

<sup>1</sup>Department of Engineering and Architecture, University of Trieste, Italy, IT (e-mail: erica.salvato@phd.units.it, fenu,emedvet,fapellegrino@units.it)

Corresponding author: Erica Salvato (e-mail: erica.salvato@phd.units.it).

**ABSTRACT** The growing demand for robots able to act autonomously in complex scenarios has widely accelerated the introduction of Reinforcement Learning (RL) in robots control applications. However, the *trial and error* intrinsic nature of RL may result in long training time on real robots and, moreover, it may lead to dangerous outcomes. While simulators are useful tools to accelerate RL training and to ensure safety, they often are provided only with an approximated model of robot dynamics and of its interaction with the surrounding environment, thus resulting in what is called the *reality gap* (RG): a mismatch of simulated and real control-law performances caused by the inaccurate representation of the real environment in simulation. The most undesirable result occurs when the controller learnt in simulation fails the task on the real robot, thus resulting in an unsuccessful *sim-to-real* transfer. The goal of the present survey is threefold: (1) to identify the main approaches to face the RG problem in the context of robot control with RL, (2) to point out their shortcomings, and (3) to outline new potential research areas.

**INDEX TERMS** Reality Gap, Reinforcement Learning, Robotics, Sim-to-real

## I. INTRODUCTION

**R**EINFORCEMENT Learning (RL) [1] allows to design controllers (often referred to as *agents*) with the capability of learning an optimal behaviour by interacting with the *environment*. The behaviour is defined in terms of state-action pairs, also known as *policy*, which is learnt through a *trial and error* process.

In some well-known RL tasks, such as pole-balancing, grid-search, or mountain car, state and action spaces of the system are small enough to allow approximating policies through tables [1], i.e., actions and states can be dealt with as finite discrete variables. However, the higher the complexity of the system to control, the more ineffective the tabular approaches become [1]. Indeed an increase in system complexity is often related to an increase of the state and action spaces dimensions, which makes a tabular approach intractable. In challenging cases, such as robot control, treating state and action as continuous variables in a compact set is a more appropriate way to deal with the problem [2, 3]. For this purpose, approximators of the policy or of some supporting element, such as value function, or of both, are required [1]. When Deep Neural Networks (DNNs)

are employed as approximators, the approach is referred to as Deep Reinforcement Learning (DRL); it allows to develop RL controllers with less manual feature engineering than classic tools (radial basis functions, tile coding, etc.) [4, 5, 6]. On the other hand, when DRL is directly employed on the robot in real-time, it results in considerably long training times. Moreover, due to the intrinsic trial and error nature of RL, a real-world training, in particular during the exploration phase of the state and action space, can lead to unsafe actions of the robot. Therefore, a way to train robots safely and quickly is needed.

Simulators allow to easily address these problems once they are provided with a model of the robot dynamics able to replicate the actual behaviour as closely as possible. In principle, simulators allow to train the controller with faster and safer procedures: once the policy has been learnt, it is transferred to the real system (*sim-to-real* transfer) [7]. However, sim-to-real transfer is only effective when the simulator is given a sufficiently accurate model of the real robot and the environment [8]; unfortunately, the more accurate the simulation, the heavier the computational cost. A less accurate simulator is therefore often preferred, although it may result

in a less effective sim-to-real transfer. The phenomenon in which a controller learnt on simulator degrades once applied on the real world is the so-called *reality gap* (RG) [9]. In the worst case, the RG leads to a failure of the policy when applied on the real world, which means a robot unable to achieve its goal.

RL is not the only approach that can be affected by the RG. Any technique in which the controller design relies on a simulator of the real system can potentially exhibit a reality gap [10]. Indeed, several works faced the RG problem in other frameworks, such as Evolutionary Computation [9, 11, 12, 13, 14, 15, 16, 17] or Model Predictive Control [18, 19, 20, 21].

However, here we focus our attention only on those works facing the RG problem on robot controllers learnt with RL. Most of the many solutions proposed in the literature, are task-dependent and/or have been tested on a specific task only. The outcomes are that: (a) generalisation is not ensured, (b) and a comparison between different approaches is not feasible.

Although a sketch of the current state of the art is already proposed in [22], here we conduct a more in-depth analysis. We introduce the main concept behind the RG in a general RL framework and we survey relevant and recent literature concerning RG in the context of robot control with RL. According to our analysis, the approaches for coping with RG in this context fall into three broad categories: domain randomisation (DR), adversarial reinforcement learning (ARL), and transfer learning (TL). Hence, we introduce a general RL framework suitable to be specialised for each of the mentioned approaches.

The aims of the present work are:

- (1) to provide a systematic picture of the literature concerning how to solve the RG problem in robot control tasks with RL;
- (2) to clarify the differences between the three main identified approaches by highlighting the relative pros and cons;
- (3) to identify new possible research areas.

The remainder of the paper is organised as follows. In Section II we describe a formal framework for RL useful to better explain the key elements of DR, ARL, and TL. In Section III we introduce RL as a means for robot control and discuss the RG in the context of robotics. In Section IV we survey and discuss the current methodologies to face the RG, relying, where possible, on the provided formalism. Finally, in Section V, we draw the conclusions.

## II. REINFORCEMENT LEARNING

Reinforcement Learning (RL) can be employed to perform optimal data-driven control without the need to rely on a mathematical model of system dynamics [1, 23, 24].

It is typically described as a Markov Decision Process (MDP); i.e., a tuple  $(X, A, p, r)$  defined by the state set  $X$ , the control set  $A$ , the transition probability  $p$ , and the immediate reward  $r \in \mathbb{R}$ . In brief, it expresses a discrete-time

stochastic control process in which a scalar  $r$  is employed to assess the quality of the controller choice in terms of task achievement.

Here, we adopt a control systems-oriented formalism similar to the one employed in [25]. However, what follows also applies to MDP settings by defining some of the following elements in terms of expected values.

### A. ADOPTED FORMALISM

A *dynamical discrete-time system*  $\Omega$  is a tuple  $(X, A, O, f, g)$  composed of the state set  $X$ , the control set  $A$ , the observation set  $O$ , a transition function  $f : X \times A \rightarrow X$ , and an observation function  $g : X \rightarrow O$ . Let  $x^{(k)} \in X$ ,  $a^{(k)} \in A$ , and  $o^{(k)} \in O$  be the state, the applied control input, and the observation, respectively, at the  $k$ -th time-instant. A dynamical discrete-time system, starting from an initial state  $x^{(0)}$  and subjected to a control sequence  $a^{(0)}, a^{(1)}, \dots$ , evolves according to the following laws:

$$x^{(k+1)} = f(x^{(k)}, a^{(k)}) \quad (1)$$

$$o^{(k+1)} = g(x^{(k+1)}) = g(f(x^{(k)}, a^{(k)})). \quad (2)$$

An *environment*  $E$  is a dynamical discrete-time system described by a tuple  $(X, A, O, f, g, h)$  composed of the same elements of  $\Omega$  plus a *reward function*  $h : X \times A \rightarrow \mathbb{R}$ . Let  $r^{(k+1)} \in \mathbb{R}$  be the reward at the  $(k+1)$ -th time-instant. An environment  $E$ , starting from an initial state  $x^{(0)}$  and subjected to a control sequence  $a^{(0)}, a^{(1)}, \dots$ , evolves according to Equations (1) and (2) and:

$$r^{(k+1)} = h(x^{(k)}, a^{(k)}). \quad (3)$$

Overall, we assume that  $g(f(\cdot)) = f(\cdot)$ , thus resulting in  $o^{(k+1)} = x^{(k+1)}$ . The following is also pertinent, with appropriate adjustments, in case of  $g(f(\cdot)) \neq f(\cdot)$  that simply means dealing with a partial observability of the state. In the MDP setting this amounts to employing a Partially Observable Markov Decision Process (POMDP), which requires to introduce the concept of *belief* [26]. However, it is not necessary for the purpose of the present work.

A controller (or *policy*) for a dynamical discrete-time system, and therefore also for an environment, is a function  $\pi : O \rightarrow A$ .

Given a discount factor  $\gamma \in [0, 1]$ , an *optimal policy*  $\pi^*$  is a policy that satisfies, for any initial state  $x^{(0)}$ :

$$\begin{aligned} \pi^* &= \arg \max_{\pi \in \Pi} \sum_{k=0}^{+\infty} \gamma^k h(x^{(k)}, \pi(o^{(k)})) \\ &= \arg \max_{\pi \in \Pi} J_\pi(x^{(0)}), \end{aligned} \quad (4)$$

where  $\Pi$  is the set of policies, while  $J_\pi(x^{(0)})$  is the *infinite horizon discounted reward* starting from  $x^{(0)}$  under the policy  $\pi$ .

We denote with  $\pi \leftrightarrow E$  the closed-loop system where  $\pi$  determines the control input to be applied on  $E$ , represented in Figure 1. In particular, according to the above definitions

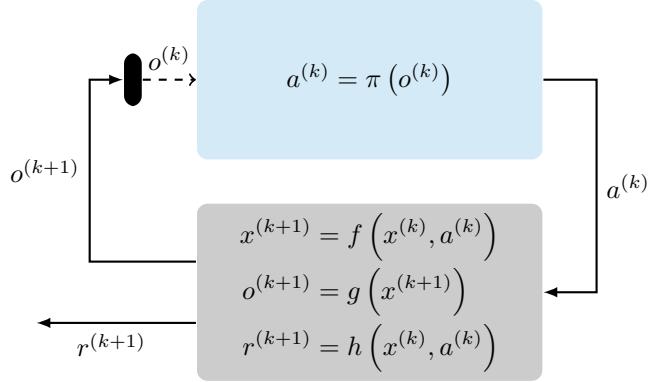


FIGURE 1: The closed-loop system  $\pi \leftrightarrow E$  where a policy (light blue block above) applies a control input  $a^{(k)}$  to an environment (gray block below) that outputs an observation  $o^{(k+1)}$  and a reward  $r^{(k+1)}$ .

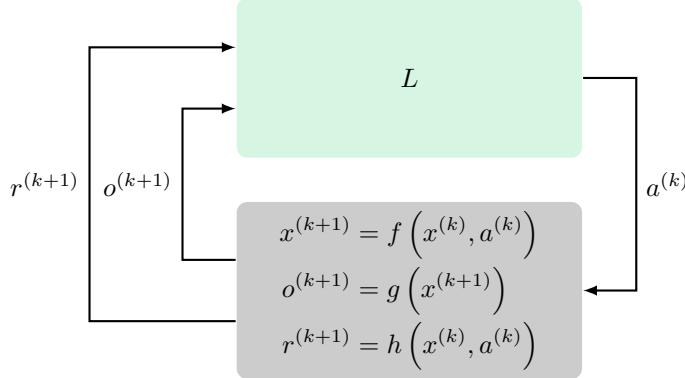


FIGURE 2: The closed-loop system  $L \leftrightarrow E$  where a RL-based agent  $L$  (green block above) applies a control input  $a^{(k)}$  to an environment (gray block below) that outputs an observation  $o^{(k+1)}$  and a reward  $r^{(k+1)}$ .

and to the assumption that  $g(f(\cdot)) = f(\cdot)$ , the scheme represents a *state feedback* control. Clearly, when  $g(f(\cdot)) \neq f(\cdot)$ , the same scheme could represent an *output feedback* control.

A *policy learning algorithm*  $L$  is an algorithm that, given an environment  $E = (X, A, O, f, g, h)$  and a *learnable policies set*  $\Pi^L \subseteq \Pi$ , outputs (learns) a policy  $\pi^L = L(E) \in \Pi^L$ . When the learning encompasses an interaction with the environment, the policy learning algorithm could be seen as an *agent*  $L$  that learns a controller  $\pi$  by interacting with an environment  $E$ . We denote by  $L \leftrightarrow E$  the resulting closed loop system (Figure 2). A policy learning algorithm is said to be *model-based*, if it relies on the knowledge of  $f$  and  $g$  (either known or identified based on collected data), or *model-free*, otherwise.

A typical RL agent can be seen as a policy learning algorithm  $L$ . Indeed, during training, it interacts with the environment  $E = (X, A, O, f, g, h)$  and updates a controller  $\pi$  by observing the consequences (in terms of reward  $r$ ) of selected control inputs. Its global goal is to learn a controller  $\pi^*$  in  $E$  which satisfies (4).

The learning procedure can either be *episodic* or in a single chapter (*non-episodic*): in the former case, the learning is performed through episodes and the state is reset in case of a failure, a goal achievement, or the achievement of the maximum episode length  $T$ . In addition, the agent can update the controller either by using data from the current policy (*on-policy*) or independently of it (*off-policy*).

RL approaches can be categorised in three main categories [1]:

- *Value-function* approaches, based on the idea of the *value* of a state (value function  $V_\pi(x)$ ) or of a state-control input pair (action-value function  $Q_\pi(x, a)$ ). The value function and the action-value function represent, respectively, the cumulative reward obtained from  $x^{(k)}$  by applying  $\pi$  and the cumulative reward obtained from  $x^{(k)}$  by first choosing an  $a^{(k)}$  and then applying  $\pi$ . The optimal policy corresponds to the optimal  $V$  or  $Q$  functions:

$$V^*(x^{(k)}) = \max_{\pi} J_{\pi}(x^{(k)}) \quad (5)$$

$$Q^*(x^{(k)}, a^{(k)}) = \max_{\pi} h(x^{(k)}, a^{(k)}) + \gamma J_{\pi}(x^{(k+1)}). \quad (6)$$

In brief, by interacting with the environment and observing rewards, either  $V^*$  or  $Q^*$  are estimated, thus leading to the optimal policy.

- *Policy-search* approaches, in which a parametrized policy  $\pi_\theta$  is defined, whose parameters  $\theta$  are updated based on the observed reward in order to maximise  $J_{\pi_\theta}$ , either employing gradient-based or gradient-free optimisation techniques [27].
- *Actor-Critic* approaches, which integrate the idea of both previous categories.  $V_\pi(x)$ , in this case, is employed as a baseline (Critic) for policy gradient optimisation (Actor).

Further technicalities are not necessary for the purpose of this work. For additional details on RL we refer readers to [1, 25, 28, 29].

### III. REINFORCEMENT LEARNING IN ROBOTICS

The motivation for using RL in robotics is to make a robot autonomous in finding an optimal policy, through *trial and error* interactions with its environment, without an explicit knowledge of the model (model-free). However, as a matter of fact, the most effective methods to date are model-based [30, 31, 32]. In addition, policy search approaches result in more efficient trainings in terms of time needed for convergence [33, 34, 35, 36, 37, 38]. Regardless of the specific RL method being used, the application of RL enabled researchers and practitioners to face different significant robotic tasks (e.g., manipulation, navigation, motion control, etc.). In Table 1, we report the tasks that are more often considered, along with a few significant research papers, some of which (those dealing with sim-to-real transferability) are surveyed in the present study.

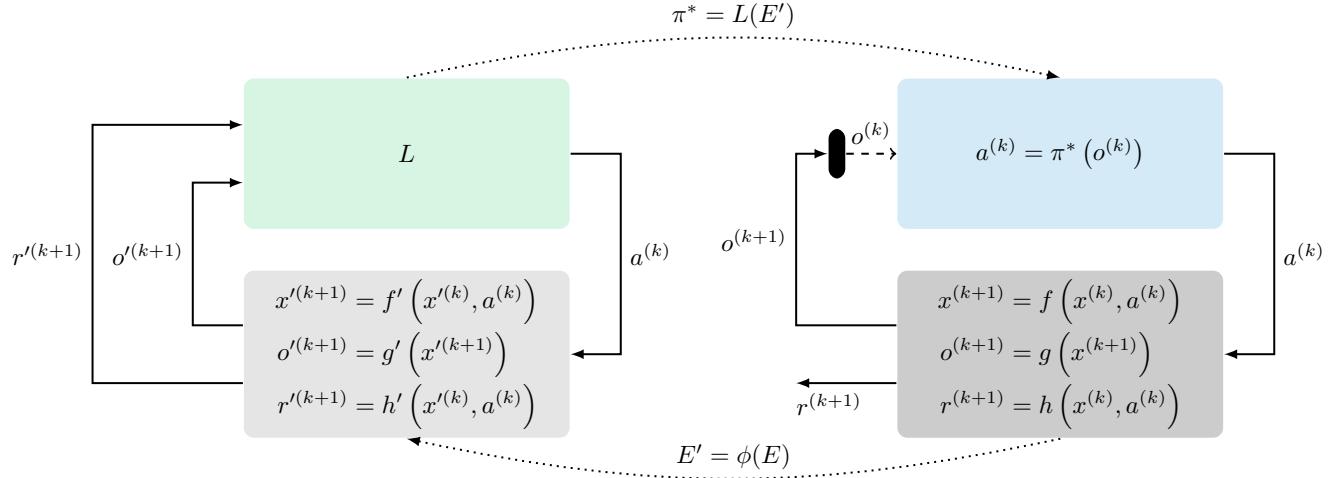


FIGURE 3: Schematic representation of the application of RL to the robot control problem using simulation. First, given the real robot and it environment  $E$  (block at bottom right), a simulator  $E'$  (block at bottom left) is obtained by modelling  $E$  using  $\phi$ . Then an agent  $L$  (block at top left) learns a policy  $\pi^*$  in simulation. Finally, the learnt policy is transferred to the real robot where it can be deployed (block at top right).

TABLE 1: Overview of the robotic tasks involving RL controllers.

Tasks	References
Navigation	[39], [40], [41], [42]
Manipulation	[43], [44], [45], [46], [47], [3], [8], [48]
Motion control	[49], [50], [51], [52], [53], [54], [55], [56], [57], [58], [59]
Locomotion	[49], [50], [51], [60], [61], [62]

When tackling a robot control problem with RL, according to the formulation and notation above provided, the designer should ensure that (a) the environment  $E$  captures the robot-surrounding environment dynamics; (b) the observations  $o \in O$  include proprioceptive measures useful to capture the robot dynamic evolution (e.g., joints position and velocity); (c) the control inputs  $a \in A$  correspond to values for the appropriate robot actuators (e.g., torques to be applied to joints or desired acceleration/speed/position, depending on the specific control system).

A first issue concerns the representation of  $V_\pi(x)$ , and  $Q_\pi(x, a)$ . While in some simple RL scenarios  $O$  and  $A$  sets can be discretised over a finite range, hence allowing a tabular representation of  $V_\pi(x)$  or  $Q_\pi(x, a)$ , in robot control problems the physical nature of observations and control inputs advocates for a finely discretised (ideally, continuous) representation of  $O$  and  $A$ . In this case, table may require a huge amount of memory and the problem gets practically intractable. An alternative and often suitable solution to escape tabular representation is function approximation, recently addressed by DRL. Here DNNs are used as function approximators of  $V_\pi(x)$ ,  $Q_\pi(x, a)$ , or directly  $\pi$ , and a loss

function  $\eta_\theta$  is designed in order to guide the training of the network itself. For a more detailed description of DRL, including recent efforts and some applications, we refer the reader to [4, 63].

A second issue, pivotal for the aim of the present work, derives from the trial and error process employed by RL and its direct application to the robot: learning time may become too long, thus unpractical, and the risk of damaging the robot, or more in general the environment, may be too high. Typically, simulators are used for addressing this issue; i.e., an *environment mapping operator*  $\phi$  is assumed to exist (albeit unknown in practice) such that  $E' = \phi(E) = (X', A', O', f', g', h')$  is a digital approximated copy of  $E$ . Intuitively,  $\phi$  corresponds to modelling a real system described by  $E$  as a simulated system described by  $E'$  such that a policy learnt on  $E'$  can be applied to  $E$ . The resulting controller design and test are, therefore, split in two distinct phases that are performed on two different environments: (i) an agent  $L$  interacts with a simulated environment  $E' = \phi(E)$  and outputs a controller  $\pi^*$ ; (ii) the resulting  $\pi^*$  is applied on the real  $E$ . In this scenario, outlined in Figure 3, the latter step, that in general can be defined as  $E'$ -to- $E$  transfer, can be renamed as *sim-to-real* transfer, given the nature of the considered  $E'$  and  $E$ .

As will be clear from the reviewed examples in the following sections, a controller learnt in simulation often exhibits performance losses when applied on the real robot and, in the worst scenario, totally fails the task. From the point of view of the designer, this issue becomes relevant when, although a safer and faster training and an effective test have been carried out on  $E'$ , the policy learnt in simulation does not achieve the goal in real world. In such a situation, the learning algorithm is affected by an intolerable RG and the corresponding learnt policy  $\pi^*$  is said to be non-transferable.

Note that, although it is not always clearly emphasised, a

mandatory step ahead of the sim-to-real transfer is to perform a test of the learnt policy on the simulator itself<sup>1</sup>. Otherwise, there is no guarantee of the learnt controller effectiveness in achieving the task even in simulation. For instance, a typical learning stop criterion consists in terminating the training when the moving average of the cumulative reward settles down. However, this empirical rule does not ensure that the learnt controller is able to correctly perform the task. The learning algorithm may have been trapped in a local maximum, and the resulting controller may have a completely unexpected behaviour in tests, even on the simulator.

If the test on the simulator is effective, possible reasons for a performance loss in a test on the real robot, and therefore for RG, are:

- (a)  $\phi$  operator is unrealistic:  $E$  and  $E'$  differences in  $f, f'$  and/or  $g, g'$  are such that, applying the same input on both environments, the resulting  $o$  and  $o'$  are different; or, possibly,  $E$  and  $E'$  differences in  $h$  and  $h'$  lead to two different optimisation problems  $J_\pi(x) \neq J_\pi(x')$ ;
- (b)  $L$  is unable to output a  $\pi^*$  sufficiently robust to possibly small and unavoidable errors in  $\phi$ .

However, the particular case in which  $h'$  is very different from  $h$  is unusual, in practice. Typically,  $h$  is properly designed and remains “the same” for both  $E$  and  $E'$  (for this reason hereinafter we consider  $h' = h$ ). Therefore, the RG can be essentially attributed to a mismatch between  $E$  and  $E'$  and the possible solutions may be: (a) improve  $E'$ , by properly adjusting  $\phi$  (not always possible because it could result in an excessive computational effort or because of lack of knowledge); (b) make the controller more robust to model errors.

Note that, in general, it is not required that the controller behaves identically when applied to  $E$  and  $E'$ , but, rather, that it is  $E'$ -to- $E$  transferable, i.e., sim-to-real transferable. In practice, this requirement translates to a (subjectively) properly bounded RG. Clearly, if an appropriate behaviour is only reached when  $E'$  is an identical copy of the robot  $E$  ( $E' = \phi(E) = E$ ), using  $E'$  rather than  $E$  has no benefit in reducing the overall learning time. But still,  $E'$  may be useful for addressing risks concerning safety.

Figure 4 summarises a generic routine for investigating the presence of RG. We start by (1) performing a training on  $E'$ ; (2) then we test the resulting policy  $\pi$  on the same  $E'$ , to check its effectiveness in reaching the task, and, only if the test ends successfully, (3) we test  $\pi$  also on  $E$ ; otherwise, we revert to (1). Once that the (3) has been executed, (4) we compare the performance obtained in the  $\pi$  tests on  $E$  and  $E'$ .

#### IV. METHODOLOGIES FOR SOLVING RG

In the present paper, we focus only on those articles that meet all the following requirements: (i) address explicitly the

<sup>1</sup>This only applies to those situations where the controller is transferred on the real robot. As shown in Section IV-C, in some scenarios, the training continues on the real robot. Hence a test procedure is not needed in these cases.

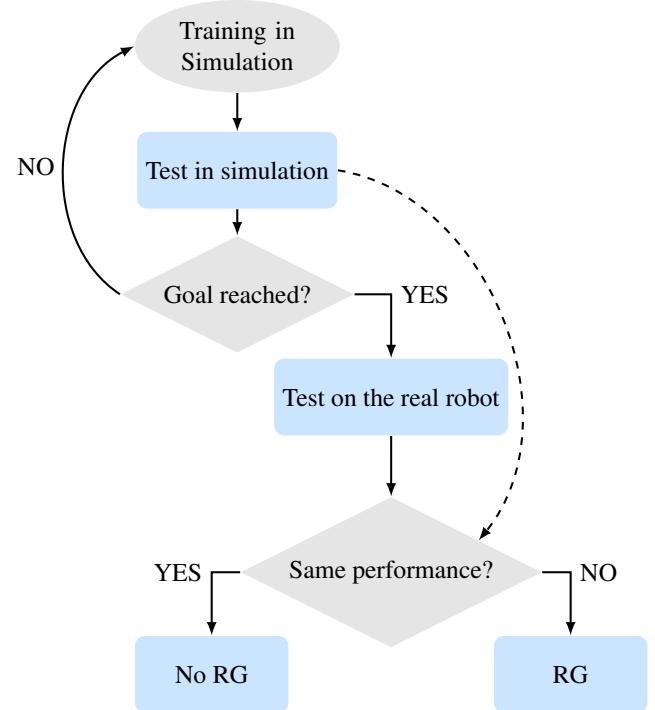


FIGURE 4: Flowchart of the routine to test the presence of the RG.

RG problem, (ii) deal with robot control applications, and (iii) employ RL techniques. We have identified three major categories of approaches for addressing the RG in this scenario: domain randomisation (DR), adversarial RL (ARL), transfer learning (TL). All the articles discussed below are summarised in Table 2 according to this categorisation. The table also shows, for each article, if the authors conducted experiments only in simulation (sim-to-sim) or (also) on real robot (sim-to-real), and specifies the employed simulators.

Due to task diversity (Figure 5 shows a visual summary of the robotic tasks) and the lack of a common theoretical framework for the RG, the surveyed articles do not present their results in a way that permit a systematic comparison. However, we provide a general formal definition for each of the previously mentioned categories, according to the formalism of Section II, which allows to understand each approach.

#### A. DOMAIN RANDOMISATION

Domain randomisation (DR) has already achieved good results in sim-to-real transfer of robotics controllers outside RL [54, 75, 76, 77, 78, 79]. The main idea behind this approach is what in control theory is called *robust control under either parametric or non parametric uncertainty* [80, 81], that is the design of controllers able to guarantee certain properties despite some tolerable parameters variations and/or noise.

We call  $\tilde{E}' = \tilde{\phi}(E)$  a corrupted simulator described by a tuple  $(\tilde{X}', \tilde{A}', \tilde{O}', Z', \Upsilon', \tilde{f}_{\xi'}, \tilde{g}_{\psi'}, h)$  in which  $Z'$  is the

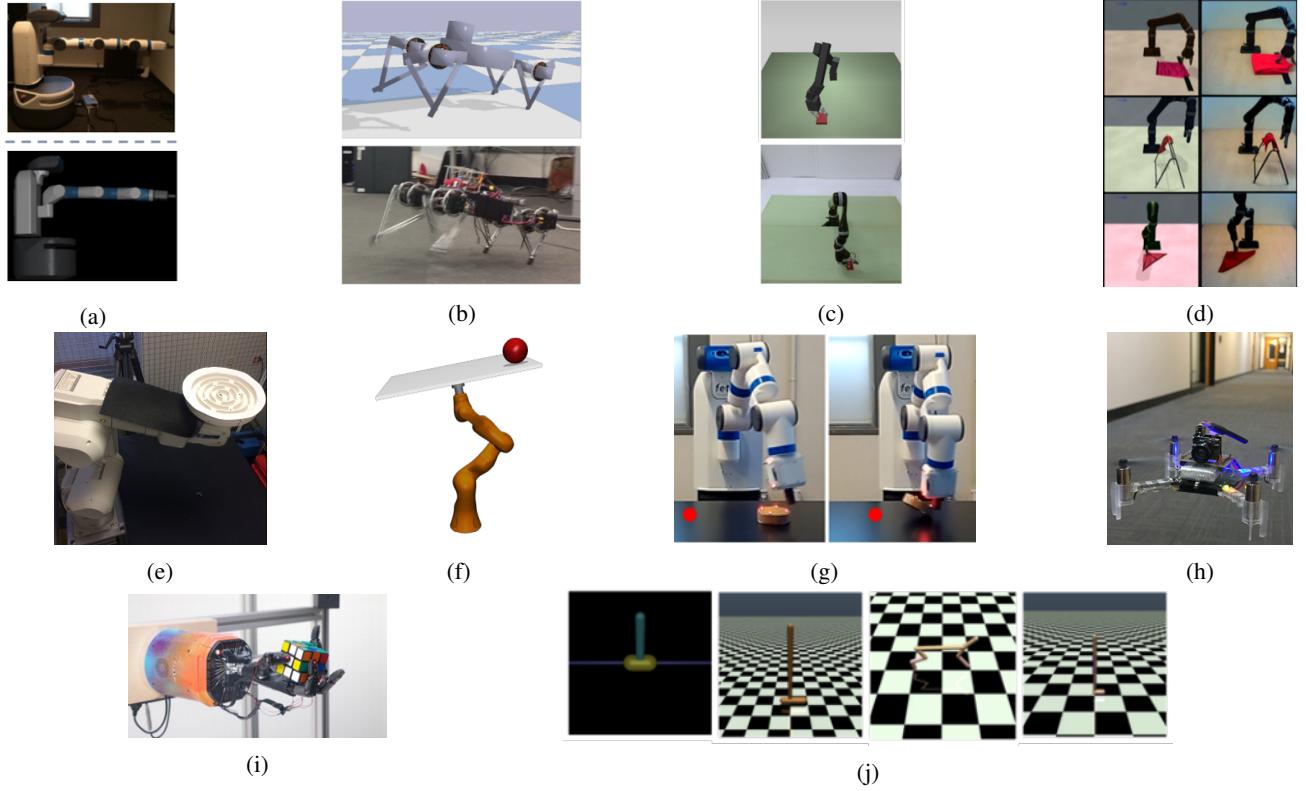


FIGURE 5: Overview of some robotic tasks considered in the surveyed articles to address the RG. (a) The Fetch robot used in [49]. (b) The Minitaur of [54]. (c) The robot employed for manipulation task in [74]. (d) The robotic arm engaged in deformable object manipulation of [45]. (e) The Marble maze game of [56]. (f) The ball on plate system used by [53]. (g) The Fetch robot used for the pushing task of [44]. (h) The quadrotor employed for the autonomous navigation task of [40]. (i) The five-finger humanoid hand used in [57]. (j) The classical Open AI Gym environment used to test in simulation several strategies [49, 50, 51, 52].

process disturbances set,  $\Upsilon'$  is the measurement disturbances set,  $f'_{\xi'} : \tilde{X}' \times \tilde{A}' \times Z' \rightarrow \tilde{X}'$  the corrupted and parametric transition function, with parameters  $\xi' \in \Xi'$ , and  $g'_{\psi'} : \tilde{X}' \rightarrow \tilde{O}'$  the corrupted and parametric observation function, with parameters  $\psi' \in \Psi'$ .

Given a parametrisation  $\xi', \psi'$ , starting from an initial state  $\tilde{x}'^{(0)}$  and subject to a control sequence  $\tilde{a}'^{(0)}, \tilde{a}'^{(1)}, \dots$ , a process disturbance sequence  $\zeta'^{(0)}, \zeta'^{(1)}, \dots$ , and a measurement disturbance sequence  $v'^{(0)}, v'^{(1)}, \dots$ , a corrupted simulator  $\tilde{E}'$  evolves according to:

$$\tilde{x}'^{(k+1)} = \tilde{f}'_{\xi'} \left( \tilde{x}'^{(k)}, \tilde{a}'^{(k)}, \zeta'^{(k)} | \xi' \right) \quad (7)$$

$$\tilde{o}'^{(k+1)} = \tilde{g}'_{\psi'} \left( \tilde{x}'^{(k+1)}, v'^{(k)} | \psi' \right) \quad (8)$$

$$\tilde{r}'^{(k+1)} = h \left( \tilde{x}'^{(k)}, \tilde{a}'^{(k)} \right). \quad (9)$$

The main idea behind DR is that, during training,  $L$  selects  $\xi'$  and  $\psi'$ , interacts with the resulting environment  $\tilde{E}'$ , and updates a controller  $\pi$  by observing the consequences (in terms of reward) of selected control inputs  $\tilde{a}'^{(k)}$ , process disturbances  $\zeta'^{(k)}$ , and measurement disturbances  $v'^{(k)}$ . Its final goal is twofold: (a) maximize the finite horizon discounted

reward in a perturbed environment (see Equation (4)) and (b) find a solution  $\pi^*$  which ensures a loss in performance lower than a threshold when applied on different domains of the same distribution. In particular, the final controller  $\pi^*$  sim-to-real transferability is here seen as a form of controller robustness obtained by training  $\pi$  in a *collection of environment models*, chosen by  $L$ , instead of a single one—Figure 6 graphically summarises this process. The resulting controller  $\pi^*$ , learnt by maximising the finite horizon discounted reward  $J_{\pi,T}$  under these conditions, is expected to be robust to perturbations. Therefore, if these perturbations are such that the  $L \leftrightarrow E$  interaction returns a  $\pi^*$  affected by a tolerable RG, the result is a sim-to-real transferable controller.

Table 3 summarizes the articles that tackle the RG using the DR approach. The table shows also the employed learning algorithms and the considered tasks.

We remark that in some of these studies the actual sim-to-real transferability is not evaluated (see Table 2); instead the controller robustness with respect to the perturbations is tested. We discuss each of the paper below.

In Sadeghi and Levine [39] authors train a vision-based navigation policy entirely in simulation, trying to use it on

TABLE 2: List of the surveyed articles, specifying the category (Cat. column). S2S and S2R columns show whether in the article the proposed technique was evaluated by performing sim-to-sim or sim-to-real experiments, respectively.

Ref.	Cat.	S2S	S2R	Simulator
[49]	TL	✓	✓	MuJoCo [64]
[39]	DR	✓	✓	Blender [65]
[43]	TL		✓	MuJoCo
[50]	DR	✓		MuJoCo
[51]	ARL	✓		MuJoCo
[60]	TL	✓	✓	SimSpark [66], Gazebo [67]
[52]	TL	✓		MuJoCo, DART [68]
[53]	DR	✓		Vortex [69], Bullet[70]
[44]	DR		✓	MuJoCo
[45]	DR		✓	PyBullet [71]
[54]	DR	✓	✓	PyBullet
[55]	TL	✓	✓	MuJoCo, PyBullet
[56]	DR		✓	MuJoCo
[40]	TL		✓	Gibson [72]
[57]	DR	✓	✓	MuJoCo
[58]	ARL	✓		TORCS [73]
[46]	DR		✓	PyBullet
[47]	TL		✓	Gazebo
[61]	DR	✓	✓	MuJoCo
[41]	DR	✓	✓	Gazebo

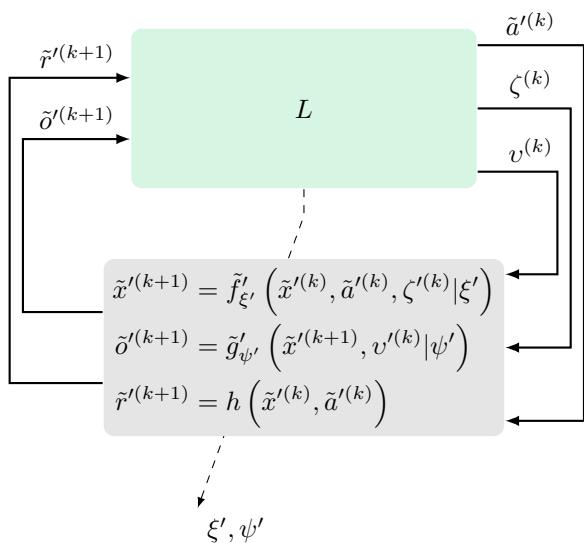


FIGURE 6: Schematic representation of  $L \leftrightarrow \tilde{E}'$  interaction in DR approaches.

a real quadrotor without performing additional real training runs. During training, at each time  $k$ , the state of the system is here represented by an indoor synthetic image  $I^{(k)}$  generated by a renderer. Images are generated in order to reproduce different hallways and a variety of environment parametric settings  $(\xi', \psi')$ . First, a Deep Convolutional Neural Network is learnt in order to predict the collision probability for each

TABLE 3: List of the surveyed articles that apply the DR approach.

Ref.	Algorithm	Tasks
[39]	DQN <sup>h</sup>	Vision-based flight
[50]	TRPO <sup>a</sup>	Inverted pendulum, Half cheetah [82], Hopper [82], Walker [82]
[53]	TRPO <sup>a</sup>	Ball on plate with robotic arm
[44]	HER <sup>b</sup> +RDPG <sup>c</sup>	Pushing task with robotic arm
[45]	DDPGfD <sup>f</sup>	Deformable object manipulation
[54]	PPO <sup>e</sup>	Trotting and galloping of quadruped
[56]	A3C <sup>d</sup>	Marble maze game with robotic arm
[57]	PPO <sup>e</sup>	Rubik's cube with robotic hand
[46]	QT-Opt <sup>g</sup>	Revision-based control task of a robotic arm for grasping
[61]	PPO	Climb and descend stairs with bipedal robot
[41]	A3C	Wheeled mobile platform navigation

<sup>a</sup>Trust Region Policy Optimisation [83], <sup>b</sup>Hindsight Experience Replay [84], <sup>c</sup>Recurrent Deterministic Policy Gradient [85],

<sup>d</sup>Asynchronous Actor-Critic Agents [86], <sup>e</sup>Proximal Policy Optimization [87], <sup>f</sup>Deep Deterministic Policy Gradient from Demonstration [88]. <sup>g</sup>Q-function Targets via Optimization [89] <sup>h</sup>Deep Q Network [90]

$I^{(k)}, a^{(k)}$ . Then, a Deep RL agent is trained for fine-tuning the previous model to provide the action-value function  $Q(I^{(k)}, a^{(k)})$ . Hallway randomisation enacts a wide variety of environments, and shows very good performance during test, both in simulation and on the real world, even with environments never seen during training. However, performance falls when the drone encounters reflective glass doors, thus resulting in a crash.

Mandlekar et al. [50] introduces an algorithm, called Adversarially Robust Policy Learning (ARPL), to teach a controller to correctly behave in presence of increasing adversarial perturbations. The agent uses a curriculum learning approach [91], in which  $\xi'$ ,  $v'^{(k)}$ , and  $\zeta'^{(k)}$  alternately assume the form of isometrically scaled versions of Fast Gradient Sign Method (FGSM) [92]. Here, the controller is parametrized by  $\theta$  ( $\pi_\theta$ ) and updated following the on-policy vanilla Trust Region Policy Optimisation (TRPO) [83]. The key idea is to use a corrupted simulator  $\tilde{E}'$  in training, and then testing the resulting  $\pi^*$  on a different corrupted simulator  $\tilde{E}'' = (X'', A'', O'', Z'', \Upsilon'', f'_{\xi''}, g'_{\psi''}, h)$  environment, obtained with different perturbations  $Z'' \neq Z'$ ,  $\Upsilon'' \neq \Upsilon'$ ,  $\Xi'' \neq \Xi'$  and  $\Psi'' \neq \Psi'$ . These perturbations are such that  $\pi^*$

is misled to provide wrong control inputs  $a''^{(k+1)}$ .

The choice of adversarial perturbations is motivated by the fact that by employing them, the resulting models are likely to generalise well [93].

The ARPL algorithm has been tested in several benchmark examples (Inverted pendulum, Half cheetah, Hopper, Walker) and seems to deliver promising results, exhibiting significant robustness. However, examples of sim-to-real controller transferability have not yet been provided.

The Simulation-based Policy Optimisation with Transferability Assessment (SPOTA) algorithm, designed in [53], uses randomised physics parameters, drawn from a probability distribution  $\xi' \sim \rho_\kappa(\xi')$  (parametrized by  $\kappa$ ), to perform a robust optimisation of the controller. In SPOTA, the controller is trained on model ensembles, according to the following 4 phases: (i) learn a candidate solution  $\pi_\theta^C$  using a TRPO updating rule; (ii) learn  $n_R$  reference solutions  $\pi_\theta^{R_j}, j = 1, \dots, n_R$  on  $n_R$  different  $\tilde{E}'$ , each obtained for different  $\xi'$  and  $\psi'$  settings; (iii) compare the performance of candidate  $C$  with that of each reference  $R_j$  in the same condition of  $R_j$ ; and, finally, (iv) decide whether or not stop the learning. The last step is carried out by introducing a Simulation Optimisation Bias (SOB) concept: an error caused by an optimistic bias of the optimisation procedure, whose existence has been proven by [94]. The authors have assumed that it can be treated as the error between the finite horizon discounted reward  $J_{\pi_\theta^{*R}, T}$  obtained by considering the reference solutions and the finite horizon discounted reward  $J_{\pi_\theta^C, T}$  of candidate solution. Taking into account that an RL approach in a stochastic setting allows to find only estimates of  $J_{\pi_\theta^{*R}, T}$  and  $J_{\pi_\theta^C, T}$ , the authors have derived an upper bound for the tolerated SOB of the candidate solution, called Upper Confidence bound on Simulation Optimisation Bias (UCSOB). In order to ensure a desired performance  $\beta$ , the final candidate solution UCSOB must be lower than  $\beta$ . In this framework, the authors have tested the algorithm by developing a controller for a ball on plane task, governed by a robotic arm, in the same simulator (obtaining satisfactory results) and varying the physics engine (with worse outcomes). However, although in this case a sim-to-sim transferability test has been carried out, a sim-to-real controller test has not been done.

Peng et al. [44] shows the effectiveness of memory-based policies (i.e., policies learnt by using past memory for future learning [95, 96, 97]) to deal with the RG, introducing DR to generalise environment dynamics. Hindsight Experience Replay [84] has been used for the purpose: a technique able to generalise over different goals using past experience as a baseline. In this case, the parameters  $\xi'$ , the measurement noise  $v'^{(k)}$ , and the time step  $\Delta t$  are sampled according to a distribution, which is a design parameter. In particular,  $\xi'$  is kept locked for an entire episode, while the remainder are varied at each time step. The proposed solution, learnt using a RDPG algorithm (off-policy), has been tested on a robotic pushing task and, when transferred to reality, shows performances comparable to those obtained in simulation,

despite poor calibration.

An improved version of DDPG [88] is adopted in [45] to solve deformable object manipulation tasks in simulation. The resulting controller transferability on the real robot is therefore tested. In particular, a robotic arm is involved in three different towel folding tasks, in which RGB images are included in the observation  $o$ . The DR is here implemented by sampling some environment values from either normal or uniform distributions around noisy ground truth estimates. Experimental results suggest that randomisation of extrinsic camera parameters (i.e., position and orientation) is particularly useful for sim-to-real transfer, since the controller has an evident sensitivity to changes of its position. Besides, they show that heavy randomisation can lead to unsuccessful transfers.

Controller sim-to-real transferability has also been tested on locomotion tasks of a Minitaur quadruped of Ghost Robotics [54]. Here authors have used Proximal Policy Optimisation (PPO) to learn  $\pi^*$  and have observed the impact of two different solutions to reduce the RG: (a) improving simulated model via system identification; (b) using randomised  $\zeta'^{(k)}, v'^{(k)}$ , and  $\xi'$  to learn robust controllers as the observation space changes. Obtained results suggest that simulators improvement is an essential requirement since, as the model becomes less adequate, not even a robust controller is able to avoid a large RG. The authors of the cited paper also pointed out that considering a large observation space does not always bring benefits. On the contrary, their evaluations have showed that controllers learnt in simulation with large observation space lead to bad results when transferred to real robot.

Van Baar et al. [56] shows the benefits of using DR and Asynchronous Actor Critic Agent (A3C) algorithm [86] (on-policy) for learning the controller, with respect to not using DR. The parameters  $\xi'$  are here randomly sampled according to a uniform distribution. Both controllers are then applied on real-world robot and the fine-tuning time required to convergence is compared. The analysed task is a Marble maze game driven by a robotic arm and the results show that there is a trade-off between controller robustness and fine-tuning steps. Controller learnt through DR requires fewer fine-tuning steps than the remaining one, further proof of an existing trade-off between efficiency and RG.

In a quite recent work [57], Automatic Domain Randomisation (ADR) is proposed in order to transfer a policy learnt in simulation on the real system, framing it in a manipulation task of Rubik's cube with a robotic hand. Here, the RL agent does not solve Rubik's cube but "only" learns, using a PPO algorithm (on-policy), how to move correctly the robotic hand in order to perform control inputs suggested by another non-AI based algorithm. What changes from standard DR idea are  $\rho_\kappa(\xi')$  and  $\rho_\kappa(v^{(k)})$  distributions (parametrized by  $\kappa$ ) that allow to randomly select  $v^{(k)}$  and  $\xi'$ . Indeed, while in other DR approaches these distributions are parametrized with fixed  $\kappa$ , in ADR  $\kappa$  changes during learning procedure. In particular, these additional environments, obtained

with different  $\kappa$ , are added to the considered collection of environment models only when a lower performance limit is reached (i.e., a fixed number of successful episodes are performed). The developed controller has been firstly tested in environments in which distributions were manually tuned, achieving good results. In addition, a sim-to-real transfer is performed, with worse results.

A Randomised-to-Canonical Adaptation Network (RCAN) is conversely introduced in [46]. The main idea is to map the observations collected on the simulated domain as well as those collected on the real domain into a common further domain called the *canonical domain*. The approach has been applied to a vision-based robot grasping task, and the canonical domain consists of extremely simplified images whose purpose is capturing just the relevant information for the task. The map is learnt by using an image-conditioned Generative Adversarial Network (cGAN) [98], able to map an image of a domain  $D$  into an adapted image of the canonical domain  $D_c$ ; i.e.,  $G : D \rightarrow D_c$ . The resulting image of  $D_c$  is then sent to the controller which is learnt by using Q-function Targets via Optimisation (QT-Opt) [89]. During training, cGAN receives randomised simulated images, sampled from the trajectories, and learns to convert them in canonical images. The resulting observations are then used by the QT-Opt to produce the policy. In the test procedure, the real-world images are mapped into canonical images and sent to the controller. The proposed approach returns excellent results, however, an effective transfer is not always achieved. In particular, when the cGAN during training is fed with images sampled only from non-successful trajectories, the final controller results in an unsatisfactory transfer.

Siekmann et al. [61] proposes a simple terrain randomization to learn robust proprioceptive controllers for bipedal robots involved in the task of climbing and descending stairs. They model the policy with a Long short-term memory (LSTM) network, for its capability of processing temporal sequences. Indeed, unlike feed-forward neural networks, LSTMs are equipped with a feedback mechanism that allows them to process sequences of input data, without treating each sample of the sequence independently. They retain useful information about earlier data points in the sequence, aiding in the processing of new data points. The authors compare the performance of three different controllers  $\pi$ : (i) A learned LSTM controller with different terrain parameters  $\psi$ , (ii) a feed-forward NN controller learned with different  $\psi'$  terrain parameters, and (iii) a LSTM controller learned on a single simulated environment. The experimental results show that the first  $\pi$  is the one with the highest overall probability of success in the task. Thus, the combination of LSTM and DR seems to be an effective solution to the problem for the tested task.

Finally, Hu et al. [41] face the reality gap of a controller involved in a wheeled robot navigation task. The proposed solution tries to render the controller robust to possible parametric errors in the model, but also to possible disturbances that corrupt its dynamics. To this end (i) a terrain

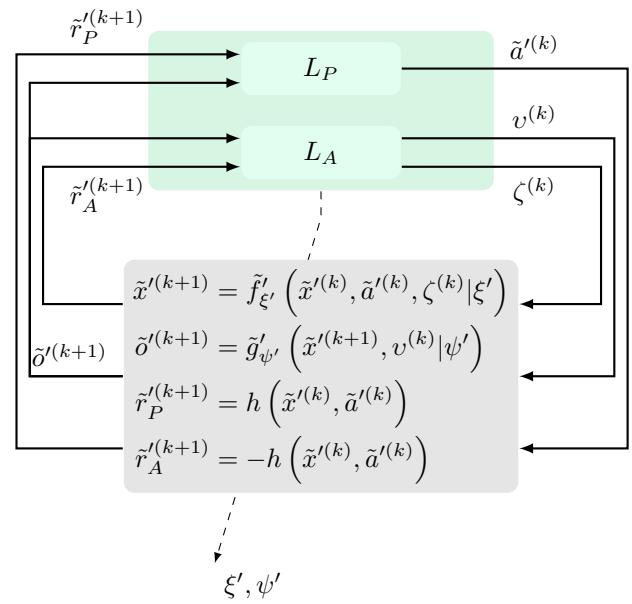


FIGURE 7: Schematic representation of  $L \leftrightarrow \tilde{E}'$  interaction in ARL approaches.

randomisation is performed by varying its viscosity and inclination  $\psi'$ , (ii) disturbances on the travel distance and on the yaw rotation  $v_1^{(k)}$  are randomly imposed to the dynamics, (iii) latency disturbances are imposed at each time-step, hypothesising that a latency between perception and movement performance occurs  $v_2^{(k)}$ , and (iv) the pose-estimation error  $v_3^{(k)}$  is also taken into account. The resulting  $\pi$ , learned entirely in simulation, results in an effective application on the real-world environment. Moreover, a comparison with some state of the art solutions for robot navigation highlights the better performance of the proposed approach in terms of success rate as well as cumulative travel distance, and time required for task execution.

## B. ADVERSARIAL RL

In the adversarial RL (ARL), the agent  $L$  is composed of two sub-agents: the *protagonist*  $L_P$  and the *antagonist*  $L_A$ . The underlying idea resembles the one behind domain randomisation: enforce robustness (and, hence, improve controller transferability) by training the controller in a collection of environment models instead of a single one. In the case of ARL, however, the diversity is obtained by training a secondary controller (the adversarial) to generate more difficult models to handle (those that minimise the cumulative reward). Figure 7 graphically summarises the process of ARL.

Given a corrupted simulator  $\tilde{E}'$  of  $E$ , defined by a tuple  $(\tilde{X}', \tilde{A}', \tilde{O}', \tilde{Z}', \Upsilon', \tilde{f}_{\xi'}, \tilde{g}_{\psi'}, h)$  and evolving according to Equations (7) to (9),  $L_P$  and  $L_A$  interact with  $\tilde{E}'$  seeking to maximise their respective discounted cumulative reward [99].

We denote with  $\tilde{r}_P^{(k+1)}$  the reward of  $L_P$ . A common

TABLE 4: List of the surveyed articles that apply the ARL approach.

Ref.	Algorithm	Tasks
[51]	TRPO <sup>a</sup>	Inverted pendulum, cheetah, Swimmer [82], Hopper, Walker 2D, Ant [82]
[58]	Ensemble DQN <sup>b</sup>	Autonomous driving

<sup>a</sup>Trust Region Policy Optimization [83], <sup>b</sup>Ensemble Deep-Q Network [104].

choice [100] is to provide  $L_A$  with a reward  $\tilde{r}_A^{(k+1)} = -\tilde{r}_P^{(k+1)}$ . As a result two controllers are learnt:

- $\pi_P$ , whose target is to maximise the cumulative reward over time, resulting in the final controller which will be tested in a  $\tilde{E}'$ -to- $E$  transfer;
- $\pi_A : \tilde{O}' \rightarrow Z' \times \Upsilon' \times \Xi' \times \Psi'$  that searches for those environment perturbations or parameters variations that *minimise* the same cumulative reward over time.

The outputs of  $L_A$  and  $\pi_A$  are perturbations and parameters variations of  $\tilde{E}'$ . In a noise-corrupted simulated environment, the observed reward and hence the discounted reward will depend on the disturbances as well as on the policy. Since, in ARL, disturbances are generated by the adversarial agent, the finite horizon discounted reward will depend on both policies. To catch this dependency we can write  $J_{\pi_P, \pi_A}(\tilde{x}'^{(0)})$ . The resulting  $L$  goal can be compactly stated as:

$$\max_{\pi_P} \min_{\pi_A} J_{\pi_P, \pi_A}(\tilde{x}'^{(0)}), \quad (10)$$

falling in a *worst-case* approach that in control theory is known as *min-max optimal control* (also referred to  $H_\infty$ -control) [101, 102, 103].

Consequently, in ARL,  $\tilde{E}'$  interacts simultaneously with  $L_A$  and  $L_P$ , thus  $L \leftrightarrow \tilde{E}'$  results in  $L_A \leftrightarrow \tilde{E}' \leftrightarrow L_P$  and the global  $L$  task is to find:

$$\pi_L^* = \arg \max_{\pi_P} \arg \min_{\pi_A} J_{\pi_P, \pi_A}. \quad (11)$$

Table 4 summarises the articles that use ARL for learning robust controllers, along with the respective learning algorithms and the considered tasks.

The ARL approach was firstly introduced by [51], with the Robust Adversarial Reinforcement Learning (RARL) algorithm. There,  $\pi_P$  (the protagonist's policy) is trained to work in presence of an adversary ( $\pi_A$ ), able to inject destabilising disturbances to environment (in particular only  $\xi'$  disturbances). The proposed solution can be summarised in two main steps that are repeated  $n_{\text{iter}}$  times: (i) learn the protagonist policy  $\pi_P$  while keeping the adversary one fixed; (ii) learn the adversary policy  $\pi_A$  while keeping  $\pi_P$  fixed. The experiments have been done on several OpenAI Gym environments [82]. In a first experiment, the authors compare mean and variance of the cumulative reward over 50 RARL

policies, obtained using different seeds and initialisation, with TRPO ones. For all tasks RARL behaves better than TRPO in terms of mean and variance. In a second experiment, [51] shows that RARL behaves better than TRPO under adversarial attacks while keeping hold the protagonist. Finally, in a third experiment, the authors introduce different  $\zeta^{(k)}$  in the test phase, and again obtain better results with respect to TRPO.

Pan et al. [58] builds on the RARL idea by introducing the Risk-Averse Robust Adversarial RL (RARARL) concept: a RARL algorithm in which the protagonist is trained to be risk-averse and the adversarial, in contrast, risk-seeking. The authors of the cited study state that “a robust policy should not only maximise long-term expected reward, but should also select actions with low variance of that expected reward”. For that purpose, they train  $\kappa$  different  $Q$ -value networks that return  $\kappa$  action-value outputs. The risk of an action is estimated by the empirical variance of these  $\kappa$   $Q$ -values ( $\text{Var}_\kappa(Q)$ ). At the beginning of each training episode one of the  $\kappa$  networks is randomly chosen and employed for control input selection during the entire episode. The protagonist and the antagonist take actions sequentially: the protagonist action-value function  $Q_{\pi_P}$  is augmented by a risk-averse term ( $\text{Var}_\kappa(Q_{\pi_P})$ ), which encourages the choice of lower variance control inputs; the adversary  $Q_{\pi_A}$ , instead, is reduced by a risk-seeking term ( $\text{Var}_\kappa(Q_{\pi_A})$ ) in order to guide it towards higher-variance outputs. The algorithm has been tested in a simulated self-driving task and obtained experimental results highlighting the better robustness of RARARL controller with respect to one subjected to random perturbation in training. During the test, control inputs are selected according to the mean value of the  $\kappa$  networks.

### C. TRANSFER LEARNING

The previously discussed approaches are aimed at controllers that, once learnt in simulation, can be directly transferred on real robots without any (or, at most, with very few) additional training steps. Basically, an agent searches for a sim-to-real transferable controller.

A different perspective is the one adopted in transfer learning (TL) approach. Indeed, its aim is not to find a solution to the RG, but rather to avoid its occurrence by means of two subsequent or simultaneous training phases (first in simulation and second in reality), penalising the resulting  $L$  efficiency.

Let us ignore for the moment the RG and suppose to be in a classic RL training scenario, described by  $L \leftrightarrow E$ . The basic idea of TL is that generalisation is possible not only *within* task but also *between* tasks [105]. Therefore, since a task can be entirely defined by an environment, considering a second different but *compatible* environment<sup>2</sup>, the controller learnt

<sup>2</sup>Two environments  $E = (X, A, O, f, g, h)$  and  $E' = (X', A', O', f', g', h')$  are compatible if and only if  $A = A'$  and  $O = O'$ . Thus, a policy for a system  $E$  is also a policy for every system  $E'$  that is compatible with  $E$ . However, a policy which is optimal for  $E$  is in general not optimal for  $E'$ .

for the first is expected to be a helpful tool to speed-up the second learning process, whether or not it involves the same agent. Thereby, in RL, in which the controller is the result of a trial and error process, TL could be employed to speed-up the learning, thus avoiding training from scratch.

Back to the RG, the idea of “recycling” policies between tasks could be useful in speeding the real robot learning procedure or, possibly, while performing a fine tuning of simulated and real agents ( $L'$  and  $L$  respectively).

In the first scenario, what has been learnt in simulation, by using  $L'$ , is reused (in its entirety or in part) in subsequent phases of real-world training, performed by using  $L$ . The expected result is a faster real-world training that bridges the discrepancy between simulator and real robot at its root, i.e., while real agent  $L$  is learning its  $\pi$ . Therefore, the transfer occurs once, and only in one direction (from simulation to reality).

In the second case, by providing some real information to simulator and taking example from it, the simulator could more realistically adapt itself to reality, thus reducing performance misalignment and, thereby, the RG. In this situation, the transfer is repeated, and in both directions (sim-to-real and real-to-sim).

Overall, we will refer to the exchanged information as  $u_{\text{exch}}^{(k)} = [u'^{(k)} \ u^{(k)}]^T$ , in which  $u'^{(k)}$  is the sim-to-real transferred information, while  $u^{(k)}$  the real-to-sim one (Figure 8).

We categorise the different TL approaches based on:

- the kind of information passed through  $u_{\text{exch}}^{(k)}$ , (e.g., weights of an image processing net, state-control input pairs, policy parameters or even the policy itself);
- the transfer timing of  $u_{\text{exch}}^{(k)}$ , either continuous (C-TL) or one shot (1-TL), the latter generally occurring at the end of the training in simulation;
- the direction of  $u_{\text{exch}}^{(k)}$  transfer, either bidirectional (2D-TL), i.e.,  $u_{\text{exch}}^{(k)} = [u'^{(k)} \ u^{(k)}]^T$ , or unidirectional (1D-TL), i.e.,  $u_{\text{exch}}^{(k)} = [0 \ u^{(k)}]^T$  or  $u_{\text{exch}}^{(k)} = [u'^{(k)} \ 0]^T$ ;
- whether there is an agent  $L$  (e.g., RL, Inverse Dynamics Neural Network (IDNN), etc.) or a controller  $\pi$  (e.g., Model Predictive control (MPC), etc.) in the real setting  $E$ .

Table 5 summarises the articles that use TL and characterises them in terms of these four factors. The table also shows the tasks and algorithms.

Christiano et al. [49] adopted TL to perform sim-to-real control input adaptation. They assume that if simulator does not exactly replicate the real robot behaviour, applying the same control input in both scenarios does not necessarily lead to same observation. However, they assume that the observation  $o'^{(k)}$  obtained in simulation, is the one that should be achieved also on the real environment. Therefore, given the simulated observation  $o'^{(k)}$ , they employ past history to discover what real control input  $a^{(k)}$  can lead to  $o^{(k)} \simeq o'^{(k)}$ . For this purpose, a neural network is trained in order to predict the control input  $a^{(k)}$  that leads to a specific  $o^{(k)}$ . In particular, it is assumed that a simulated-based policy  $\pi'$ ,

a forward dynamic robot model  $F$ , and a sequence  $\tau_i = (o^{(0)}, a^{(0)}, \dots, o^{(i-1)}, a^{(i-1)}, o^{(i)})$  of  $i$  real observations and  $i - 1$  real control inputs are known. Here, TRPO is used to learn  $\pi'$ , but any  $L$  agent (not necessarily RL) could be used for the purpose. While training the policy  $\pi$ , the policy  $\pi'$  returns a control input  $a'^{(i)}$  based on the provided history  $\tau_i$ . However, rather than being applied to the real robot, it is sent to  $F$  and the resulting observation  $o'^{(i)}$  is provided, together with  $\tau_i$ , as inputs to  $L$ . Finally, the learnt policy  $\pi$  provides the control input  $a^{(i)}$  that results in  $o^{(k)}$  similar to  $o'^{(k)}$ . Therefore, here the communication is continuous and in both directions  $u_{\text{exch}}^{(i)} = [u'^{(i)} \ u^{(i)}]^T = [o'^{(i)} \ \tau_i]^T$ . It is worth remarking that what is actually learnt in this case is an inverse dynamic model (implemented as a neural network) that must be employed in conjunction with the controller learnt in simulation. The authors of the cited study evaluate their approach first in a sim-to-sim scenario with several OpenAI Gym environments and then in a sim-to-real scenario based on a Fetch robot. Results highlight the effectiveness of such method and the relatively low number of samples required for convergence. However, the authors assume that, when the consequences of a control input applied in simulation differ from those applied on the real robot, real observations should match simulated ones, and that is not always true.

In [60], the Grounded Action Transformation (GAT) algorithm is proposed to learn a humanoid bipedal locomotion policy. Inspired by the Grounded Simulation Learning (GSL) idea, introduced in [109], they try to reproduce it in a RL framework, additionally improving some aspects. GSL is based on two main principles: *grounding* and *guide*. The former refers to making the simulator  $E'$  closer to the real robot  $E$ , by properly modifying some parameters of  $E'$  on the basis of data collected from  $E$ . The latter consists in having an expert able to guide the optimisation algorithm in finding the proper parameters of  $E'$  to be tuned. In practice, given an evaluation function  $J_{\text{eval}}$ , such as a penalty function (for example the opposite of the reward), a policy  $\pi$  is applied to the real robot in order to collect end-effector trajectories  $\mathcal{D}$ . By performing  $\mathcal{D}$  both in simulation and in reality, and collecting the resulting real end-effector trajectories, an optimisation problem is solved in order to find those  $E'$  parameters able to minimise the Kullback-Leibler divergence between the probability of observing the same trajectories in the two cases. The resulting  $E'$  is therefore used in order to find a set of candidate policies  $\Pi_C$  trying to minimise  $J_{\text{eval}}$ . The optimal policy is the  $\pi \in \Pi_c$  such that  $J_{\text{eval}}$  is minimised once performed on the real robot  $E$ . However, the above procedure is aimed at finding the correct values of the  $E'$  parameters. Conversely, GAT introduces an action transformation function  $a'^{(k)} = m(a^{(k)})$ , learnt in a supervised fashion, able to map each action  $a^{(k)} \in A$  into an action  $a'^{(k)} \in A'$ . In particular, a forward robot dynamics model is trained to compute the  $x^{(k+1)}$  resulting from  $a^{(k)}$ . The inverse robot dynamics, instead, is trained to find the simulated action  $a'^{(k)}$  able to lead the simulator in  $x'^{(k+1)} = x^{(k+1)}$ . The resulting

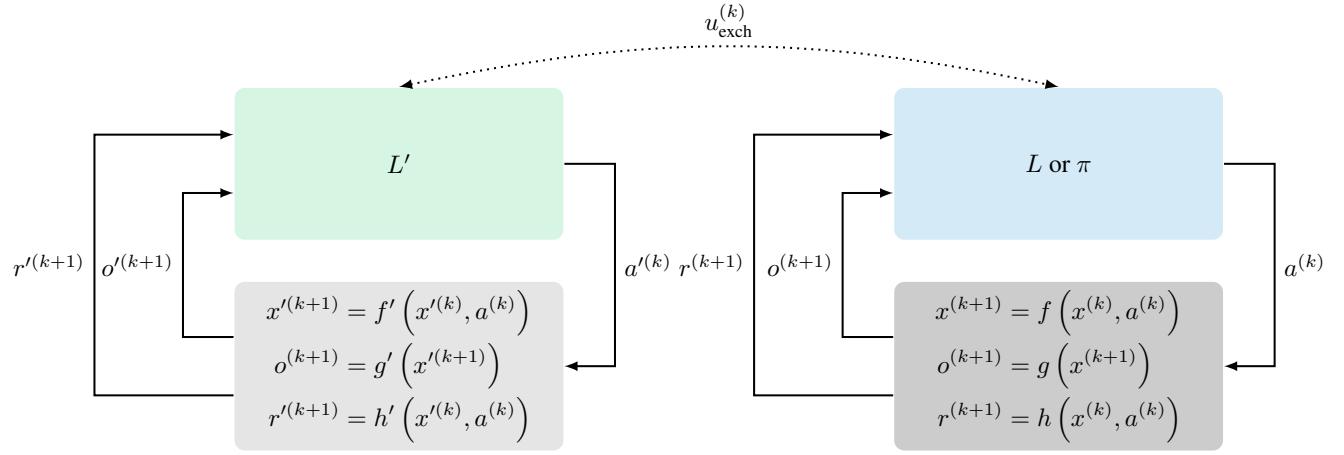


FIGURE 8: Schematic representation of TL approaches.

TABLE 5: List of the surveyed articles that apply the TL approach.

Ref.	Algorithm	Transf. info. $u_{\text{exch}}^{(k)}$	Timing	Dir.	$L, \pi$ in $E$	Tasks
[49]	TRPO <sup>a</sup>	$a'(k)$	C-TL	2D-TL	IDNN <sup>e</sup>	Reacher, Half cheetah, Hopper, Humanoid, Fetch robot, Trajectory control
[43]	A3C <sup>c</sup>	NN activation functions	1-TL	1D-TL	RL	Jaco robot manipulation
[60]	CMA-ES <sup>b</sup>	$a(k), a'(k)$	C-TL	2D-TL	RL	Humanoid bipedal locomotion
[52]	TRPO <sup>a</sup>	$o'^{(k)}, a'^{(k)}, o^{(k)}, a^{(k)}$	C-TL	2D-TL	RL	Reacher, Hopper 2D, Reacher 2D
[55]	PPO <sup>g</sup>	$\pi'^*$	1-TL	1D-TL	RL	Pusher, Striker, Ergo Reacher, ErgoShield
[40]	DQL <sup>d</sup>	NN weights	1-TL	1D-TL	MPC <sup>f</sup>	Quadrotor collision avoiding
[47]	DQL <sup>d</sup>	$O', Q(x'^{(k)}, a'^{(k)})$	1-TL	1D-TL	RL	Robot object manipulation

<sup>a</sup>Trust Region Policy Optimization [83], <sup>b</sup>Covariance Matrix Adaptation-Evolutionary Strategy [106], <sup>c</sup>Asynchronous Actor-Critic Agents [86], <sup>d</sup>Double Q-Learning [107], <sup>e</sup>Inverse Dynamics Neural Network, <sup>f</sup> Model Predictive Control [108]

procedure leads to  $u_{\text{exch}}^{(k)} = [u'^{(k)} \ u^{(k)}]^T = [a'^{(k)}, a^{(k)}]^T$ . Both sim-to-sim and sim-to-real experiments provide good results; however, as authors point out, the drawback of using a supervised method for  $m(\cdot)$  learning is that policies are no longer effective when there are changes between training and testing distributions. Moreover, neglecting the contact dynamics can lead to simulation bias.

Wulfmeier et al. [52] proposed Mutual Alignment Transfer Learning (MATL), a method that relies on a Generative Adversarial Network (GAN) [110]. The main idea is similar to [49]: enforcing a similarity between observations  $o'^{(k)}$  and  $o^{(k)}$ . Indeed, although a control sequence achieving the goal in simulation may not produce the same effects in the real environment, the corresponding sequence of simulated observations, if reproduced on the real system, can lead to task accomplishment. For the purpose, here, simulator  $E'$  and real robot  $E$  work in parallel as generators and interact with two different agents ( $L'$  and  $L$  respectively). A discriminator  $D$ , instead, is employed (and trained along with  $L'$  and  $L$ ) to classify the environment from which a sequence of

observations  $\tau_\kappa$ , provided as input, has been collected.  $L$  and  $L'$  are trained not only to maximise their respective environment reward, but also to mislead the discriminator, based on the assumption that the more the discriminator is misled, the more “aligned” the observations are. Therefore, each time a sequence of observations  $\tau_\kappa$  is collected,  $D$  will receive it as input and will output the probability  $D(\tau_\kappa)$  that it was generated by  $E'$ . A term  $\log(D(\tau_\kappa))$  is respectively added and subtracted to  $E$  and  $E'$  rewards thus encouraging misleading actions (here TRPO is the employed algorithm). The proposed solution results in an experience exchange between  $L$  and  $L'$  and a consequent alignment of the collected observations. To exploit the simulator,  $L'$  is updated  $M$  times more frequently than  $L$ , thus accelerating learning. Thereby, in this case  $u_{\text{exch}}^{(k)}$  is the  $D$  output and it is equal to  $[-\log(D(\tau_\kappa)) \ \log(D(\tau_\kappa))]^T$ . Wulfmeier et al. [52] evaluate their approach on various RL tasks: rllab [111], OpenAI Gym, and DartEnv [82]. Results show that MATL is able to work with significantly different environments of same simulator in which only parameters variation is performed. Less

encouraging results are reported when employing different simulators.

A different solution was proposed in [43], where progressive nets [74] are employed for sim-to-real information transfer. Here, by exploiting the capability of those nets to learn a tasks sequence through lateral connections, simulated knowledge can be used to avoid training from scratch on real robot. A progressive net is composed of  $l$  “columns” in which the  $i$ -th column represents an independent network of  $\kappa$  hidden activations. Each  $j$ -th activation  $\text{act}_{j,i}$  of the  $i$ -th column is a function of the same column  $j - 1$ -th activation ( $\text{act}_{j-1,i}$ ) and of the  $j - 1$ -th activation of all the previous  $m < i$  columns ( $\text{act}_{j,1}, \dots, \text{act}_{i-1}$ ). Rusu et al. [43] propose to use this tool in order to learn a simulated controller  $\pi'$ , via A3C algorithm, in the first column and subsequently transfer its knowledge on real robot by means of lateral connections headed towards the second column, which represents the real agent  $L$ . Then, the second net training ( $L$ ) begins. Therefore, letting  $s$  be a simulated column,  $u_{\text{exch}}^{(k)} = [u'^{(i)} \ 0]^T = [(\text{act}_{1,s}, \dots, \text{act}_{k,s}) \ 0]^T$ . The authors evaluate their approach on a robot manipulation task on Jaco arm [112] in which a visual target must be reached. The performances are compared with those obtained using a fine-tuning approach showing the superiority of progressive nets.

In [55], a Neural Augmented Simulation (NAS) approach is used to reduce the RG. A Long Short Term Memory is trained on the differences between simulated and real robot, and used to adapt the simulator on the basis of real world data. The policy  $\pi$  resulting from a Proximal Policy Optimisation algorithm is learnt on the simulated environment, whose next state at each step is adjusted by using the correction term  $\Delta$  provided by the LSTM. The resulting policy, therefore, associates to each estimated value of the real robot state  $x^{(i)} = x'^{(i)} + \Delta$  an action  $a'^{(i)} = a^{(i)}$ . The transfer is in this case only from the simulator to the real robot and  $u_{\text{exch}}^{(k)} = [u'^{(i)} \ 0]^T = [\pi \ 0]^T$ . The NAS has been tested in a sim-to-sim and a sim-to-real transfer. For the former, authors create an artificial RG by varying some parameters of two different simulated robotics environments of Open AI Gym, one of which was considered as the real environment. For the sim-to-real transfer, two Poppy Ergo Jr robots [113] have been used in a ErgoShield task, in which an *attacker* (one of the two involved robots) is controlled to touch as often as possible the shield attached to the end-effector of a *defender* (the other robot). The defender is able to move the shield in random poses. Experimental results show good performance both in sim-to-sim and sim-to-real transfer. Moreover, since a policy-specific fine-tuning is not required, the method can be appropriate for multi-task robotic applications.

In [40], conversely, they propose to learn a RL controller in simulation (using a deep Q-Learning approach) in which the first stages represent a visual perception module, parametrized by vector  $\theta_{\text{VP}}$ . Therefore, by keeping the weights fixed, they use this module to work with real-data and predict rewards for  $h$  planned control inputs by learning a DNN. The predictor is trained, by means of a real-world

data-set, in order to minimise reward prediction error. In the second phase, the real application, the predictor is used by an MPC controller. Hence, at each step, the MPC controller computes a sequence of  $h$  control inputs which maximises the expected discounted predicted reward within an horizon  $h$ . In this context,  $u_{\text{exch}}^{(k)} = [u'^{(k)} \ 0]^T = [\theta_{\text{VP}} \ 0]^T$ . As prescribed by the MPC approach, only the first control input is applied; hence, the process is repeated. Kang et al. [40] consider a nano aerial vehicle collision avoidance task to assess the proposed solution. Moreover, the authors compare it with other approaches: simulation only, simulation with fine-tuning, simulation with fine-tuning and perception fixed, real world only, supervised and unsupervised. Their solution outperforms all others tested, and shows the best result in terms of pre-collision time.

Yuan et al. [47] performs an action-value function adaptation in a supervised fashion. Here, a Baxter robot is asked to solve a nonprehensile rearrangement task, i.e., the problem of pushing an object into a predefined goal pose. The proposed procedure consists of three sequential steps: (1) learning, in simulation, an optimal action-value function  $Q_{\text{sim}}$  able to select the best action  $a'$  to perform when an image of the scene is provided as observation  $o'$ , (2) collecting a data-set of real and simulated observation pairs  $(o, o')$ , and (3) using it, along with the pre-trained  $Q_{\text{sim}}$ , in order to create a  $Q_{\text{real}}$  useful to adapt the agent for a real world application. In the former, a deep-Q network is used to approximate  $Q_{\text{sim}}$ . In the second step, starting from real scenes  $o$ , the obstacle and the portable object positions are used to recreate the same scenes in simulation  $o'$ . In particular, randomly setting  $o^0$ , the RL agent (learnt in step (1)) is applied to the real robot in order to collect a set of  $O_{\text{real}} = o^0, o^1, o^2, \dots$  from which the respective simulated counterpart set  $O_{\text{sim}} = o'^0, o'^1, o'^2, \dots$  is created. The resulting data-set, composed of  $(O_{\text{real}}, O_{\text{sim}}, Q_{\text{sim}})$  is used in order to learn a  $Q_{\text{real}}$  able to minimise a loss function defined as  $r + \gamma Q_{\text{sim}} - Q_{\text{real}}$ . Three different strategies have been deployed: (a) train  $Q_{\text{real}}$  keeping the  $Q_{\text{sim}}$  structure but retraining the network in its entirety; (b) use  $Q_{\text{sim}}$  as baseline and adapt only the parameters of the convolutional layers; (c) add two new fully connected layers to increase the flexibility of the network and learn their parameters and those of the convolutional layers. Therefore in this case  $u_{\text{exch}}^{(k)} = [u'^{(k)} \ 0]^T = [(O_{\text{sim}}, Q_{\text{sim}}), 0]^T$ . Authors compare their results with the one obtained by using the same domain randomisation idea of [75]. Experimental results show that their approaches surpass the one of [75] in terms of performances and, in particular, (c) turns out to be the best. However, collecting data both in simulation and in reality in order to build the data-set used to learn the  $Q_{\text{real}}$  may be costly and time-consuming.

## D. DISCUSSION AND PROMISING IDEAS

Despite several attempts found in the literature to make sim-to-real transferable controllers, many of which associate the idea of robustness with that of sim-to-real transferability (DR and ARL), the lack of uniformity of the considered tasks does

not allow to determine which solution is the more appropriate in terms of the RG. Besides, a considerable fraction of the proposed approaches were not experimentally evaluated in an actual sim-to-real scenario. A controller robust to certain model disturbances or parametric variations, is not necessarily sim-to-real transferable. In fact, if these variations and disturbances do not correctly represent the simulator inaccuracies with respect to reality, it might result not sim-to-real transferable. This suggests that an interaction with the real system during training is still needed: to this respect, TL approaches appear promising. On the other hand, the TL approaches here surveyed often lacked an assessment of the robustness. Moreover, since TL requires two successive (or simultaneous) training phases, it may exhibit low efficiency.

Some mixed approaches exist, that borrow ideas from DR, ARL, and TL. A first attempt in merging DR and TL is proposed in [114]. Here the authors propose to learn a policy in a randomised simulation and to adapt the distribution of simulation parameters on the basis of a real-world performance.

A promising research direction to tackle the RG problem could be the meta-learning strategy application in RL, in order to quickly adapt experience gained in simulation on the real system [115, 116]. In Meta-RL, given a distribution over tasks, the agent learns an adaptive policy that maximises the expected reward for a new task from the distribution. A recent work [117] has shown the great ability of this approach to generalise in environments totally different from those used during the training.

Another promising solution to avoid a direct RL training on the real robot seems to be the Probabilistic Inference for Learning Control (PILCO) proposed by [118]. Here, a probabilistic model of the system dynamics is learnt incorporating uncertainty by using only some trial on the real system and a policy is learnt through it. This solution allows to avoid the RG in the first place, by training a simulator with few real interactions and using it for the trial and error procedure. However, although the potential usefulness of PILCO and Meta-RL to cope with the RG, the former underestimates state uncertainty at future time steps [119], thus possibly leading to a decrease in performance; the latter, on the other hand, is computational demanding and needs a high number of real-world evaluations [120].

## V. CONCLUSIONS AND OPEN CHALLENGES

Training a RL robotic controller in real-time in its actual environment is a costly process. While simulators can alleviate the problem, the approximations often present in the employed models play a crucial role in determining the effectiveness of the learnt controllers. The more accurate the model of the robot (and the surrounding), the more effective the controller but, at the same time, the greater the computational cost of learning it. When performances of the controller in the simulator and in real robot are different, a RG exists and the controller is not sim-to-real transferable.

In the present article, we provided a formal framework for the RG and reviewed the most relevant existing methods aiming to achieve sim-to-real transferable controllers in robotics RL applications. We surveyed the literature concerning RL and RG and categorised the approaches as: domain randomisation, adversarial reinforcement learning, and transfer learning. Moreover, we described them in detail according to the proposed framework and in terms of the employed algorithms, involved tasks, and evaluation methods.

We conclude commenting on some significant open challenges. Each one of the examined approaches appears tailored to a specific task, and its applicability to other, potentially different, tasks is not clear. A general task-independent approach able to guarantee an effective sim-to-real transferability of the controller is still missing. With this in mind, we believe that a significant open problem is that of providing a proper index able to reveal and quantify the RG. Indeed, being able to characterise and quantify the RG would (i) enable systematic comparison among different techniques, hence favouring the advancement of research, and (ii) allow to use the measure of RG directly as an optimisation objective, hence putting the transferability as a direct goal in the learning of RG-aware controllers. Another significant open problem is the sample efficiency of the developed approaches. As highlighted in [121], RL is very data-intensive. The computation effort required for learning a RL agent involved in a complex task can be huge, thus limiting the practical applicability of such methods. The surveyed approaches lead to learning paradigms that unquestionably aggravate this issue, especially in those cases in which an interaction with different environment domains is involved.

## REFERENCES

- [1] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [2] J. Kober, J. A. Bagnell, and J. Peters, “Reinforcement learning in robotics: A survey,” *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1238–1274, 2013.
- [3] A. A. Shahid, L. Roveda, D. Piga, and F. Braghin, “Learning continuous control actions for robotic grasping with reinforcement learning,” in *2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, 2020, pp. 4066–4072.
- [4] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, “Deep reinforcement learning: A brief survey,” *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.
- [5] N. Kohl and P. Stone, “Policy gradient reinforcement learning for fast quadrupedal locomotion,” in *IEEE International Conference on Robotics and Automation, 2004. Proceedings. ICRA’04. 2004*, vol. 3. IEEE, 2004, pp. 2619–2624.
- [6] A. Y. Ng, A. Coates, M. Diel, V. Ganapathi, J. Schulte, B. Tse, E. Berger, and E. Liang, “Autonomous inverted helicopter flight via reinforcement learning,” in *Experi-*

- imental robotics IX.* Springer, 2006, pp. 363–372.
- [7] M. Breyer, F. Furrer, T. Novkovic, R. Siegwart, and J. Nieto, “Flexible robotic grasping with sim-to-real transfer based reinforcement learning,” *ArXiv e-prints*, 2018.
- [8] E. Salvato, G. Fenu, E. Medvet, and F. A. Pellegrino, “Characterization of modeling errors affecting performances of a robotics deep reinforcement learning controller in a sim-to-real transfer,” in *2021 44rd International Convention on Information, Communication and Electronic Technology (MIPRO)*, 2021 in press.
- [9] N. Jakobi, P. Husbands, and I. Harvey, “Noise and the reality gap: The use of simulation in evolutionary robotics,” in *European Conference on Artificial Life*. Springer, 1995, pp. 704–720.
- [10] S. Höfer, K. Bekris, A. Handa, J. C. Gamboa, M. Mozifian, F. Golemo, C. Atkeson, D. Fox, K. Goldberg, J. Leonard, C. Karen Liu, J. Peters, S. Song, P. Welinder, and M. White, “Sim2real in robotics and automation: Applications and challenges,” *IEEE Transactions on Automation Science and Engineering*, vol. 18, no. 2, pp. 398–400, 2021.
- [11] J. B. Pollack, H. Lipson, S. Ficici, P. Funes, G. Hornby, and R. A. Watson, “Evolutionary techniques in physical robotics,” in *International Conference on Evolvable Systems*. Springer, 2000, pp. 175–186.
- [12] S. Nolfi, D. Floreano, O. Miglino, and F. Mondada, “How to evolve autonomous robots: Different approaches in evolutionary robotics,” in *Artificial life iv: Proceedings of the fourth international workshop on the synthesis and simulation of living systems*. MIT Press, 1994, pp. 190–197.
- [13] D. Floreano and J. Urzelai, “Evolution of plastic control networks,” *Autonomous robots*, vol. 11, no. 3, pp. 311–317, 2001.
- [14] C. Hartland and N. Bredeche, “Evolutionary robotics, anticipation and the reality gap,” in *2006 IEEE International Conference on Robotics and Biomimetics*. IEEE, 2006, pp. 1640–1645.
- [15] J.-B. Mouret and K. Chatzilygeroudis, “20 years of reality gap: a few thoughts about simulators in evolutionary robotics,” in *Proceedings of the Genetic and Evolutionary Computation Conference Companion*. ACM, 2017, pp. 1121–1124.
- [16] S. Koos, J.-B. Mouret, and S. Doncieux, “The transferability approach: Crossing the reality gap in evolutionary robotics,” *IEEE Transactions on Evolutionary Computation*, vol. 17, no. 1, pp. 122–145, 2012.
- [17] S. Kriegman, D. Blackiston, M. Levin, and J. Bongard, “A scalable pipeline for designing reconfigurable organisms,” *Proceedings of the National Academy of Sciences*, vol. 117, no. 4, pp. 1853–1859, 2020. [Online]. Available: <https://www.pnas.org/content/117/4/1853>
- [18] A. S. Badwe, R. D. Gudi, R. S. Patwardhan, S. L. Shah, and S. C. Patwardhan, “Detection of model-mismatch in mpc applications,” *Journal of Process Control*, vol. 19, no. 8, pp. 1305–1313, 2009.
- [19] M. Kano, Y. Shigi, S. Hasebe, and S. Ooyama, “Detection of significant model-plant mismatch from routine operation data of model predictive control system,” *IFAC Proceedings Volumes*, vol. 43, no. 5, pp. 685–690, 2010.
- [20] S. Selvanathan and A. Tangirala, “Diagnosis of poor control loop performance due to model- plant mismatch,” *Industrial & Engineering Chemistry Research*, vol. 49, no. 9, pp. 4210–4229, 2010.
- [21] L. E. Olivier and I. K. Craig, “Model-plant mismatch detection and model update for a run-of-mine ore milling circuit under model predictive control,” *Journal of Process Control*, vol. 23, no. 2, pp. 100–107, 2013.
- [22] W. Zhao, J. P. Queraltà, and T. Westerlund, “Sim-to-real transfer in deep reinforcement learning for robotics: a survey,” in *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE, 2020, pp. 737–744.
- [23] R. Bellman, “On the theory of dynamic programming,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 38, no. 8, p. 716, 1952.
- [24] L. Busoniu, R. Babuska, B. De Schutter, and D. Ernst, *Reinforcement learning and dynamic programming using function approximators*. CRC press, 2017.
- [25] D. P. Bertsekas, “Reinforcement learning and optimal control,” *Athena Scientific*, 2019.
- [26] A. C. Rodriguez, R. Parr, and D. Koller, “Reinforcement learning using approximate belief states,” in *Advances in Neural Information Processing Systems*, 2000, pp. 1036–1042.
- [27] H. Van Hasselt, “Reinforcement learning in continuous state and action spaces,” in *Reinforcement learning*. Springer, 2012, pp. 207–251.
- [28] M. P. Deisenroth, G. Neumann, and J. Peters, *A Survey on Policy Search for Robotics*, 2013.
- [29] B. Recht, “A tour of reinforcement learning: The view from continuous control,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 2, pp. 253–279, 2019.
- [30] C. G. Atkeson and S. Schaal, “Robot learning from demonstration,” in *ICML*, vol. 97. Citeseer, 1997, pp. 12–20.
- [31] P. Abbeel, A. Coates, M. Quigley, and A. Y. Ng, “An application of reinforcement learning to aerobatic helicopter flight,” in *Advances in neural information processing systems*, 2007, pp. 1–8.
- [32] H. Durrant-Whyte, N. Roy, and P. Abbeel, *Learning to Control a Low-Cost Manipulator Using Data-Efficient Reinforcement Learning*. MITP, 2012, pp. 57–64. [Online]. Available: <https://ieeexplore.ieee.org/document/6301026>
- [33] V. Gullapalli, J. A. Franklin, and H. Benbrahim, “Ac-

- quiring robot skills via reinforcement learning,” *IEEE Control Systems Magazine*, vol. 14, no. 1, pp. 13–24, 1994.
- [34] H. Miyamoto, S. Schaal, F. Gandolfo, H. Gomi, Y. Koike, R. Osu, E. Nakano, Y. Wada, and M. Kawato, “A kendama learning robot based on bi-directional theory,” *Neural networks*, vol. 9, no. 8, pp. 1281–1302, 1996.
- [35] R. Tedrake, T. W. Zhang, and H. S. Seung, “Learning to walk in 20 minutes,” in *Proceedings of the Fourteenth Yale Workshop on Adaptive and Learning Systems*, vol. 95585. Beijing, 2005, pp. 1939–1412.
- [36] J. Peters and S. Schaal, “Learning to control in operational space,” *The International Journal of Robotics Research*, vol. 27, no. 2, pp. 197–212, 2008.
- [37] ———, “Natural actor-critic,” *Neurocomputing*, vol. 71, no. 7-9, pp. 1180–1190, 2008.
- [38] K. Chatzilygeroudis, V. Vassiliades, F. Stulp, S. Calinon, and J. Mouret, “A survey on policy search algorithms for learning robot controllers in a handful of trials,” *IEEE Transactions on Robotics*, pp. 1–20, 2019.
- [39] F. Sadeghi and S. Levine, “Cad2rl: Real single-image flight without a single real image,” *arXiv preprint arXiv:1611.04201*, 2016.
- [40] K. Kang, S. Belkhale, G. Kahn, P. Abbeel, and S. Levine, “Generalization through simulation: Integrating simulated and real data into deep reinforcement learning for vision-based autonomous flight,” in *2019 International Conference on Robotics and Automation (ICRA)*, May 2019, pp. 6008–6014.
- [41] H. Hu, K. Zhang, A. H. Tan, M. Ruan, C. Agia, and G. Nejat, “A sim-to-real pipeline for deep reinforcement learning for autonomous robot navigation in cluttered rough terrain,” *IEEE Robotics and Automation Letters*, vol. 6, no. 4, pp. 6569–6576, 2021.
- [42] Y. Zhu, Z. Wang, C. Chen, and D. Dong, “Rule-based reinforcement learning for efficient robot navigation with space reduction,” *IEEE/ASME Transactions on Mechatronics*, pp. 1–1, 2021.
- [43] A. A. Rusu, M. Vecerík, T. Rothörl, N. Heess, R. Pascanu, and R. Hadsell, “Sim-to-real robot learning from pixels with progressive nets,” in *1st Annual Conference on Robot Learning, CoRL 2017, Mountain View, California, USA, November 13-15, 2017, Proceedings*, 2017, pp. 262–270. [Online]. Available: <http://proceedings.mlr.press/v78/rusu17a.html>
- [44] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, “Sim-to-real transfer of robotic control with dynamics randomization,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 1–8.
- [45] J. Matas, S. James, and A. J. Davison, “Sim-to-real reinforcement learning for deformable object manipulation,” in *2nd Annual Conference on Robot Learning, CoRL 2018, Zürich, Switzerland, 29-31 October 2018, Proceedings*, 2018, pp. 734–743. [Online]. Available: <http://proceedings.mlr.press/v87/matas18a.html>
- [46] S. James, P. Wohlhart, M. Kalakrishnan, D. Kalashnikov, A. Irpan, J. Ibarz, S. Levine, R. Hadsell, and K. Bousmalis, “Sim-to-real via sim-to-sim: Data-efficient robotic grasping via randomized-to-canonical adaptation networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 627–12 637.
- [47] W. Yuan, K. Hang, D. Kräig, M. Y. Wang, and J. A. Stork, “End-to-end nonprehensile rearrangement with deep reinforcement learning and simulation-to-reality transfer,” *Robotics and Autonomous Systems*, vol. 119, pp. 119–134, 2019.
- [48] R. Liu, F. Nageotte, P. Zanne, M. de Mathelin, and B. Dresp-Langley, “Deep reinforcement learning for the control of robotic manipulation: a focussed mini-review,” *Robotics*, vol. 10, no. 1, p. 22, 2021.
- [49] P. Christiano, Z. Shah, I. Mordatch, J. Schneider, T. Blackwell, J. Tobin, P. Abbeel, and W. Zaremba, “Transfer from simulation to real world through learning deep inverse dynamics model,” *arXiv preprint arXiv:1610.03518*, 2016.
- [50] A. Mandlekar, Y. Zhu, A. Garg, L. Fei-Fei, and S. Savarese, “Adversarially robust policy learning: Active construction of physically-plausible perturbations,” in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 3932–3939.
- [51] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta, “Robust adversarial reinforcement learning,” in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR.org, 2017, pp. 2817–2826.
- [52] M. Wulfmeier, I. Posner, and P. Abbeel, “Mutual alignment transfer learning,” in *Proceedings of the 1st Annual Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, S. Levine, V. Vanhoucke, and K. Goldberg, Eds., vol. 78. PMLR, Nov. 2017, pp. 281–290. [Online]. Available: <http://proceedings.mlr.press/v78/wulfmeier17a.html>
- [53] F. Muratore, F. Treede, M. Gienger, and J. Peters, “Domain randomization for simulation-based policy optimization with transferability assessment,” in *Conference on Robot Learning*, 2018, pp. 700–713.
- [54] J. Tan, T. Zhang, E. Coumans, A. Iscen, Y. Bai, D. Hafner, S. Bohez, and V. Vanhoucke, “Sim-to-real: Learning agile locomotion for quadruped robots,” in *Proceedings of Robotics: Science and Systems*, Pittsburgh, Pennsylvania, June 2018, pp. 10–20.
- [55] F. Golemo, A. A. Taiga, A. Courville, and P.-Y. Oudeyer, “Sim-to-real transfer with neural-augmented robot simulation,” in *Conference on Robot Learning*, 2018, pp. 817–828.
- [56] J. Van Baar, A. Sullivan, R. Cordorel, D. Jha,

- D. Romeres, and D. Nikovski, "Sim-to-real transfer learning using robustified controllers in robotic tasks involving complex dynamics," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 6001–6007.
- [57] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino, M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng, Q. Yuan, W. Zaremba, and L. Zhang, "Solving rubik's cube with a robot hand," *arXiv preprint*, 2019.
- [58] X. Pan, D. Seita, Y. Gao, and J. Canny, "Risk averse robust adversarial reinforcement learning," in *2019 International Conference on Robotics and Automation (ICRA)*, May 2019, pp. 8522–8528.
- [59] K. G. S. Apuroop, A. V. Le, M. R. Elara, and B. J. Sheu, "Reinforcement learning-based complete area coverage path planning for a modified htri hex robot," *Sensors*, vol. 21, no. 4, 2021. [Online]. Available: <https://www.mdpi.com/1424-8220/21/4/1067>
- [60] J. P. Hanna and P. Stone, "Grounded action transformation for robot learning in simulation," in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [61] J. Siekmann, K. Green, J. Warila, A. Fern, and J. Hurst, "Blind bipedal stair traversal via sim-to-real reinforcement learning," *arXiv preprint arXiv:2105.08328*, 2021.
- [62] Z. Li, X. Cheng, X. B. Peng, P. Abbeel, S. Levine, G. Bersek, and K. Sreenath, "Reinforcement learning for robust parameterized locomotion control of bipedal robots," *arXiv preprint arXiv:2103.14295*, 2021.
- [63] S. S. Mousavi, M. Schukat, and E. Howley, "Deep reinforcement learning: an overview," in *Proceedings of SAI Intelligent Systems Conference*. Springer, 2016, pp. 426–440.
- [64] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 5026–5033.
- [65] Blender Community, "Blender: Open Source 3D modeling suit." <http://www.blender.org/>, 2020, online.
- [66] Sourceforge, "SimSpark," <http://simspark.sourceforge.net/>, 2020, online.
- [67] Open Source Robotics Foundation, "GAZEBO: Robot simulation made easy." <http://www.gazebosim.org/>, 2020, online.
- [68] J. Lee, M. X. Grey, S. Ha, T. Kunz, S. Jain, Y. Ye, S. S. Srinivasa, M. Stilman, and C. K. Liu, "Dart: Dynamic animation and robotics toolkit." *J. Open Source Software*, vol. 3, no. 22, p. 500, 2018.
- [69] CM Labs, "Vortex Studio Real-time simulation and visualization software for system-level modeling of mechatronics and mechanical equipment." <https://www.cm-labs.com/vortex-studio/>, 2001, online.
- [70] E. Coumans and Y. Bai, "Bullet physics sdk," <http://bulletphysics.org>, 2016–2019.
- [71] ———, "Pybullet, a python module for physics simulation for games, robotics and machine learning," <http://pybullet.org>, 2016.
- [72] F. Xia, A. R. Zamir, Z. He, A. Sax, J. Malik, and S. Savarese, "Gibson env: Real-world perception for embodied agents," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2018, pp. 9068–9079.
- [73] B. Wymann, E. Espié, C. Guionneau, C. Dimitrakakis, R. Coulom, and A. Sumner, "Torcs, the open racing car simulator," *Software available at http://torcs.sourceforge.net*, vol. 4, no. 6, 2000.
- [74] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, "Progressive neural networks," *arXiv preprint arXiv:1606.04671*, 2016.
- [75] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sep. 2017, pp. 23–30.
- [76] S. James, A. J. Davison, and E. Johns, "Transferring end-to-end visuomotor control from simulation to real world for a multi-stage task," in *1st Annual Conference on Robot Learning, CoRL 2017, Mountain View, California, USA, November 13-15, 2017, Proceedings*, 2017, pp. 334–343. [Online]. Available: <http://proceedings.mlr.press/v78/james17a.html>
- [77] I. Mordatch, K. Lowrey, and E. Todorov, "Ensemblecio: Full-body dynamic motion planning that transfers to physical humanoids," in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2015, pp. 5307–5314.
- [78] J. M. Wang, D. J. Fleet, and A. Hertzmann, "Optimizing walking controllers for uncertain inputs and environments," *ACM Transactions on Graphics (TOG)*, vol. 29, no. 4, p. 73, 2010.
- [79] A. Vandesompele, G. Urbain, H. Mahmud, F. wyffels, and J. Dambre, "Body randomization reduces the sim-to-real gap for compliant quadruped locomotion," *Frontiers in Neurorobotics*, vol. 13, p. 9, 2019. [Online]. Available: <https://www.frontiersin.org/article/10.3389/fnbot.2019.00009>
- [80] S. Bhattacharyya, "Robust control under parametric uncertainty: an overview and recent results," *Annual Reviews in Control*, vol. 44, pp. 45–77, 2017.
- [81] J. Ackermann, *Robust control: the parameter space approach*. Springer Science & Business Media, 2012.
- [82] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.
- [83] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International conference on machine learning*, 2015, pp.

- 1889–1897.
- [84] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, O. P. Abbeel, and W. Zaremba, “Hindsight experience replay,” in *Advances in Neural Information Processing Systems*, 2017, pp. 5048–5058.
  - [85] N. Heess, J. J. Hunt, T. P. Lillicrap, and D. Silver, “Memory-based control with recurrent neural networks,” *arXiv preprint arXiv:1512.04455*, 2015.
  - [86] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” in *International conference on machine learning*, 2016, pp. 1928–1937.
  - [87] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.
  - [88] M. Večerík, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. Heess, T. Rothörl, T. Lampe, and M. Riedmiller, “Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards,” *arXiv preprint arXiv:1707.08817*, 2017.
  - [89] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, and S. Levine, “Scalable deep reinforcement learning for vision-based robotic manipulation,” in *Proceedings of The 2nd Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, A. Billard, A. Dragan, J. Peters, and J. Morimoto, Eds., vol. 87. PMLR, 29–31 Oct 2018, pp. 651–673. [Online]. Available: <http://proceedings.mlr.press/v87/kalashnikov18a.html>
  - [90] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, “Human-level control through deep reinforcement learning,” *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
  - [91] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, “Curriculum learning,” in *Proceedings of the 26th annual international conference on machine learning*. ACM, 2009, pp. 41–48.
  - [92] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. [Online]. Available: <http://arxiv.org/abs/1412.6572>
  - [93] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. J. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” in *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. [Online]. Available: <http://arxiv.org/abs/1312.6199>
  - [94] B. F. Hobbs and A. Hepenstal, “Is optimization optimistically biased?” *Water Resources Research*, vol. 25, no. 2, pp. 152–160, 1989.
  - [95] A. W. Moore and C. G. Atkeson, “Prioritized sweeping: Reinforcement learning with less data and less time,” *Machine learning*, vol. 13, no. 1, pp. 103–130, 1993.
  - [96] S. P. Choi and D.-Y. Yeung, “Predictive q-routing: A memory-based reinforcement learning approach to adaptive traffic control,” in *Advances in Neural Information Processing Systems*, 1996, pp. 945–951.
  - [97] J. Peng and R. J. Williams, “Efficient learning and planning within the dyna framework,” *Adaptive Behavior*, vol. 1, no. 4, pp. 437–454, 1993.
  - [98] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, “Image-to-image translation with conditional adversarial networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.
  - [99] W. Uther and M. Veloso, “Adversarial reinforcement learning,” Technical report, Carnegie Mellon University, 1997. Unpublished, Tech. Rep., 1997.
  - [100] M. L. Littman, “Markov games as a framework for multi-agent reinforcement learning,” in *Machine learning proceedings 1994*. Elsevier, 1994, pp. 157–163.
  - [101] D. L. Ma, S. H. Chung, and R. D. Braatz, “Worst-case performance analysis of optimal batch control trajectories,” in *1999 European Control Conference (ECC)*, Aug. 1999, pp. 3256–3261.
  - [102] D. L. Ma and R. D. Braatz, “Worst-case analysis of finite-time control policies,” *IEEE Transactions on Control Systems Technology*, vol. 9, no. 5, pp. 766–774, Sep. 2001.
  - [103] J. C. Doyle, K. Glover, P. P. Khargonekar, and B. A. Francis, “State-space solutions to standard  $h_{\text{sub}} 2/$  and  $h_{\text{sub}} \infty$  / control problems,” *IEEE Transactions on Automatic Control*, vol. 34, no. 8, pp. 831–847, Aug. 1989.
  - [104] I. Osband, C. Blundell, A. Pritzel, and B. Van Roy, “Deep exploration via bootstrapped dqn,” in *Advances in neural information processing systems*, 2016, pp. 4026–4034.
  - [105] M. E. Taylor and P. Stone, “Transfer learning for reinforcement learning domains: A survey,” *Journal of Machine Learning Research*, vol. 10, no. Jul, pp. 1633–1685, 2009.
  - [106] N. Hansen, “The cma evolution strategy: a comparing review,” in *Towards a new evolutionary computation*. Springer, 2006, pp. 75–102.
  - [107] H. Van Hasselt, A. Guez, and D. Silver, “Deep reinforcement learning with double q-learning,” in *Thirtieth AAAI conference on artificial intelligence*, 2016.
  - [108] G. Kahn, A. Villaflor, B. Ding, P. Abbeel, and S. Levine, “Self-supervised deep reinforcement learning with generalized computation graphs for robot

- navigation,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 1–8.
- [109] A. Farchy, S. Barrett, P. MacAlpine, and P. Stone, “Humanoid robots learning to walk faster: From the real world to simulation and back,” in *Proceedings of the 2013 international conference on Autonomous agents and multi-agent systems*, 2013, pp. 39–46.
- [110] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [111] Y. Duan, X. Chen, R. Houthooft, J. Schulman, and P. Abbeel, “Benchmarking deep reinforcement learning for continuous control,” in *International Conference on Machine Learning*, 2016, pp. 1329–1338.
- [112] V. Maheu, P. S. Archambault, J. Frappier, and F. Routhier, “Evaluation of the jaco robotic arm: Clinico-economic study for powered wheelchair users with upper-extremity disabilities,” in *2011 IEEE International Conference on Rehabilitation Robotics*, 2011, pp. 1–5.
- [113] Poppy Station, “Poppy Project: Open Source Robotics Platform.” <https://www.poppy-project.org/en/>, 2020, online.
- [114] F. Muratore, C. Eilers, M. Gienger, and J. Peters, “Bayesian domain randomization for sim-to-real transfer,” *arXiv preprint arXiv:2003.02471*, 2020.
- [115] N. Schweighofer and K. Doya, “Meta-learning in reinforcement learning,” *Neural Networks*, vol. 16, no. 1, pp. 5–9, 2003.
- [116] A. Gupta, R. Mendonca, Y. Liu, P. Abbeel, and S. Levine, “Meta-reinforcement learning of structured exploration strategies,” in *Advances in Neural Information Processing Systems*, 2018, pp. 5302–5311.
- [117] L. Kirsch, S. van Steenkiste, and J. Schmidhuber, “Improving generalization in meta reinforcement learning using learned objectives,” *arXiv preprint arXiv:1910.04098*, 2019.
- [118] M. Deisenroth and C. E. Rasmussen, “Pilco: A model-based and data-efficient approach to policy search,” in *Proceedings of the 28th International Conference on machine learning (ICML-11)*, 2011, pp. 465–472.
- [119] M. P. Deisenroth, D. Fox, and C. E. Rasmussen, “Gaussian processes for data-efficient learning in robotics and control,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 2, pp. 408–423, 2013.
- [120] T. Hospedales, A. Antoniou, P. Micaelli, and A. Storkey, “Meta-learning in neural networks: A survey,” *arXiv preprint arXiv:2004.05439*, 2020.
- [121] S. Tu and B. Recht, “Least-squares temporal difference learning for the linear quadratic regulator,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 5005–5014.



and Robotics.

ERICA SALVATO received the B.Sc. degree in Electronic Engineering in 2015 from the University of Messina (Italy) and the M.Sc. degree in Electrical and Control Systems Engineering in 2018 from the University of Trieste (Italy). She is currently a PhD student of the Department of Engineering and Architecture at the University of Trieste. Her research focuses on the Artificial Intelligence application as a systems control tool but includes also Control Theory, Machine Learning,



GIANFRANCO FENU received the “Laurea” degree (M.Sc. degree) in Electronic Engineering and the PhD degree in Information Engineering from the University of Trieste, respectively in 1996 and 2001. From October 1997 to March 1998 he was visiting PhD student at the Fachgebiet Mess- und Regelungstechnik (Institute for Measurement and Automatic Control) of the University of Duisburg (Germany), led by Prof. Paul M. Frank. From November 1999 he is an Assistant Professor in the

Engineering and Architecture of University of Trieste. His research interests include Control Theory, Fault Diagnosis, Machine Learning, and Robotics.



ERIC MEDVET received the degree in Electronic Engineering cum laude in 2004 and the PhD degree in Computer Engineering in 2008, both from the University of Trieste, Italy, where he is currently an Associate Professor in Computer Engineering and the director of the Evolutionary Robotics and Artificial Life Lab. Besides Evolutionary Robotics, his research interests include Genetic Programming and applications of Machine Learning.



FELICE ANDREA PELLEGRINO was born in Conegliano (Italy) in 1974, received the Laurea degree in Engineering from the University of Udine, Italy, in June of 2000 and the Ph.D. degree from the same university in May of 2005. From 2006 he is with Department of Engineering and Architecture of University of Trieste, currently as Associate Professor. He is Associate Editor of the IEEE Control Systems Letters and member of the Conference Editorial Board of the IEEE Control Systems Society and of the European Control Association. His research interests include Control Theory, Machine Learning, and Computer Vision.

...