

A Reinforcement Learning Architecture That Transfers Knowledge Between Skills When Solving Multiple Tasks

Paolo Tommasino, Daniele Caligiore, Marco Mirolli, and Gianluca Baldassarre

Abstract—When humans learn several skills to solve multiple tasks, they exhibit an extraordinary capacity to transfer knowledge between them. We present here the last enhanced version of a bio-inspired reinforcement-learning (RL) modular architecture able to perform skill-to-skill knowledge transfer and called transfer expert RL (TERL) model. TREL architecture is based on a RL actor-critic model where both actor and critic have a hierarchical structure, inspired by the mixture-of-experts model, formed by a gating network that selects experts specializing in learning the policies or value functions of different tasks. A key feature of TREL is the capacity of its gating networks to accumulate, in parallel, evidence on the capacity of experts to solve the new tasks so as to increase the responsibility for action of the best ones. A second key feature is the use of two different responsibility signals for the experts' functioning and learning: this allows the training of multiple experts for each task so that some of them can be later recruited to solve new tasks and avoid catastrophic interference. The utility of TREL mechanisms is shown with tests involving two simulated dynamic robot arms engaged in solving reaching tasks, in particular a planar 2-DoF arm, and a 3-D 4-DoF arm.

Index Terms—Autonomous robotics, bio-inspired modular neural architecture, catastrophic interference, cumulative learning, functioning and learning responsibility signals, mixture-of-expert networks, reaching tasks, transfer reinforcement learning (TRL).

I. INTRODUCTION

HUMANS, in particular children, learn multiple skills in a cumulative fashion, from simple to progressively more complex ones [1]. A striking feature of this cumulative learning process is its increasing speed. Thus, for example, children initially learn basic reaching and grasping

skills in months [2], [3], but later rapidly develop a repertoire of variants of the basic patterns in increasingly shorter times [4], [5].

This paper is in part motivated by the aim to understand the mechanisms through which children learn skills in an increasingly fast fashion. In this respect, our leading hypothesis is that the increasing learning speed observed in children relies on the transfer of knowledge from already learned skills to new skills to be acquired. The increasing learning rates might thus result from the fact that, thanks to knowledge transfer, learning processes can focus on acquiring only the aspects of the new skills that are novel with respect to the already acquired ones.

In particular, this paper describes a bio-inspired modeling architecture, developed within a reinforcement learning (RL) framework, that can be used to study skill-to-skill knowledge transfer. The system is called Transfer Expert RL (TERL) model. TREL has two capabilities that we deem essential for obtaining an increasingly fast learning of multiple skills (note that here we will use the term “skill” as a synonym of the RL “policy”):

- 1) the capacity to transfer knowledge from already acquired skills to new skills to be learned [6];
- 2) the capacity to store knowledge of the newly acquired skills so that it does not damage the knowledge on the already acquired ones; in other words, it can avoid “catastrophic interference” (or “catastrophic forgetting” [7], [8]).

The architecture, as further specified in Section II-A, is “bio-inspired” in that it is based on principles suggested by behavioral and brain mechanisms operating in organisms, and at the same time meets some constraints also faced by organisms when acquiring multiple skills.

The challenge faced here contributes to the developmental robotics overall objective of endowing robots with “developmental programs” supporting a prolonged autonomous development [9]–[14]. In particular, within this overall objective we face here the subproblem of how robots could acquire multiple skills in increasingly efficient ways.

The model we illustrate and analyze here is related to the machine learning literature on transfer RL (TRL) investigating how transferring knowledge between different domains and tasks (see [6], [15] for two reviews). The architecture described here focuses on a specific important problem of TRL concerning the development of algorithms capable of automatically

Manuscript received September 1, 2015; revised December 3, 2015 and April 1, 2016; accepted July 12, 2016. Date of publication October 18, 2016; date of current version June 10, 2019. This work was supported by the European Commission through the 7th Framework Programme (FP7/2007–2013), in particular by the ICT Challenge 2 “Cognitive Systems and Robotics” through the Project “IM-CLeVeR—Intrinsically Motivated Cumulative Learning Versatile Robots” under Grant Agreement ICT-IP-231722.

P. Tommasino is with Robotics Research Centre, Nanyang Technological University, Singapore (e-mail: paolo001@ntu.edu.sg).

D. Caligiore, M. Mirolli, and G. Baldassarre are with the Laboratory of Computational Embodied Neuroscience, Istituto di Scienze e Tecnologie della Cognizione, Consiglio Nazionale delle Ricerche, 00185 Rome, Italy (e-mail: daniela.caligiore@istc.cnr.it; marco.mirolli@istc.cnr.it; gianluca.baldassarre@istc.cnr.it).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCDS.2016.2607018

selecting already solved “source tasks” and the skills acquired to solve them, and to use such skills as starting points to best solve new “target tasks.” We will refer to this problem as the “source-task selection problem.” Here we sought solutions to the source-task selection problem for an important class of tasks: the tasks share the same primitive actions and the same environment transition function, but they involve different reward functions depending on the tasks [6], [15]. This class of tasks is biologically relevant as it reflects a common situation faced by real animals where they have to satisfy different needs (e.g., hunger, thirst, etc.): these needs generate different tasks in the same environment and these tasks are solved using the same set of “primitive actions” (i.e., muscle activations). Needs generate tasks as they represent reward functions that produce a reward signal when a need is present and is satisfied with the attainment of a suitable resource [16].

TERL architecture has been developed to its current form through a series of predecessor models. The idea of using the mixture-of-expert model to implement an actor-critic RL model was first proposed in [17] who used it to compose skills and to implement skill-based transfer between complex tasks (see Section IV). The idea of using the mixture-of-experts model principles and actor-critic models was also used in [18] and [19] to implement automatic task decomposition and multitask learning and was tested with a simple navigation task involving a 2-D scenario and discrete actions. The architecture was successively developed to work with continuous actions to control a 2-D dynamic simulated arm engaged in solving multiple reaching tasks [20] but its capacity to implement skill-to-skill transfer was not investigated. The latter system was further developed in [21] by decoupling the responsibility signals of experts used for functioning and for learning: this decoupling greatly enhances the capacity of the system to face the skill-to-skill TRL problem (see Section IV). Such work was the first to use the “TERL” acronym to indicate the system. Two previous works [21], [22] used TERL to show that the principles of the mixture-of-experts system applied to RL can be used to model and investigate the processes of *assimilation* and *accommodation* studied in developmental psychology [1].

The system described in this paper represents the culmination of these previous efforts into a new architecture that resembles the previous systems in many ways but also presents relevant innovations:

- 1) a systematic justification of the model ingredients from a computational perspective, and their link to relevant features of brain (Section II-A);
- 2) a fully revised formal description of the model (Section II-B);
- 3) a refinement of the new mechanism for decoupling the experts’ responsibility signals used for functioning and for learning, and other improvements (e.g., noise generation and a best-expert freezing mechanism; Section II-B): these mechanisms allow the system to train multiple experts (that we call “background copies”) for each task, so some of them can be later recruited to solve new tasks while avoiding catastrophic interference;

- 4) a systematic study and quantification of the skill-to-skill transfer capabilities of TERL involving sequential trials (Section III-B) or interleaved trials (Section III-D) related to different tasks;
- 5) tests on the scalability of the system to a large number of sequential tasks (Section III-C);
- 6) tests with the 4 DoF redundant robotic arm in addition to the planar arm (Section III-E);
- 7) a throughout discussion of other systems related to skill-to-skill transfer, from the option-framework systems to systems belonging to the MOSAIC-model family, and a clarification of their differences with respect to TERL (Section IV).

The paper is closed with the summary of the main results presented and the illustration of possible future developments of the system (Section V).

II. TERL: RATIONALE, ARCHITECTURE, FUNCTIONING, AND LEARNING

A. Principles Used to Build TERL

The principles followed to build TERL are inspired by the way in which animals’ brain and behavior might solve the source-task selection problem. At the same time, these principles allow TERL to face such a problem by fulfilling some constraints undergone by animals. This is intended to facilitate the future use of TERL as a model of multiple-skill learning in animals.

1) *Continuous States and Actions, and Function Approximation:* TERL is a RL system using function approximators. Function approximation is needed for RL to work in realistic setups, e.g., with embodied models and robots involving continuous state and action spaces [23]–[25]. The system described here is suitable to tackle such continuous state and action spaces (but not discrete actions as it is based on mechanisms mixing different actions). In particular, TERL uses linear function approximators (here also called experts) to encode policies and value function estimations. Linear function approximators have the computational advantages of simplicity of implementation, learning speed, convergence properties, and stability [24], [26]. However, they cannot solve linearly-separable problems: this limitation is commonly solved, as here, by recoding the input space with kernel functions [24], [27] and by leveraging modularity.

From a biological perspective, building systems capable of working with continuous state and action spaces captures the constraint faced by organisms that interact with the world through sensors and actuators involving continuous output and input signals. The use of linear function approximators increases the biological plausibility of models as they can use learning rules based on locally-available information [28].

2) *Actor-Critic Architecture:* TERL overall architecture is based on the *actor-critic* RL model [24], [29]. Actor-critic models use separate data structures for storing the action policy and the value function learned on the basis of the temporal difference (TD) algorithm [24]. The use of different data structures to represent policies and value functions was done as these often require different segmentations of the

problem space. Moreover, since the system “compiles” the gathered knowledge related to policies and evaluations into neural-network data structures (respectively mapping states into actions or into state-values), it is less memory demanding and more computationally efficient than other transfer models directly storing experience in the form of “state, action, state, and reward” tuples (we will further explain this in Section IV reviewing some examples of the latter models). Actor–critic models, which explicitly represent policies, are also suitable to face continuous action problems where drawing the policy from value functions is costly or complex [26].

From a biological perspective, the actor–critic architecture is considered one of the best models of the brain mechanisms underlying trial-and-error learning in organisms, in particular of basal ganglia [30]–[33]. Moreover, the TD learning rule at the core of the architecture learning processes accurately reproduces the dynamics of phasic dopamine, a neuromodulator playing a major role in trial-and-error learning of organisms [34], [35].

3) *Mixture-of-Experts to Solve the Source-Task Selection Problem:* The actor and the critic components of TERL are each based on the mixture-of-experts neural-network model [36], [37]. This model, proposed to solve supervised learning problems, was modified to work within a RL framework in the models preceding TERL (starting from [17] and [18]). The use of the modified principles of the mixture-of-expert architecture is at the basis of the capacity of TERL to solve the source-tasks selection problem. For this purpose, the actor and the critic are each formed by a hierarchy having two levels:

- 1) low-level “experts,” learning the policies or value functions needed to solve the novel tasks;
- 2) a high-level “gating network,” learning to select the experts (the term “hierarchy” is used here to refer to the feature of the architecture for which a high-level module controls the selection of lower-level modules).

The two gating networks of the actor and the critic select experts encoding already acquired skills to solve a new task on the basis of the accumulation of evidence related to their actual capacity to solve such new task. This accumulation of evidence is done in parallel for all available experts, so it scales up well with the number of experts. We shall see that this process results in an effective mechanism capable of quickly identifying the best expert for a given task, so after an initial phase of exploration the system will assign all the responsibility for action to that expert.

4) *No Prior Information About the Similarity Between Source and Target Tasks:* The scenario that was used here to test TERL capacity to solve the source-task selection problem has an important feature: the system is not given any prior information about the similarity between the source and the target tasks, so the only information it can use to solve the source-task selection problem is to sample the performance of the already acquired skills in the new target task. In particular, the model is tested here with robot arms that have to learn to solve different reaching tasks where the target object is located in different positions; these positions have different degrees of similarity between them, so they offer different

opportunities for knowledge transfer: the robot is, however, informed only on the task identity (i.e., at different trials it is told “this is task 1” or “this is task 2,” etc.), but not on the position of the target object, so the only way it has to select for reuse the already acquired skills is to try them out in the new target task. We adopted this demanding condition to be sure to find algorithms that are able to best use the knowledge on the performance of the already-acquired skills in the new task. Indeed, in some domains this is the only available information. As we further discuss in Section IV, the mechanisms examined here to exploit this type of information can be integrated with other mechanisms capable of exploiting information about the similarity between tasks that is available in some domains.

Biologically, the condition where the agent has no prior knowledge on the similarity between the source and target tasks captures the situations in which animals have information about their internal needs but have no clue on the external world conditions where they have to satisfy them [38], [39]. In these cases, the animals know that they have to solve different tasks (as they have to satisfy different needs, e.g., extinguish hunger or thirst), but they do not know if the behaviors to learn to solve them have to be similar or different (for example, if they have to perform similar behaviors to “collect” fruits with high content of nutrient or water; or if they have to perform very different behaviors to collect nutrient fruits or water from a lake). Instead, the condition in which the agent has information about the similarity between source and target tasks reflects situations where the tasks are defined in terms of “goals” (i.e., desired states of the external world) rather than “needs”: in this case, the similarities between the goals of tasks can be heuristically used to infer that the behaviors to learn to solve them might be similar or different [40]. Note that the restrictive condition faced here, where the system has no information about the similarity between tasks, is also relevant for robots as it reflects situations where the robot tasks are defined in terms of very abstract goals or in terms of internal variables of the robot (e.g., “satisfy the user’s request to drink” or “increase your energy level”), rather than specific goals (e.g., “reach the glass located in position x_1, y_1 on the table” or “reach the glass located in position x_2, y_2 on the table”); indeed, abstract goals and internal variables give little or no information about the similarity between the behaviors needed to solve the tasks, a situation captured by the restrictive condition used here.

5) *Redundant Experts to Face Catastrophic Interference:* The system is based on multiple possibly redundant, low-level experts, which are individually very simple but collectively capable of solving multiple, complex tasks. Computationally, this solution has several advantages such as fault tolerance and the possibility of using simple linear experts [41]–[43]. A further important advantage of this solution, exploited here, is the possibility to reduce catastrophic interference. Here the problem of catastrophic interference is generated by the requirement that the system has to use the same experts during its whole learning life, so knowledge of new tasks can disrupt knowledge of already solved tasks. This requirement captures an important constraint undergone by brain, namely the fact

that it cannot generate “new experts” (new neural resources) to face new tasks (note that for simplicity here we do not consider the developmental processes involving the physical structuring of the brain in time guided by epigenetic programs).

TERL faces the catastrophic interference problem on the basis of a novel mechanism: the decoupling of the responsibility signals used to select the experts for functioning and the responsibility signals used to modulate their learning. This innovation, introduced by TERL models and departing from other models based on the mixture-of-experts principles, allows the regulation of the learning rate of experts as desired, in particular independently of the size of their responsibility for acting, or evaluating states, in the task (the decoupling is explained in detail in Section II-B). This decoupling allows the formation of copies of the best policy expert and value expert trained when solving a given task. These “copy experts” can then be used to solve new tasks without disrupting the capacity of previously-trained experts to solve previous tasks.

In this respect, we mentioned above that the system is informed on the “identity” of the task being faced, so the reader might wonder what are the advantages of duplication with respect to a straightforward solution creating one new expert for each new task. These advantages are several (see also [42]). First, the system could be used to model the brain as this has most neural resources since birth and so it must be endowed with mechanisms to use them when learning an increasing number of tasks. Second, having all modules since the beginning allows the system to use experts already trained with similar tasks to solve new tasks rather than starting from scratch: this is at the basis of the transfer learning capabilities of the system. Third, the presence of copy experts gives robustness to the system, e.g., in case of failure of some experts. Fourth, although not done here, using the same set of experts from the beginning might allow the system to progressively organize them in space (e.g., within a 2-D grid) on the basis of their specialization and temporal activation [43].

We also anticipate that since the task tackled by the system changes repeatedly the model has to suitably regulate the exploratory noise. Indeed, noise should progressively decrease with the learning of a task and increase again in correspondence to new tasks. Since this problem was not central for this paper, we solved it with a simple approach that uses high levels of noise only in correspondence to tasks that the system finds difficult to solve and where a high degree of exploration is needed, and a low level of noise for tasks that are readily solved by the system.

From a biological perspective, although we do not have direct empirical evidence for neural duplication (i.e., the fact that a piece of knowledge encoded in a neural structure is reproduced in another neural structure; see [44]) we do have evidence for the partially modular organization of some areas of brain [41]:

- 1) motor cortex is based on neural columns, where different assemblies of columns might participate to encode multiple repertoires of skills, from very similar to very different [45], [46];
- 2) basal ganglia are organized in channels supporting the trial-and-error learning and selection of different actions and mental contents [47], [48].

B. TERL Architecture and Functioning

1) *Architecture Components*: The architecture of TERL, shown in Fig. 1, is based on two main components: 1) the *actor*, responsible for learning the action policy and 2) the *critic*, responsible for learning to approximate the value function. Each of the two components is formed by a number of *expert networks*, learning the policy or the value function, and a *gating network*, learning to select the experts given the task (note that here for simplicity the number of experts of the actor and of the critic is the same but it could be different). The functioning and learning processes of the architecture are now explained in detail (the Appendix gives its parameters).

2) *Inputs and Outputs of the Architecture Components*: The experts receive as input the current environment state encoded in an expanded space of features (here Gaussian basis functions with centers equally distributed over the input space, see below for details) and have no information about the goal or task. In particular, each environmental state is encoded with an I -dimensional vector of continuous variables $\mathbf{s} = \langle s_1, s_2, \dots, s_i, \dots, s_I \rangle$. In the experiments reported here, \mathbf{s} encoded the arm joint angles, so $I = 2$ in the 2-D planar arm simulations, and $I = 4$ in the simulations involving the 3-D humanoid robot arm. The state vector is expanded into a D -dimensional vector of continuous variables (features) $\mathbf{f} = \langle f_1, f_2, \dots, f_d, \dots, f_D \rangle$ where $D \gg I$. Features are computed through normalized Gaussian basis functions

$$f_d = \frac{e^{-\frac{\|\mathbf{s} - \mathbf{s}_d\|^2}{\sigma_f^2}}}{\sum_{d=1}^D e^{-\frac{\|\mathbf{s} - \mathbf{s}_d\|^2}{\sigma_f^2}}} \quad (1)$$

where \mathbf{s}_d is the preferred state vector of feature d and σ_f^2 is the width of the Gaussians (in degrees). The preferred vectors of the features lay on the vertices of a regular grid overlapped with the state space. Notice that the experts are not informed on the task to solve but only on the posture of the arm, so different tasks require different experts. This limited input eased the analysis of the transfer capabilities of the system, and reflects the organization of the brain where primary motor cortex is mainly informed as to the limb posture but not as to the overall tasks the system is accomplishing [49]–[51]. This issue is further discussed in Section V.

The output of the actor experts is a J -dimensional vector $\mathbf{a} = \langle a_1, a_2, \dots, a_j, \dots, a_J \rangle$ encoding the controlled continuous variables. Here, the vector encoded the angles of the arm joints, so $J = 2$ in the 2-D planar arm simulations, and $J = 4$ in those involving the 3-D humanoid robot arm. As further explained below, these desired angles, varying at each time step, are used as the equilibrium points (EPs) of proportional derivative (PD) controllers that produce the torques controlling the dynamic-arm joints.

The gating networks are informed only as to the identity of the task to pursue. The tasks used here are of the type “reach and touch object A in space” through a dynamic arm. This tells the gating networks which task is being solved (e.g., task A, task B, etc.) but it does not furnish any information about the similarity between tasks. Formally, the current task identity is encoded as a K -dimensional vector

Critic (C)

Actor (A)

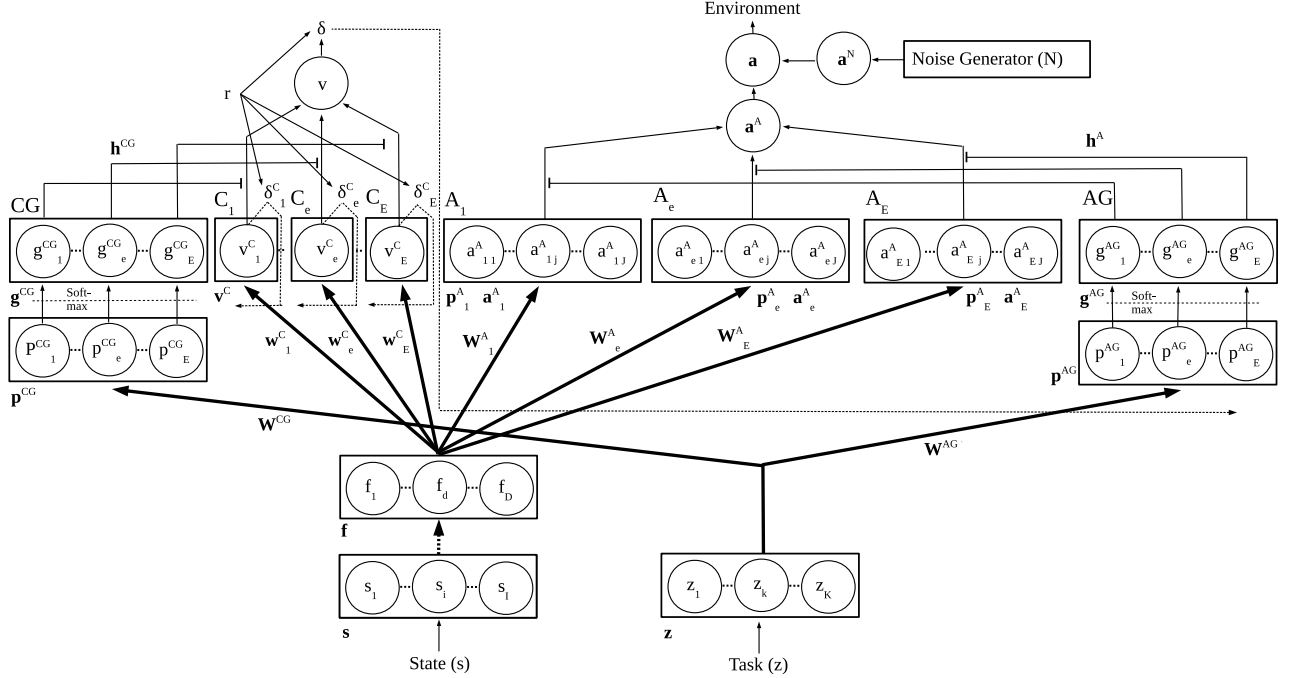


Fig. 1. Architecture of TERL, formed by a modular critic and a modular actor. The figure gives an overview of the elements composing the system that are explained in detail in the main text. The critic is formed by a gating network (CG) and a number of experts (C_e). Also the actor is formed by a gating network (AG) and a number of experts (A_e). Thin arrows: one-to-one or many-to-one hardwired connections set to one (for information relay); dotted bold arrow: all-to-all hardwired connections (for computing features); thin dashed arrows: learning signals that cause the update of all crossed connections (notice how the TD-errors of the critic experts are locally used to train themselves, whereas the critic selector does not need any TD-error to learn as it learns based on the mixture-of-experts supervised learning rule; the general TD-error is used to train the actor experts and the actor selector); bold arrows: trained all-to-all connections; and flat-head arrows: one-to-one hardwired multiplicative gating connections set to one (for carrying the responsibility signals). Letters indicate the symbols used in the mathematical description of the system in the text. Bold letters represent vectors (for simplicity, the vectors of biases are omitted).

$\mathbf{z} = \langle z_1, z_2, \dots, z_k, \dots, z_K \rangle$. The \mathbf{z} vectors, as many as the tasks, are orthogonal and are formed by binary elements all activated with 0 with the exception of one element, activated with 1, corresponding to the current task. As mentioned above, the gating networks are only informed on the task but not on the posture of the arm. As for the input to the experts, this limited input eased the analysis of the transfer capabilities of the system and reflects the organization of the brain where higher cortical areas are mainly informed on abstract goals and overall motivations [52]–[54].

The output of the gating networks are two E -dimensional vectors, where E represents the number of experts: $\mathbf{g}^{CG} = \langle g_1^{CG}, g_2^{CG}, \dots, g_e^{CG}, \dots, g_E^{CG} \rangle$ for the critic gating network and $\mathbf{g}^{AG} = \langle g_1^{AG}, g_2^{AG}, \dots, g_e^{AG}, \dots, g_E^{AG} \rangle$ for the actor gating network. These vectors encode the responsibility signals of the respective experts.

Finally, the system receives a reward $r = 1$ when the task is accomplished with success (here when the arm succeeds to touch the target object). The system receives a reward $r = 0$ at any other time step.

3) *Actor Gating Network—Functioning:* The actor gating network, denoted with AG, receives as input the task identity \mathbf{z} and returns as output the responsibility signals of the experts. These are encoded in the vector \mathbf{g}^{AG} with elements g_e^{AG} corresponding to experts $e = 1, 2, \dots, E$. The responsibility signal of an expert represents the “prior probability” that such expert is the best one to solve the current task among the

actor experts. The priors are used: 1) to set the contribution of each expert to produce actions and 2) to establish, after a suitable transformation (see below), the size of its learning update.

The activation potential of the output units of AG, denoted with the vector \mathbf{p}^{AG} with elements p_e^{AG} , is computed as

$$\mathbf{p}^{AG} = \mathbf{W}^{AG} \cdot \mathbf{z} \quad (2)$$

where \mathbf{W}^{AG} is the matrix of the connection weights of the network. The gating networks of the model do not have a bias as they should not have a tendency to select specific default experts.

As in the mixture-of-experts model, the prior responsibilities g_e^{AG} are computed with the soft-max function

$$g_e^{AG} = \frac{e^{p_e^{AG}/\kappa}}{\sum_{e=1}^E e^{p_e^{AG}/\kappa}} \quad (3)$$

where κ is the temperature regulating the slope of the soft-max. The soft-max guarantees that $\sum_{e=1}^E g_e^{AG} = 1$ so the responsibility signals can be interpreted as probabilities.

4) *Actor Experts—Functioning:* The actor experts, denoted with A_e , are a set of E neural networks each of which gets as input the state features \mathbf{f} and encodes a possible action with logistic output units. The activation potential of the output units of expert A_e , denoted with the vector \mathbf{p}_e^A with elements p_{ej}^A

(e.g., corresponding to the robot joint angles), is computed as

$$\mathbf{p}_e^A = \mathbf{W}_e^A \cdot \mathbf{f} + \mathbf{b}_e^A \quad (4)$$

where \mathbf{W}_e^A is the connection weight matrix of the expert A_e and \mathbf{b}_e^A the vector of biases. The activation of the output units of expert A_e , denoted with the vector \mathbf{a}_e^A with elements α_{ej}^A , are computed with a logistic function

$$\alpha_{ej}^A = \frac{1}{1 + e^{-\mathbf{p}_{ej}^A}}. \quad (5)$$

As we shall see below, \mathbf{a}_e^A is mixed with noise and so should be interpreted as the mean of a probability distribution of the system action.

5) *Global Actor Action*: The global action \mathbf{a}^A of the actor (before the addition of noise) is formed by mixing the actions \mathbf{a}_e^A of experts based on their priors

$$\mathbf{a}^A = \sum_{e=1}^E \left[g_e^{\text{AG}} \cdot \mathbf{a}_e^A \right]. \quad (6)$$

To foster exploration, the action of the system to be executed, denoted with \mathbf{a} , is obtained by mixing the global actor action \mathbf{a}^A with a noisy action, denoted with \mathbf{a}^N , produced by the noise generator component explained below. Thus, at a specific time t

$$\mathbf{a}_t = u_t \cdot \mathbf{a}_t^A + (1 - u_t) \cdot \mathbf{a}_t^N \quad (7)$$

where u_t is a variable regulating the exploitation-exploration level in different moments of the trial and explained below in the paragraph on the noise generator. This equation implies that with high values of u_t the performed action is strongly based on the actor action, whereas with low values it is strongly based on noise.

6) *Critic Gating Network—Functioning*: The critic gating network, denoted with CG receives as input the task identity \mathbf{z} and returns as output the responsibility signals of the experts of the critic. These responsibility signals are encoded in the vector \mathbf{g}^{CG} with elements g_e^{CG} corresponding to the experts e . Similar to what happens for the actor, these responsibility signals represent the prior probability that each expert is the best one in estimating the value function, given the current task and the current policy, and are used for the critic functioning and, after a suitable transformation illustrated below, for its learning.

The activation potential of the output units of CG, encoded in the vector \mathbf{p}^{CG} with elements p_e^{CG} , is computed as

$$\mathbf{p}^{\text{CG}} = \mathbf{W}^{\text{CG}} \cdot \mathbf{z} \quad (8)$$

where \mathbf{W}^{CG} is the matrix of connection weights of the network.

As with the actor, the prior responsibilities g_e^{CG} are computed with the soft-max function

$$g_e^{\text{CG}} = \frac{e^{p_e^{\text{CG}}/\kappa}}{\sum_{e=1}^E e^{p_e^{\text{CG}}/\kappa}}. \quad (9)$$

7) *Critic Experts—Functioning*: The critic experts, denoted with C_e , are a set of E neural networks each of which gets as input the state features \mathbf{f} and returns the state value with a linear output unit. The activation v_e^C of the output unit of expert C_e is computed as

$$v_e^C = \mathbf{w}_e^C \cdot \mathbf{f} + b_e^C \quad (10)$$

where \mathbf{w}_e^C is a row vector encoding the connection weights of the expert and b_e^C its bias. Note how the actor experts use a logistic function to produce the output so this can be mapped to a limited range of the robot's joint angles (see Sections III-B and III-E). The critic experts instead use a linear output unit so they can produce state values outside the (0, 1) range.

8) *Critic—Global Value*: The global value of the critic, denoted with v , is computed by mixing the values v_e^C of the experts based on their priors

$$v = \sum_{e=1}^E \left[g_e^C \cdot v_e^C \right]. \quad (11)$$

C. TERL Learning

1) *Critic—Global TD-Error*: The global value v_t at time t is used to compute the global TD-error δ_t of the critic

$$\delta_t = \begin{cases} 0 & \text{if } t = 0 \\ (r_t + \gamma \cdot v_t) - v_{t-1} & \text{if } 0 < t < T \\ r_t - v_{t-1} & \text{if } t = T \end{cases} \quad (12)$$

where r_t is the reward at time t , and γ is a discount factor. In this formula, v_t is set to zero at the end of the trial ($t = T$) to take into account the episodic nature of the RL problems considered here. As we shall see below, δ_t , related to the performed actions, is used to train the actor experts and the actor gating network.

2) *Critic Experts—TD-Error*: The TD-error of each critic expert, denoted with $\delta_{e_t}^C$, is computed as follows on the basis of the expert value $v_{e_t}^C$ at time t :

$$\delta_{e_t}^C = \begin{cases} 0 & \text{if } t = 0 \\ (r_t + \gamma \cdot v_{e_t}^C) - v_{e_{t-1}}^C & \text{if } 0 < t < T \\ r_t - v_{e_{t-1}}^C & \text{if } t = T. \end{cases} \quad (13)$$

This value represents the expert error in estimating the current state value and so it is computed on the basis of the expert current and past values and the experienced reward. For this reason, as we shall see below, it is used both to train the critic experts and to update the prior probability estimate g_e^{CG} that each critic expert is the best to evaluate the actor policy used to solve the current task.

3) *Actor Gating Network—Learning*: We now show how the estimate of the goodness of experts in solving the new task is updated on the basis of a process of accumulation of information and the TD-error. At each time step, each expert's responsibility signal is updated on the basis of the new evidence (likelihood) that the expert contributed to determine the executed action \mathbf{a} in the current task (recall that such action is noisy due to the \mathbf{a}^N component). Formally, the likelihood $l_{e_t}^{\text{AG}}$ is computed on the basis of a Gaussian function: this makes

the likelihood of each expert inversely related to the distance between that expert action \mathbf{a}_e^A and the executed noisy action \mathbf{a}

$$l_{e_t}^{AG} = e^{-\frac{1}{2\sigma_{AG}^2} \|\mathbf{a}_{t-1} - \mathbf{a}_{e_{t-1}}^A\|^2} \quad (14)$$

where σ_{AG}^2 is the Gaussian width parameter. Thus, the more similar the expert action to the executed action, the higher the likelihood that it had a large responsibility in generating it.

The likelihood is then used to compute the posterior probabilities with the Bayes rule

$$h_{e_t}^{AG} = \frac{l_{e_t}^{AG} \cdot g_{e_{t-1}}^{AG}}{\sum_{e=1}^E [l_{e_t}^{AG} \cdot g_{e_{t-1}}^{AG}]}. \quad (15)$$

Notice that this formula implies that $\sum_{e=1}^E [h_e^{AG}] = 1$ and that $h_{e_t}^{AG} > g_{e_{t-1}}^{AG}$ for experts whose action \mathbf{a}_e^A was more similar to the executed action \mathbf{a} relative to the action of all other experts.

The update of the connection weights of the actor gating network AG, and hence the responsibility signals of the actor experts, is based on the difference between the posteriors and the priors but also on the overall TD-error δ

$$w_{ek_t}^{AG} = w_{ek_{t-1}}^{AG} + \eta^{AG} \cdot \delta_t \cdot (h_{e_{t-1}}^{AG} - g_{e_{t-1}}^{AG}) \cdot z_{k_{t-1}} \quad (16)$$

where η^{AG} is a learning rate. The formula, which is based on the mixture-of-expert formula changed to take into consideration the RL framework used here (in particular, the TD-error), changes the connection weights so that the prior responsibility signal of experts, $g_{e_{t-1}}^{AG}$, gets closer to or away from the posterior responsibility signal, $h_{e_{t-1}}^{AG}$, depending on the sign and size of the TD-error, δ_t . In particular, the formula implies that the responsibility signals of experts whose actions are similar to the executed action (hence have a high $h_{e_{t-1}}^{AG}$) are increased, but only if the executed action produced a *positive TD-error*, i.e., if the executed action had positive effects. The responsibility signals of experts that produced an action dissimilar from the executed one are instead decreased. On the contrary, in the case of a *negative TD-error* the responsibility signals are updated in the *opposite direction*, yielding a lower responsibility for experts whose action was more similar to the executed one as the latter had negative effects. The responsibility signals of experts that produced an action dissimilar from the executed one are instead increased. Overall, the rule ensures that the responsibility signals are progressively increased for the experts that contribute to achieve the highest reward in the current task as the effect of an accumulation of evidence on the fact that they are better than the other experts in solving such task.

4) *Actor Experts—Learning*: The mixture-of-experts model uses the posterior probabilities to scale the learning rate of each expert so that the best experts are not only assigned the highest responsibility during functioning, but also an update proportional to such posteriors. One of the most important departures of the current system from this strategy is based on the decoupling of the responsibility signals used for functioning and those used for learning. Indeed, a number of pilot experiments using several variants of the mixture-of-experts strategy showed that the classical strategy leads not only to a desirable rapid increase of the responsibility of the best expert

toward the maximum possible value of one, but also to a rapid decrease to zero of the learning rate of all other experts. In our case, the latter outcome is detrimental because when the system has to learn a new task similar to a previously acquired one, it recruits the expert trained to solve the latter, modifies it, and so loses the capacity to solve it (catastrophic forgetting). Instead, the use of fixed responsibilities for learning allows the training of multiple experts similar to the best one: the size of those responsibilities allows the regulation of the number and learning rate of the copy experts. These copy experts can be later recruited to solve similar tasks without destroying the capacity to solve the previously solved source tasks.

To implement this idea, we ranked the experts on the basis of the decreasing value of their priors g_e^{AG} and then assigned them fixed learning rates based on the ranks. In particular, here we used the experts' ranks, encoded with $g_e^{rAG} \in \{0, 1, 2, \dots, E-1\}$, to compute their learning responsibility signals g_e^{lAG} as follows:

$$g_e^{lAG} = \frac{\zeta^{-g_e^{rAG}}}{\sum_{e=1}^E [\zeta^{-g_e^{rAG}}]} \quad (17)$$

where ζ was a coefficient. For example, setting $\zeta = 6$, as done here, implies that the learning responsibility signals g_e^{lAG} are, when ordered by rank, equal to the following values: (0.834, 0.139, 0.023, 0.004, 0, 0, 0, 0, ...)

Note that:

- 1) the algorithm will train one main expert and a given number of copies (here three) of it, whereas other experts will not be trained so avoiding disrupting their knowledge;
- 2) the experts that are trained are those that have the highest priors, i.e., those that are the best in solving the current task;
- 3) responsibility signals different from those used here can be used, including some with $\sum_{e=1}^E [g_e^{lAG}] \neq 1$, to obtain a desired number of background expert copies and set their learning rates. For example, one might establish a fixed learning rate for the first k experts. The higher the number of background experts copies are formed, the higher the redundancy of the system and all related advantages and disadvantages.

The connection weights of experts are then updated on the basis of g_e^{lAG} as follows:

$$\begin{aligned} m_{ejd_t}^A &= (a_{j_t} - a_{ej_t}^A) \cdot (a_{ej_t}^A \cdot (1 - a_{ej_t}^A)) \cdot f_{d_t} \\ w_{ejd_t}^A &= w_{ejd_{t-1}}^A + \eta^A \cdot \delta_t \cdot g_{e_{t-1}}^{lAG} \cdot m_{ejd_{t-1}}^A \end{aligned} \quad (18)$$

where $(a_{ej_t}^A \cdot (1 - a_{ej_t}^A))$ is the derivative of the logistic transfer function, and η^A is a learning rate. The rule is based on the gradient-descent formula of the mixture-of-experts model applied to neural logistic output units, but the size of the update is also weighted by the RL TD-error. The rule implies the following:

- 1) the size of the update is higher for experts that produce an action more similar to the action actually executed (i.e., a higher $|a_{j_t} - a_{ej_t}^A|$);
- 2) \mathbf{a}_e^A moves toward \mathbf{a} when $0 < \delta_t$ as the noisy action \mathbf{a} has been better than expected;

- 3) \mathbf{a}_e^A moves away from \mathbf{a} when $\delta_t < 0$ as the noisy action \mathbf{a} has been worse than expected;
- 4) if δ_t is close to zero, for example when the system has converged to good solutions, the action is not updated because the noisy action \mathbf{a} did not perform better or worse than expected;
- 5) if some of the expert copies developed for the task are later recruited to learn a similar new task, the capacity to solve the current task tend to not be impaired as the main expert used to solve it will not be recruited (see also the subsection “additional mechanisms to preserve the best experts” below).

5) *Critic Gating Network—Learning*: We now show how the responsibility signals of the experts of the critic are updated on the basis of mechanisms that accumulates evidence on the capacity of experts to face the given task similarly to what is done in the mixture-of-experts model for supervised learning problems. This is possible as the problem of learning the value function can still be seen as a supervised learning problem with the difference being that the value to use as the desired output is not given and so, as common in RL [24], it has to be formulated on the basis of the reward signal and the estimate of the next state value.

The new evidence (likelihood) that an expert is the best one to estimate the value function is computed as a Gaussian function of that expert TD-error, δ_e^C , as this represents its error in predicting the incoming rewards

$$l_{e_t}^{CG} = e^{-\left(\frac{(\delta_{e_t}^C)^2}{2 \cdot \sigma_{CG}^2}\right)} \quad (19)$$

where σ_{CG}^2 is the Gaussian width.

Similarly to the actor gating network, the likelihood is used to compute the posterior probabilities with the Bayes rule

$$h_{e_t}^{CG} = \frac{l_{e_t}^{CG} \cdot g_{e_{t-1}}^{CG}}{\sum_{e=1}^E [l_{e_t}^{CG} \cdot g_{e_{t-1}}^{CG}]}. \quad (20)$$

The formula implies that $h_{e_t}^{CG}$ is higher than $g_{e_{t-1}}^{CG}$ for experts that have a low TD-error in comparison to the TD-error of the other experts. Also here $\sum_{e=1}^E [h_e^{CG}] = 1$.

The update of the connection weights of the critic gating network is then computed as in the mixture-of-expert model

$$w_{ek_t}^{CG} = w_{ek_{t-1}}^{CG} + \eta^{CG} \cdot (h_{e_t}^{CG} - g_{e_{t-1}}^{CG}) \cdot z_{k_{t-1}} \quad (21)$$

where η^{CG} is a learning rate. The formula implies that the responsibility signals of experts with relatively smaller TD-errors are increased while those of experts with larger TD-errors are decreased, leading to larger responsibilities for experts that produce more accurate estimates of the state values for the current task. Overall, as for the actor experts, the rule ensures that the responsibility signals are progressively increased for the experts that contribute to achieving the highest reward in the current task as the effect of an accumulation of evidence that they are better than the other experts in solving such task.

6) *Critic Experts—Learning*: Due to the same reasons explained for the actor experts, we also implemented a decoupling between the functioning and the learning responsibility

signals of the critic experts. The expert ranks g_e^{rCG} were computed as for the actor and used in the same way (17) to compute the learning responsibility signals g_e^{lCG} used to modulate the learning rate of the experts.

Each expert was then trained using the TD-learning RL formula and its own TD-error, δ_e^C , to enforce the self-consistency of the value estimates produced by the expert

$$w_{ed_t}^C = w_{ed_{t-1}}^C + \eta^C \cdot \delta_{e_t}^C \cdot g_e^{lCG} \cdot f_{d_{t-1}}. \quad (22)$$

7) *Additional Mechanisms to Preserve the Best Experts*: Note that in some conditions that are more challenging for catastrophic forgetting, additional mechanisms can be used to further reduce interference between tasks. For example, in the “sequential conditions” tested here (Section III-B) we found it very useful to use a mechanism for which when the prior of an expert overcomes a certain (high) threshold, indicating that much evidence has been accumulated that such expert is the best for the current task, its learning rate is reduced to low values (here to zero for simplicity). Such low learning rate facilitates the preservation of the best experts found for previous tasks thus leading the system to update especially the background copies to solve the new tasks. Notice that alternative mechanisms might be investigated in the future to preserve the best expert found for a given task.

8) *Noise Generator and Executed Noisy Action*: One of the major problems of RL is the regulation of exploratory noise, also known as the exploration-exploitation dilemma. The heuristic solution most commonly used in the literature is to progressively lower exploration noise with learning so as to augment the exploitation versus exploration while the system becomes progressively more skilled (i.e., progressively better at achieving reward) [24]. Here we adapted a very simple mechanism that is based on this idea, but note that the other mechanisms of TERL could work with any other method one might use for noise generation. The mechanism takes into consideration the fact that the new task to solve can be very similar, or even identical, to already solved tasks: in these cases, exploratory noise has to be very low since the initial trials of learning of the new task because the system already has some experts that can be readily used to solve it. For this purpose, we started each trial with low noise and then increased it during the trial. In this way, if the system is capable of solving the task quickly, it is not disturbed by noise; instead, if it takes a long time to solve it during the trial, then noise increases causing high exploration. Overall, the mechanism implies that with not-yet-learned tasks noise decreases progressively, on average, over the tasks (as in standard RL problems). Instead, with already-learned or easy tasks, for example with tasks very similar to already solved ones, noise is immediately low or decreases fast, on average over the trials, with learning.

To implement this idea in detail, we divided each trial into two phases. In a first “exploitation phase” ($t = 0, 1, \dots, M$) noise is very low. In a second “exploration phase” ($t = M + 1, M + 2, \dots, T'$, where T' is the trial timeout; note that T , the trial duration, is possibly shorter than T' in case of reward accomplishment) noise gets progressively higher. Notice that M should be set to a value equal or higher than the time assumed to be sufficient to solve the task (a similar

parameter has to be set when using the progressive decrease of noise mechanism commonly employed in the RL literature and that could not be used here). In detail, the mechanism is implemented by setting the value of u_t regulating the mixture between the noisy action and the actor action [see (7)] as follows:

$$u_t = \begin{cases} v & \text{if } 0 < t \leq M \\ u_{t-1} - \beta u_{t-1} & \text{if } M < t \leq T' \end{cases}$$

where v is the value of u in the exploitation phase and β is a time constant regulating the dynamics of u in the exploration phase. A small amount of noise is still present in the initial exploitation phase to have some refinement of the policy even during such phase, whereas β is set to a value that implies that noise rapidly increases during the second exploration phase.

As done in [55], using a noise filter is important to control dynamic robotic arms because their physical inertia tends to cancel out white noise and also avoids issuing jerky commands to the robot. For this reason, the noisy action \mathbf{a}^N that is mixed with the actor action to produce the executed action (7) is produced by a noise generator component, denoted with N , through a first order filter

$$\mathbf{a}_t^N = \mathbf{a}_{t-1}^N + \tau(-\mathbf{a}_{t-1}^N + \mathbf{n}_t) \quad (23)$$

where τ represents the filter time constant ranging in $(0, 1)$ and \mathbf{n}_t is a vector having each element randomly drawn on the basis of a uniform probability distribution in $[-\epsilon, +\epsilon]$ at each time step.

III. RESULTS

A. Measuring Transfer Quality

To measure the potential for transfer of TERL, its performance was compared with two alternative systems furnishing the upper and lower bounds of learning speed curves (see [6]):

- 1) a system that learns to solve the new task based on a certain source task: this may result in an advantage or a disadvantage for learning the new task depending on its similarity to the source task;
- 2) a system that learns to solve any new task from scratch, without the possibility of skill transfer.

The two systems are also important as they view the multiple-task problem as either a set of separated Markov decision processes (MDPs; [24]), each corresponding to a single task, or as a whole MDP, where the problem is to maximize the reward given multiple tasks. The possibility of viewing the problem in these two ways is important for distinguishing the system discussed here from other related systems (see Section IV).

The first system, called “SIM” (which stands for “simple system”), is simply formed by one critic expert and one actor expert and no gating networks (each of the two experts is a linear function approximator getting as input the Gaussian features \mathbf{f} used to encode the arm posture). SIM is capable of learning only a policy mapping $\pi_k : S \times A \rightarrow [0, 1]$ to solve one specific MDP problem k denoted with $\langle S, A, T, R_k \rangle$ (where S denotes the set of states, A the set of actions, T the transition function mapping the current state and action to a probability distribution over the next states, and R_k the

reward function of task k mapping the current state, action and resulting state to a reward value). This, and the fact that SIM did not receive an input on the solved task, implies that when it learns a new task it fully forgets the previously learned one. Most tests illustrated below required the solution of tasks in sequence: in this condition, SIM allows the identification of the advantages and disadvantages of attempting to solve each new task *starting from the previous one*. SIM used identical parameters as the experts of TERL.

The second system, called “EXP” (which stands for “system with an input expanded with information about the identity of the task”), as SIM is again formed only by one actor expert and one critic expert and no gating networks. However, in EXP such experts have enough computational capabilities to solve all tasks. In particular, in EXP the two experts receive information not only on the arm posture (features \mathbf{f}), but also on the task identity (\mathbf{z}): EXP uses the information about the task identity as an additional input to the experts together with the posture features (hence forming an “expanded” input) rather than as information sent to the gating networks (not present in EXP) as in TERL. Technically, the task identity is used as a further dimension of the input problem space so that the Gaussian basis functions encoding the features (\mathbf{f}) are reproduced for a number of times equal to Z (the number of tasks) and activated only when the corresponding feature z is active. For this reason, in the tests EXP allows the identification of the advantages/disadvantages of solving new tasks *starting from scratch* without any opportunity of transferring knowledge from previously acquired ones (no generalization). At the same time, the system allows the measure of performance in complete absence of catastrophic interference. EXP used identical parameters as the experts of TERL.

Differently from SIM, but the same as EXP, TERL is informed of the task at hand. However, unlike EXP, TERL uses the information regarding the task identity to take advantage of the fact that different tasks share the same states, actions, and transition function to transfer knowledge from acquired tasks to new tasks to be learned.

The minimal performance that we require from TERL is that it converges toward the best MDP solution, i.e., *performance optimality* [56]. The core of the skill-transfer problem is, however, captured by *learning optimality* [56], namely the fact that an algorithm achieves a high (in theory the best) learning speed on new tasks thanks to the previously acquired ones. The literature on TRL measures this capacity for transfer through three metrics [6], [15]:

- 1) the *jumpstart*—i.e., a higher initial performance when learning a new task based on transfer with respect to a nontransfer condition;
- 2) the *learning speed*—i.e., a faster learning with respect to a nontransfer condition;
- 3) the *asymptotic improvement*—i.e., an improvement in asymptotic performance (i.e., the performance when the learning process achieves its maximum after a sufficiently long training). This is relevant only when a difficult task can be learned only after having previously learned other tasks, otherwise this criterion is not distinct from performance optimality.

TABLE I

EXPECTED PERFORMANCE OF SIM AND EXP, AND DESIRED PERFORMANCE OF TERL, WHEN LEARNING A NEW TASK AFTER HAVING PREVIOUSLY LEARNED A TASK REQUIRING THE *Same* SKILL, A *Similar* SKILL, OR A *Different* SKILL. IN PARTICULAR, SIM STARTS FROM THE PREVIOUS SKILL, SO IT HAS AN ADVANTAGE IF TRANSFER IS POSSIBLE BUT AT THE SAME TIME ALWAYS LOSES THE CAPACITY TO SOLVE THE PREVIOUSLY SOLVED TASK (CATASTROPHIC FORGETTING). EXP LEARNS EACH TASK FROM SCRATCH, SO IT CANNOT HAVE ANY ADVANTAGE OF TRANSFER NOR ANY DISADVANTAGE IN TERMS OF CATASTROPHIC INTERFERENCE. TERL WAS DESIGNED TO PRODUCE THE BEST OUTCOME IN ALL CONDITIONS: THAT IS, TO TAKE ADVANTAGE OF TRANSFER WHENEVER POSSIBLE, TO AVOID WASTING TIME WHEN TRANSFER IS NOT POSSIBLE, AND TO ALWAYS AVOID CATASTROPHIC FORGETTING. *L*: EXPECTED LEARNING SPEED; *C.f.*: EXPECTED EFFECTS OF CATASTROPHIC FORGETTING; *N*: “NEUTRAL” PERFORMANCE, AS IN LEARNING FROM SCRATCH (WITHOUT TRANSFER); *G*: GOOD PERFORMANCE, BENEFITING OF TRANSFER OR HAVING NO CATASTROPHIC FORGETTING; AND *B*: BAD PERFORMANCE, WORSE THAN STARTING FROM SCRATCH, OR PRESENCE OF CATASTROPHIC FORGETTING. A BOLD LETTER INDICATES THE BEST OUTCOME IN THE GIVEN CONDITION

Prev. task	SIM		EXP		TERL	
	L.	C.f.	L.	C.f.	L.	C.f.
Same	G	G	N	G	G	G
Similar	G	B	N	G	G	G
Different	B	B	N	G	N	G

Here we required not only that TERL approaches performance optimality (solution of the tasks), but also that it approaches the performance of SIM when transfer is advantageous and the performance of EXP when starting from scratch is the best thing to do, while at the same time avoiding catastrophic forgetting. A comparison between the expected performance of SIM and EXP and the desired performance of TERL in different conditions is summarized in Table I.

B. Tests With the Planar Arm: Sequential Learning of Tasks

TERL was tested with two setups requiring simulated robotic dynamical arms to reach targets positions in space. These setups were chosen because of the following:

- 1) they involve continuous state and action spaces and the control of dynamical plants;
- 2) they allow a useful visualization and analysis of the performance of the system, in particular the parallel visualization of the postures/movements of different actor experts;
- 3) they involve limb movements, a category of biological phenomena that could be investigated with TERL in future work.

The first set-up involved a 2-D 2 DoF simple simulated dynamic arm working on a plane containing four “object” goals each having a radius of 3 cm (Fig. 2). The planar arm was very useful for developing the algorithm and also facilitates the explanation of its functioning (see examples below). The second setup involves a 3-D, redundant, 4 DoF simulated dynamic robotic arm and allowed the test of the capacity of the system to scale up to more complex tasks (Section III-E).

Fig. 2 shows the 2-D arm used to test the model. The arm was formed by two links: an upper arm measuring 25 cm

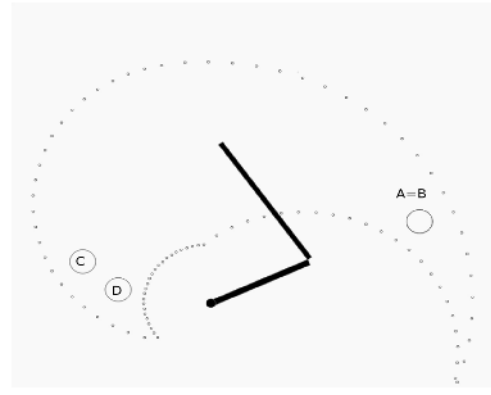


Fig. 2. Dynamic 2-D arm and the four target objects A, B, C, and D (A = B indicates that the location of A coincides with the location of B). Dots represent the borders of the workspace established by the length of the arm links and the range of its joints. The work space is asymmetric as the “elbow” joint range is asymmetric with respect to the upper arm as in a humanoid robot arm.

and a forearm measuring 35 cm. The movement range of the joints was set to $[-100^\circ; +30^\circ]$ for the shoulder (0° corresponding to the upper limb located forward the robot and the angle was measured anticlockwise) and to $[0^\circ; +160^\circ]$ for the elbow (0° corresponding to the straight arm and the angle was measured anticlockwise).

The dynamics of the arm were simulated based on the following equations [57]:

$$\begin{aligned}
 q_s &= (I_s + I_e + 2M_e L_s S_e \cos \theta_e + M_e L_s^2) \ddot{\theta}_s \\
 &\quad + (I_e + M_e L_s S_e \cos \theta_e) \ddot{\theta}_e \\
 &\quad - M_e L_s S_e (2\dot{\theta}_s + \dot{\theta}_e) \dot{\theta}_e \sin \theta_e + B_s \dot{\theta}_s \\
 q_e &= (I_e + M_e L_s S_e \cos \theta_e) \ddot{\theta}_s + I_e \ddot{\theta}_e \\
 &\quad + M_e L_s S_e \dot{\theta}_s^2 \sin \theta_e + B_e \dot{\theta}_e
 \end{aligned} \tag{24}$$

where q_s and q_e are the actuated torques of the shoulder and elbow joints, respectively, and the parameters M , L , S , I , and B are, respectively, the mass, the length, the distance from the center of mass to joint, the rotational inertia of links, and the coefficient of viscosity (these parameters were set to $\{0.9, 0.25, 0.11, 0.065, 0.08\}$ for the shoulder joint and to $\{1.1, 0.35, 0.15, 0.1, 0.08\}$ for the elbow joint as in [57]). The equations were integrated with a fourth order Runge–Kutta method using a time step of 0.01 s.

The arm had two actuated DoF: one for the shoulder joint (θ_s) and one for the elbow joint (θ_e). In particular, the robot controllers used here set the desired posture (or “desired angles,” or “action”) of the arm, and actor experts represented the mappings from sensory input to desired postures. A PD controller was used to generate the torque of each of the two arm joints [58]

$$q = K_p(\theta_{des} - \theta) - K_d \dot{\theta} \tag{25}$$

where θ and θ_{des} are, respectively, the current and desired joint angles, and K_p and K_d are, respectively, the proportional and damping gains (K_p was set to 25 and K_d to 4 for both joints).

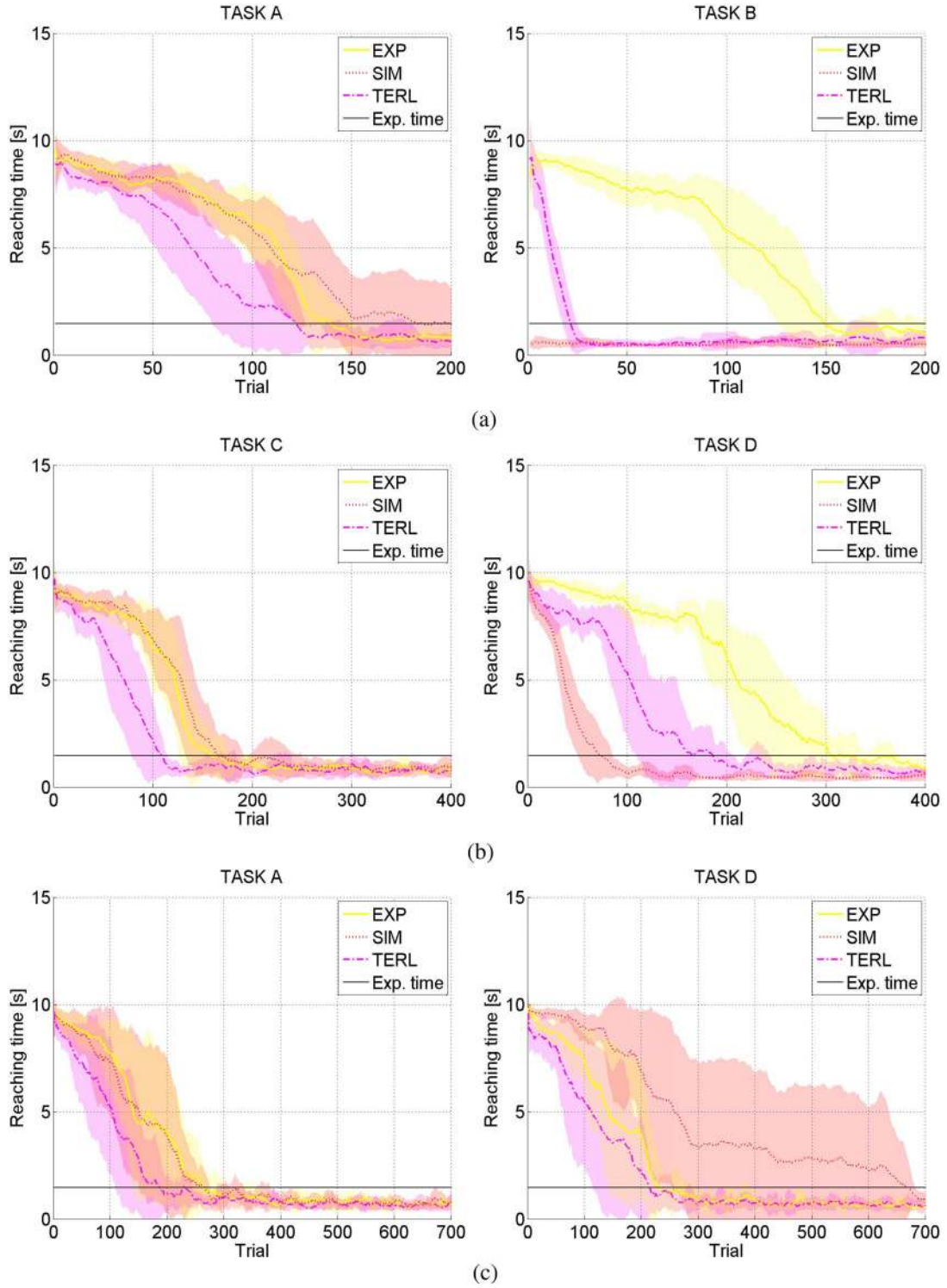


Fig. 3. Performance (reaching time) during learning in three sequential tests in each of which the systems TERL, SIM, and EXP learn a first task for n trials (200; 400; 700 in the three conditions; left panels) and then a second task for further n trials (right panels). Each curve represents the average and standard deviation of ten repetitions of each experiment. (a) Test 1: test with two tasks requiring the same sensorimotor mapping (tasks A and B). (b) Test 2: test with two tasks requiring a similar sensorimotor mapping (tasks C and D). (c) Test 3: test with two tasks requiring a different sensorimotor mapping (task A and D). “Exp. time” represents the “exploitation time” of the noise generator (as explained in detail in Section II-B, this represents an initial period of time at the beginning of each trial during which exploratory noise is kept at a minimum to foster exploitation).

The objects shown in Fig. 2 allow the specification of four reaching tasks, each requiring that the system learns to reach one specific object. We will call these tasks “task A,” “task B,” etc., depending on the target object. The solution of each task

requires the system to acquire a skill that allows the system to reach one of these objects starting from any initial posture. All trials involving the solution of the tasks started with the arm set at a random posture. Each trial terminated either with the

achievement of the goal (target object) or with a time out of 6 s (i.e., 600 steps of 0.01 s each; 0.01 is also the integration time step used for dynamical equations of the arm). The models got a reward of 1 when the hand touched any point of the target object, and 0 otherwise. This implies that the system can work with tasks defined in terms of reaching a set of states, rather than a single state, to the extent that a suitable reward function can be associated with them.

The tasks were used to build three tests corresponding to the three conditions of the rows of Table I. In each, the system had to first learn a task during n trials, and then another task for further n trials (this is what we call here a “sequential test”). The first test involves first learning task A and then task B each for 200 trials: as these tasks require the same sensorimotor mapping this allows us to check if TERL is able to reuse experts used to solve task A when learning to solve task B and how efficient it is in doing so (by comparing TERL with SIM, which performs an instantaneous transfer, this allows us to quantify the “overhead” cost of solving the task with the hierarchical architecture of TERL rather than with a simple expert). We set the number of trials for this and other experiments through pilot experiments that showed the amount of learning allowing the best tested systems for each condition to achieve steady-state performance. The second test involves first learning task C and then task D each for 400 trials: as these tasks require similar sensorimotor mappings, TERL should exhibit a transfer advantage when learning D in comparison to EXP that learns all tasks from scratch. Finally, the third test involves learning first task A and then task D each for 700 trials: as these tasks require very different sensorimotor mappings, TERL should realize that no transfer has to be attempted, and show a performance similar to EXP, that learns from scratch, and superior to SIM, that starts to learn task D from the different sensorimotor mapping learned for task A.

The avoidance of catastrophic forgetting can be ascertained by comparing the performance in the first task (e.g., C) after learning it, with the performance in the same task after the system has learned the second task (e.g., D) and without relearning the first one. The system is robust to forgetting if the performance in solving the first task does not decrease after learning the second task.

Overall, the tests involve these challenges.

- 1) The experts have to learn to associate, at each time step of the trial, a suitable desired posture (output) to the current posture of the robot (input): the dynamics of the arm will then generate the actual trajectory of the arm given the desired EPs selected by the experts; no cost is given for performing movements, but due to the discount of the reward (see Section II-B) the RL algorithm tries to find trajectories of EPs that minimize the time taken by the dynamic arm to touch the target object. This aspect is not further discussed here, but in previous work we have shown how one expert of the type used here can learn to produce quite complex equilibrium-point trajectories to minimize such time and possibly move around obstacles [59], [60].

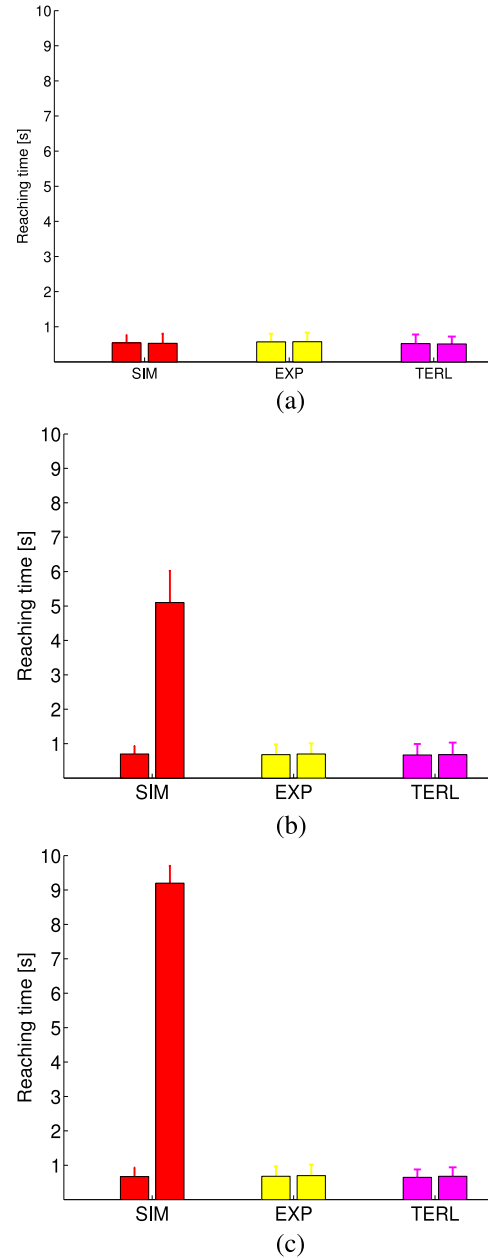


Fig. 4. Test on catastrophic forgetting. Each graph shows the average and standard deviation of reaching time in the first task before (bars on the left) and after (bars on the right) learning the second task. (a) Same sensorimotor mapping (task A then task B). (b) Similar sensorimotor mapping (task C then task D). (c) Different sensorimotor mapping (task A then task D).

- 2) The gating networks are only informed about the identity of the currently solved task (task A, task B, etc.) and on this basis they have to learn to select the best expert(s) to solve it.
- 3) The management of the expert copies by the learning algorithms have to support the exploitation of the opportunities for transfer between tasks while at the same time avoiding catastrophic interference.

The initial connection weights of the expert actors were randomly drawn from a uniform distribution ranging in $[-0.2, 0.2]$, and the connection weights of the critic experts were set to zero for SIM, EXP, and TERL. The connection

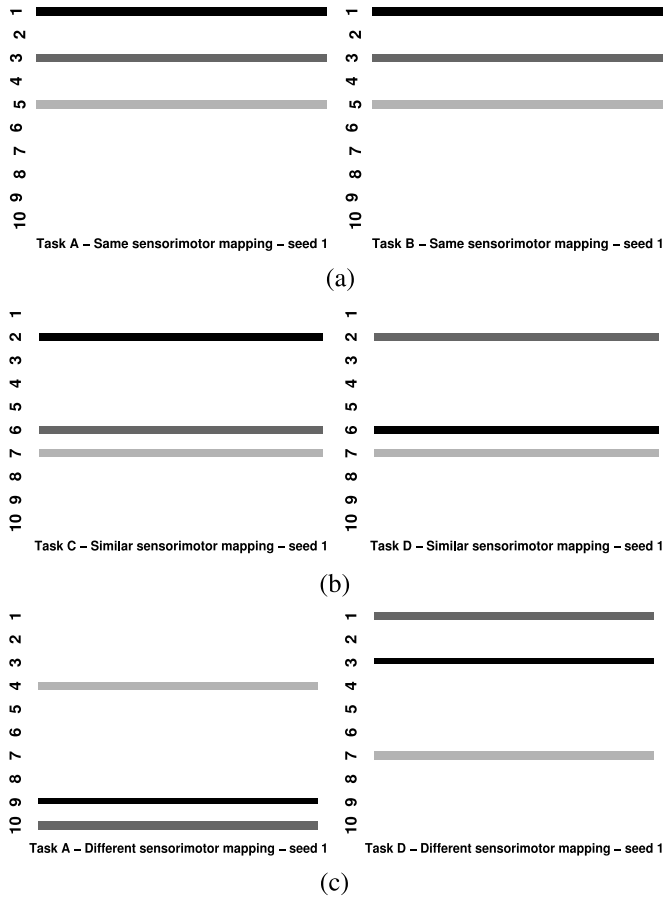


Fig. 5. Use of the actor experts by TERL in three tests, each involving the solution of two tasks in sequence, characterized by (a) same sensorimotor mapping, (b) similar sensorimotor mapping, and (c) different sensorimotor mapping. Each of the three graphs shows the three highest prior responsibility signals of the ten experts in one trial after learning the first task (left panels) and in one trial after learning the second task (right panels). In each graph the highest, second highest, and third highest priors are, respectively, marked with a black, dark gray, and light gray strip, while all other priors are not marked (white). The data refer to one specific system (seed 1); in other replications of the experiment, or if the analysis is repeated for the critic experts, the results are the same but of course involve different experts. Note that this type of graphs give information on the selection of experts during the whole trial: here the graphs simply show that TERL reliably selected specific experts during the whole trial; in other conditions the graphs are very useful to analyze expert-based models similar to TERL, especially at the beginning of the learning process or when using gating networks having highly varying input.

weights of TERL gating networks were set to zero. The actor and the critic of TERL were each formed by ten experts each. This number was chosen to give the system enough resources to form main experts (i.e., experts with the highest functioning responsibilities) and copies for the three tasks and to study the selective recruitment of redundant resources by the system's gating networks. Other parameter values of TERL are summarized in the Appendix.

Fig. 3 illustrates the learning curves of SIM, EXP, and TERL in the three tests. In particular, for each test the figure reports the learning curves related to the first and second task forming each test.

The results of the first test, formed by two tasks involving the same sensorimotor mapping, show that TERL manages to quickly discover that in order to solve the new (second) task B

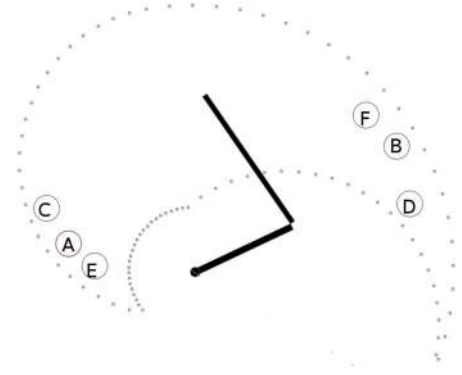


Fig. 6. Configuration of targets (tasks) used to test the scaling-up capacities of TERL.

it can reuse the skill previously acquired by solving task A. The overhead cost for this discovery is rather low (compare TERL performance with SIM performance in task B). TERL has also some advantage on the other models in the task A solved as first: the reason for this is that TERL experts are initialized to incorporate different initial behaviors from which the selection mechanism can draw one that is most appropriate for the current task (this is further explained on the basis of some figures in Section III-C).

The results of the second crucial test, formed by two tasks involving similar sensorimotor mappings, show that TERL is capable of discovering that in order to solve the second task (task D) it can start from the skill previously acquired by solving task C and this produces a notable advantage in its learning speed. The difference in performance between EXP (starting from scratch) and TERL in task D shows the transfer advantage for TERL. The difference in performance between TERL and SIM (transferring immediately) in task D shows the overhead cost that TERL pays to understand the possibility of exploiting transfer.

Finally, the results of the third test, formed by two tasks involving very different sensorimotor mappings, show that TERL manages to avoid using the skill previously learned for task A to solve the new task D as the two are quite different. The difference in performance between SIM and TERL in task D shows the cost that TERL avoids by not attempting to transfer from the different previously-learned task A.

While having these transfer strengths, TERL also manages to avoid catastrophic forgetting. Fig. 4 shows the comparison between: 1) the performance of the three models in the first task (in each of the three tests) right after such task has been learned; and 2) the performance right after the model has learned the second task (with the learning rates set to zero, so as to avoid relearning). The comparison of the performance in the first task, before and after learning the second task, allows the evaluation of how the learning of the second task interferes with the performance of the first one.

As expected, SIM suffers catastrophic forgetting both when the second task is similar but not identical to the first one and when the two tasks are very different. Again as expected, EXP performs very well since it learns every task from scratch, and

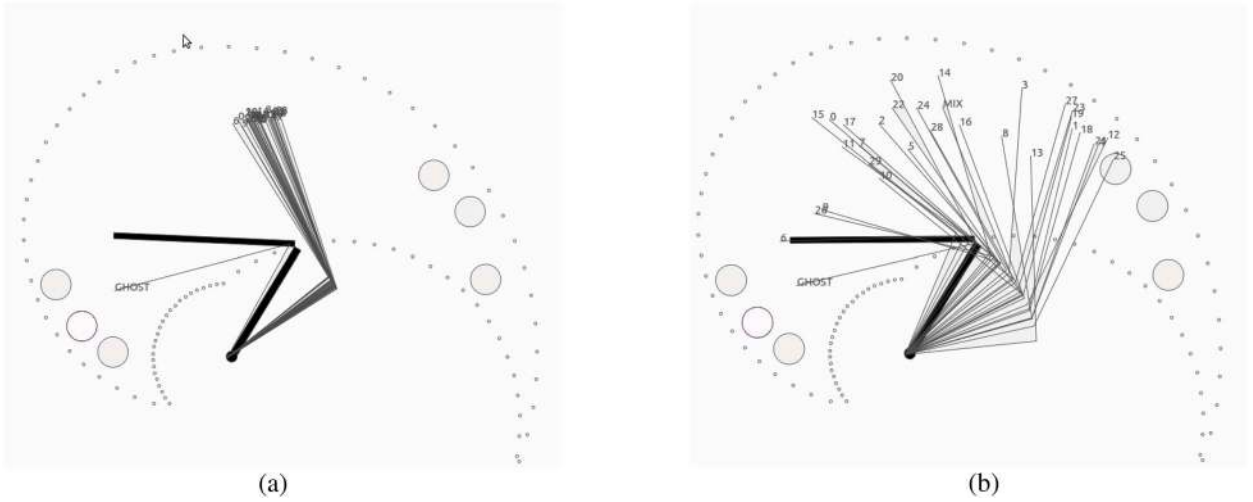


Fig. 7. Snapshot of the postures (actions) suggested by the actor experts in the (a) low scattering condition (TERL) and (b) high scattering condition (TERLs). The thick stylized 2-D arm indicates the current arm posture whereas the other thin stylized 2-D arms represent the arm postures suggested by the 30 actor experts (the number on each arm drawing indicates the expert). The thin arm marked with “MIX” represents the overall arm posture suggested by the mixture of experts; the thin arm marked with “GHOST” indicates the arm posture suggested by the exploration noise added to the mixture posture.

previous acquired skills are not affected by the learning of the new ones. What is most important is that, as expected TERL does not suffer of the problem of catastrophic forgetting in any of the three tests. The reason is that when a new task is the same or similar to the previously learned task, TERL uses a copy of the previously acquired skill to solve the new task thus not damaging the capacity of the best expert trained to solve the previous task. When the new task is different, TERL does not transfer but rather uses a completely new expert, thus avoiding interference.

To show the actual formation and use of background copies by TERL, Fig. 5 shows the responsibilities that it assigns to actor experts during learning in the three tests involving same Fig. 5(a), similar Fig. 5(b) or different Fig. 5(c) sensorimotor mappings (left panels show responsibilities of actor experts for the first task of each test; right panels for the second one; data are qualitatively similar for the responsibility signals of the critic experts, data not shown). The graphs show that, when the two tasks require the same sensorimotor mapping [Fig. 5(a)], TERL uses the same experts with the same responsibilities, i.e., it reuses the skill developed for task A to solve task B. When the two tasks require similar sensorimotor mappings [Fig. 5(b)], TERL solves the second task, e.g., task D, by reusing a copy expert developed in background while solving the previous task, e.g., task C. Thus TERL is able to exploit previously acquired knowledge to solve new tasks, and at the same time to avoid catastrophic interference, by recruiting and modifying expert copies developed for previously solved tasks. Finally, when the two tasks require very different sensorimotor mappings [Fig. 5(c)], TERL uses different experts as it realizes that trying to transfer knowledge would be useless or even detrimental.

C. Tests With the Planar Arm: Scaling-Up to Several Tasks and Experts

The experiments reported in this section test the capacity of TERL to scale up to learn a larger number of tasks with

respect to the previous sections, and in particular show how the mechanism of recruitment of already trained experts, or of still untrained experts to solve new tasks, continues to work with a large number of available experts. The experiments also show the capacity of TERL to recruit a restricted subset of experts among the multiple available ones, for both functioning and learning, based on the responsibility signals. For this purpose, we trained the system in sequence with the six tasks shown in Fig. 6, each for 500 trials, using 30 experts for the critic and 30 for the actor (we set this large number to test the scaling capabilities of the system; the number of experts of the critic and actor do not need to be the same).

Before illustrating the results of these experiments, we also indicate how they were also used to quantify the advantage of TERL over other systems due to the possibility of differentially initializing its experts. For this purpose, in these experiments TERL was also tested in a version where the connection weight of the bias of each actor expert network was randomly generated in $[-1, +1]$ with the effect that the initial actions (here desired arm postures, or “EPs,” pursued by the robot PD or PID) of all actor experts were more uniformly distributed in the work space (as usual, all other connection weights were set to zero). In the graphs we will refer to this condition with “TERLs” (“s” refers to the “high initial scattering”) to distinguish it from the condition with a low initial scattering used so far where the bias connection weight was drawn in $[-0.1, +0.1]$ to have some differentiation. Fig. 7 shows the postures of the actor experts in the low and high scattering condition in the initial phase of training. Having initially scattered actor experts is advantageous as it produces different initial postures/behaviors and this facilitates the selection of them by the information-accumulation mechanism of the actor gating network. Although low, the initial (small) scattering of the actor experts used in TERL also explains the advantage it has with respect to SIM and EXP in learning the first task in the tests shown in Fig. 3. Notice that the possibility of exploiting the best skill from a rich initial repertoire of

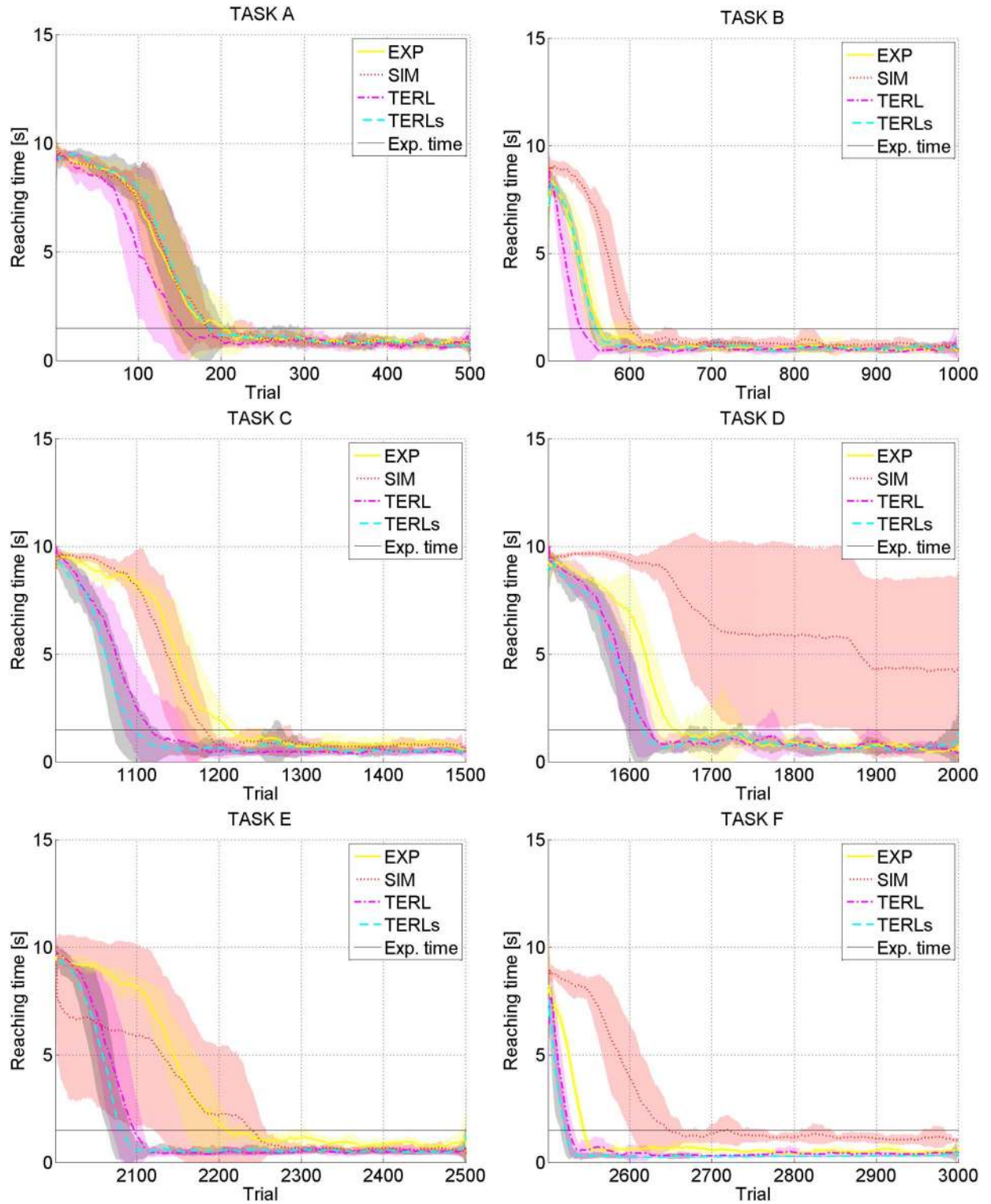


Fig. 8. Reaching time of SIM, EXP, TERL, and TERLs during learning of six tasks in sequence. The test involved learning one after the other the targets A, B, C, D, E, and F shown in Fig. 6, each for 500 trials. Each curve and shadow represent the average and standard deviation computed over ten different repetitions of the simulation.

skills is general as TERL samples in parallel the goodness of all available experts for the current task (since it compares the action of all experts with the action actually performed) and so it rapidly focuses the selection on the best available one. Without such a mechanism (e.g., as is in SIM and EXP) the

initial repertoire of experts, and in general any repertoire of experts available at a later time, could not be probed to isolate the best expert usable to solve the current task. Note that TERL rather than TERLs was used throughout this paper to avoid that the additional advantage given by scattering confounded the

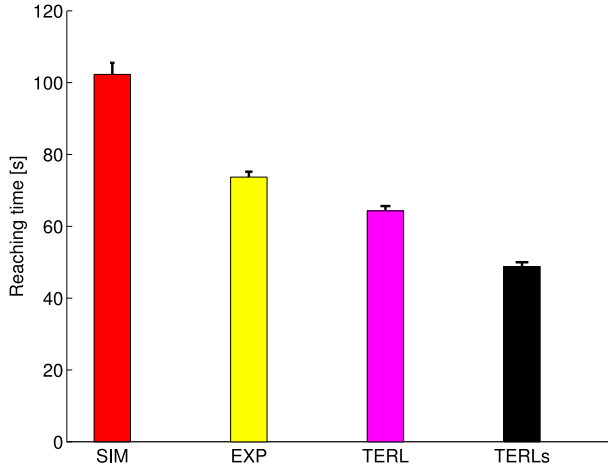


Fig. 9. Overall performance of SIM, EXP, TERL, and TERLs during learning of six tasks in sequence: sums of the integrals, one for each model, of the mean curves of the six graphs plotted in Fig. 8.

advantage of the system based on the mechanisms identifying the best experts from which to transfer knowledge (below we present some other results on TERLs alongside TERL to show its advantages).

Figs. 8 and 9 show the results of the test on the six tasks learned in sequence. TERL performs better than both EXP and SIM architectures. This confirms the capacity of TERL to scale up to a higher number of experts, 30 in this case, and to exploit transfer opportunities when this is possible (e.g., when learning tasks C and E after task A, and tasks D and F after task B). When the actor experts start with a higher scattering of the initial EPs (TERLs), the system has an even higher performance due to the advantages mentioned above.

A larger number of experts involves the following effects on computational costs.

- 1) The functioning and learning responsibility signals of experts are computed in parallel on the basis of the action/value errors of the experts: this implies that the number of tests to do with a task in order to evaluate the experts does not increase with the number of experts.
- 2) The computation of the activation of experts increases linearly with the number of experts (but notice that the activation of different experts is independent, so it could be implemented in parallel hardware).
- 3) The learning processes involve only a given subset of experts for each step (here eight in total) so it does not depend at all on the total number of experts (recall that thanks to the functioning/learning decoupling, the number of experts learning at each trial is fixed and depends on how many background copy experts one wants to train, not on the total number of available experts). Importantly, this also means that the computation time needed by the system functioning and learning is fully independent of the experience already acquired, contrary to most TRL systems that become slower with the accumulation of experience (see Section IV).

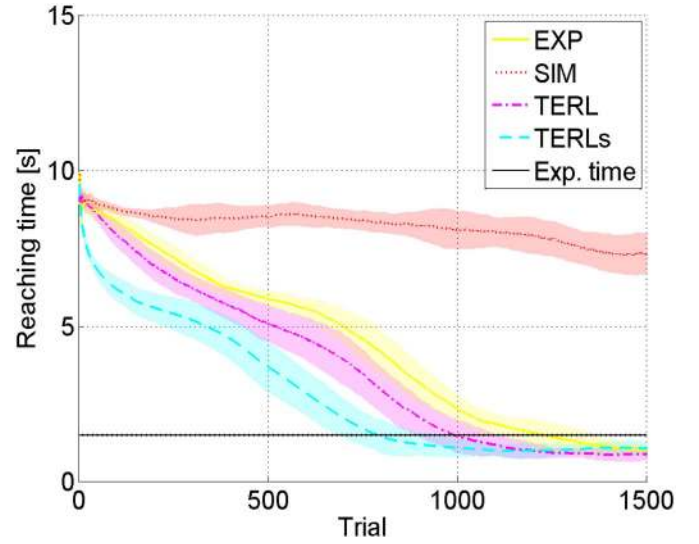


Fig. 10. Performance of the different systems during learning (trials) in the interleaved test where six tasks changed at each trial and for several times. The performance was measured as the time taken by the systems to reach and touch the target during a trial. The six tasks are those of Fig. 6. Curves are averages of ten repetitions of the simulations.

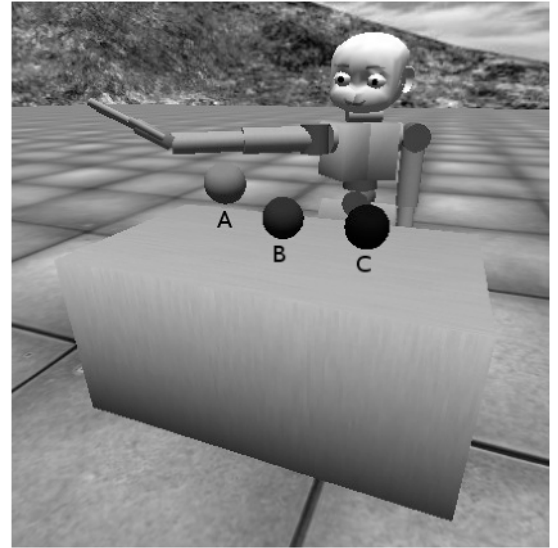


Fig. 11. iCub robot and the environment used to test TERL. The picture refers to the beginning of a trial. Object A represents an obstacle, whereas objects B and C represent the target objects that the robot has to reach.

D. Tests With the Planar Arm—Interleaved Tasks

We also tested the capacity of TERL and TERLs to solve multiple tasks when these are learned in an interleaved fashion rather than in a sequence of blocks each focused on a task, as done in the previous tests. For this purpose, the six tasks considered in the previous section were learned during 1500 trials where at each trial the task to be learned was randomly selected. As before, TERL and TERLs used 30 experts for both actor and critic.

Fig. 10 shows the results of the test. The first interesting point is the advantage of TERL on EXP, which learns from scratch. Although not large, this advantage is important as it shows that when TERL learns similar tasks at the same time

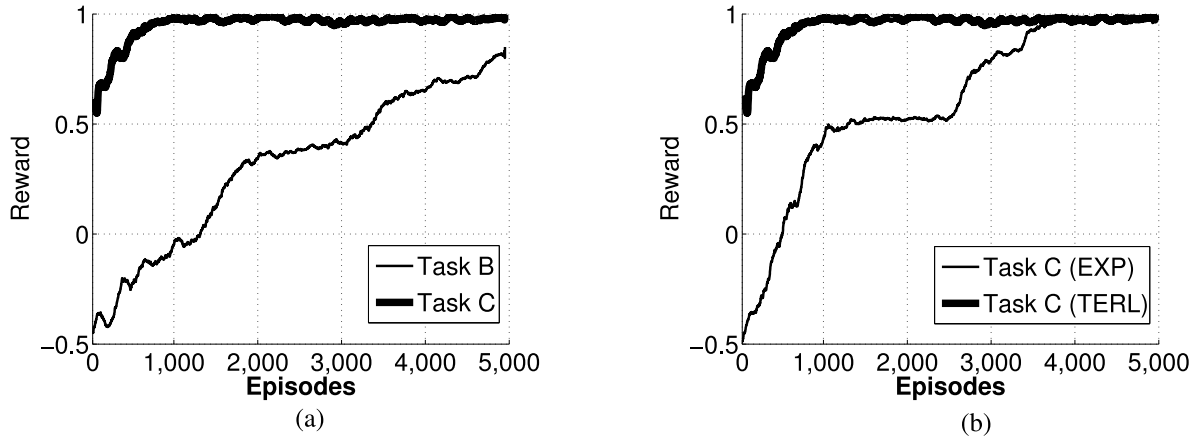


Fig. 12. Performance of TERL and EXP where the systems are first trained on task B for 5000 trials and then on task C for another 5000 trials. (a) Reward acquired by TERL during the sequential learning of the two 3-D tasks requiring a similar sensorimotor mapping (tasks B and C): although task B and task C are learned sequentially, the learning curves are represented in the same graph to ease comparison. (b) Reward acquired by TERL and EXP during learning of task C and showing the advantage of TERL over EXP due to its capacity for transfer. Each curve represents an average over ten repetitions of the experiment, and data are smoothed over 100 trials. Notice that TERL solving task B (a) and EXP solving task C (b) have an initial negative reward indicating that they often hit the obstacle.

with interleaved trials the information gained is shared among the experts. For example, when the system has to learn to reach tasks A, C, and E starting from EPs initially set at the center of the work space, the experts that it forms can initially be the same, so that up to a certain point modifying them for solving one task also improves them to accomplish similar tasks. A second relevant point is the large advantage of TERLs versus all other models, confirming that the initial differentiation of the experts can significantly boost learning speed. Last, for all conditions a higher performance also tends to produce a lower statistical variation of the results of different repetitions of the tests.

E. Test With 3-D Redundant Simulated 4-DOF Arm (iCub)

This section illustrates a test to evaluate whether TERL was capable of exploiting between-skill transfer while avoiding catastrophic forgetting in more complex conditions. The test required in particular to control the 3-D 4 DoF simulated robotic dynamic arm of the humanoid robotic platform *iCub*, an open-source robot built for studying cognitive development in humans [61].

Fig. 11 shows the simulated setup used in this test formed by the *iCub* robot and a 3-D environment containing three spherical target objects (in the simulator, these three objects can be anchored to the world without an actual physical support). Each arm of the *iCub* has 16 joints: three for the shoulder (J_{0-2}), one for the elbow (J_3), three for the wrist (J_{4-6}), and nine for the hand (J_{7-15}).¹ Here we used TERL to control the movements of four joints of the right arm, in particular: 1) J_0 , the “shoulder pitch,” responsible for the front-back movement when the arm is aligned with gravity; 2) J_1 , the “shoulder roll,” affecting the adduction-abduction movement of the arm; 3) J_2 , the “shoulder yaw,” affecting the yaw movement when the arm principal axis is aligned with gravity; and 4) J_3 , the joint related to the elbow. During the simulation J_0 could assume

values in the range $[-80^\circ; -15^\circ]$, J_1 in the range $[10^\circ; 110^\circ]$, J_2 in the range $[-10^\circ; 75^\circ]$, and J_3 in the range $[20^\circ; 85^\circ]$. The positions of the remaining joints were set at fixed values ($J_4 = -10^\circ$; $J_5 = -30^\circ$; $J_9 = 80^\circ$; $J_{6-8} = J_{10-15} = 0^\circ$). The torso joint affecting the yaw with respect to the vertical axis was fixed to -30° . All trials started with the arm set at a fixed posture so that the obstacle was in the way of the targets at each trial: $J_0 = -90^\circ$, $J_1 = 100^\circ$, $J_2 = 90^\circ$, and $J_3 = 6^\circ$.

The three spherical objects in the environment had a diameter of 3.5 cm and were set in front of the robot (Fig. 11). The objects allowed the implementation of two reaching tasks, each requiring that TERL learn how to control the right arm of the *iCub* in order to reach either object B (“task B”) or object C (“task C”) while avoiding hitting the “obstacle” A.

The system had the same architecture and functioning of the system of the previous sections, and was trained as in the sequential test involving the 2-D arm with a few differences. First, all trials involving the solution of the tasks started with the arm set at a fixed posture at the right of object A. This made the task more difficult as object A was always an obstacle for reaching the two targets B and C (Fig. 11) and so the robot had to perform curved trajectories around the obstacle to reach the targets. Second, each trial ended when the *iCub* hit any one of the three objects with the hand, or after a timeout of 8 s (in the 2-D tests the time out was 6 s; as before, the integration time-step of the model and robot equations was 0.01 s). The longer trial duration allowed a longer exploration necessary as the model had to discover a more complex trajectory to reach the targets while controlling a redundant arm (redundancy required the algorithm to autonomously converge to one possible solution among all those explored). Third, when the *iCub* hand hit the obstacle the model got a negative reward signal set to -0.5 , whereas if it touched the target object received a reward equal to 1. The reward was 0 otherwise. Fourth, TERL was endowed with 10 experts for both actor and critic. Last, the values of a few parameters of the model were changed to take into account the different setup, in

¹http://wiki.icub.org/wiki/ICub_joints

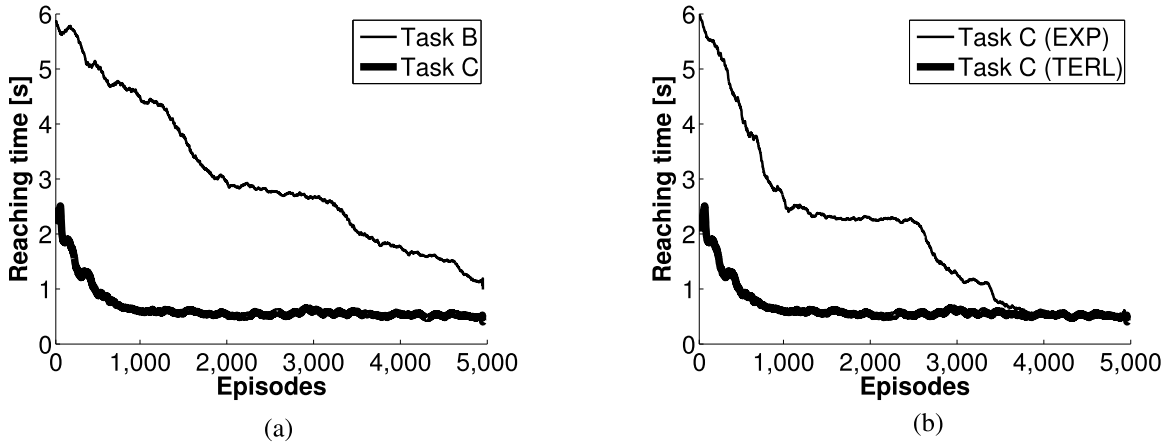


Fig. 13. (a) Reaching time of TERL during the sequential learning of the two 3-D tasks requiring a similar sensorimotor mapping (tasks B and C). TERL first learns task B for 5000 trials and then learns the second task C for other 5000 trials. (b) Reaching time of TERL and EXP during learning of task C. Each curve represents an average over ten repetitions of the experiment, and data are smoothed over 100 trials with a moving average.

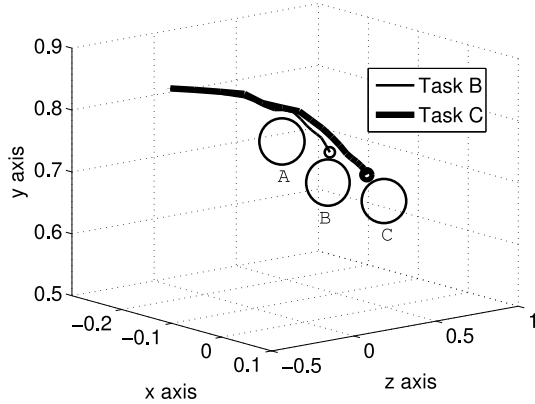


Fig. 14. Trajectories of the iCub hand produced by TERL after learning two 3-D tasks requiring a similar sensorimotor mapping (tasks B and C while avoiding A, see Fig. 11). The trajectories refer to one run of the simulation but other runs produced qualitatively similar results.

particular the higher number of controlled degrees of freedom that required a different number of output units, and thus a different number of Gaussian functions used by the algorithm as features (see Table II in the Appendix).

In the test, the system first learned task B and then task C, each for 5000 trials. As the tasks require a similar sensorimotor mapping, they allowed us to evaluate whether TERL was capable of reusing the experts employed to solve task B when learning to solve task C. To show the capacity of TERL to exploit transfer we also tested EXP in the same experiment involving first learning of task B and then of task C.

Fig. 12(a) shows that TERL learns the second task (task C) much faster than the first task (task B) as it can partially transfer the ability acquired for the first task to the second one. The fact that the advantage in task C is actually due to TERL’s transfer capacity rather than a different difficulty of the task C with respect to task B is demonstrated by the fact that TERL is much faster in learning task C than EXP [Fig. 12(b)]. The same result is shown by the reaching time instead of the reward (Fig. 13).

The curve trajectories performed by TERL to solve the obstacle avoidance tasks B and C are shown in Fig. 14. These curves show how the sensorimotor mappings needed to solve the two tasks are similar which explains the higher performance of TERL capable of taking advantage of this.

IV. DISCUSSION: TERL COMPARISON WITH OTHER MODELS

This section compares TERL with other models illustrating its similarities and novelties with respect to them and the different subproblems of TRL they face.

A. Option Framework and Off-Policy Learning

Several systems from the RL literature perform transfer of knowledge between tasks by relying on the option framework (see [62] for a review). An option is a data structure that encapsulates a policy (a skill), an initiation set (states where the option can be selected), and termination condition (establishing when the option execution terminates, e.g., when some subgoal states are achieved). Options can be used as building blocks to solve different complex tasks through their suitable combination [63]. Option-based systems have been mainly used to face compositionality problems rather than skill-to-skill knowledge-transfer problems as here. These are two very different problems. Compositionality problems concern how assemble multiple skills to solve complex tasks that require multiple options (skills) to be solved. For example, the option framework allows the formation of policies that select multiple options in sequence, together with primitive actions, to accomplish complex goals. Compositionality transfers knowledge from a (complex) task to a new (complex) task by reusing whole skills (or part of them), rather than from acquired skills to a newly learned skill as required in skill-to-skill transfer problems. For example, the complex skill of pressing a button and the complex skill of “scratching with the finger” might share a very similar component-skill of “reaching with the finger” performed before performing, respectively, the component-skills “pushing with the finger” or “scratching.”

A pioneering model of this type is *compositional Q-learning model* [17], relevant here also because it proposed the idea of adapting some principles of the mixture-of-experts model to a RL context. As another example, it has been shown [64] that the transfer of whole options between different grid-world environments can be facilitated when options are based on the agent's sensations (e.g., egocentric position of objects) rather than on information related to the specific problem (e.g., absolute spatial relations between objects).

A key idea that can be exploited within the option framework [62], and more generally in systems formed by multiple RL experts, is "off-policy learning" [63]. Off-policy learning is employed by several TRL systems reviewed below and allows RL algorithms, for example Q-learning [24], [65], to train a "target policy" (and/or a value function) on the basis of actions selected by another "behavior policy." Off-policy learning has been used for "intraoption learning" to let options learn from their partial execution or from the execution of actions by other components of the system [62]. A system that makes massive use of off-policy learning is the Horde architecture [66]. This architecture is also related to TERL as it is formed by multiple individually simple experts ("demons") that can be used for prediction or control purposes. Each demon learns an approximation (e.g., based on features and linear functions) of a *generalized value function* (GVF). A GVF stores knowledge and learns to predict a relevant element through mechanisms analogous to those used to estimate RL conventional value functions, but it can refer to elements different from rewards such as the activation of a robot sensor. A demon is based on a reward function, a termination function, a terminal-reward function, a policy (which is given in the case of prediction) and learns the related GVF. Such functions can be structured so that, for example, a demon can be used to learn to predict the time when a particular sensor will be activated, or to learn a policy to accomplish a certain activation of a sensor through an off-policy RL algorithm. TERL shares with Horde the idea of behavior or value prediction based on multiple experts and their parallel learning, and the background objective of life-long learning, but differs from it in many respects, in particular it pivots on transfer learning to improve learning speed whereas Horde relies on off-policy learning.

B. Models From the Literature on Transfer Reinforcement Learning

Among the several systems proposed within the literature on TRL, we focus here on those closer to TERL, i.e., those that address the source-task selection problem by explicitly reasoning about "libraries" of already solved tasks to decide from which of them to transfer knowledge to the new task [6]. One of these systems, called *policy reuse Q-learning* (PRQ-learning) [67], uses off-policy learning, based on Q-learning [65], to transfer knowledge. In particular, it uses "source policies" (i.e., previously learned policies) from which to possibly transfer knowledge to train (offline) a new optimal policy for each new task: the source policy to use for transfer is selected at each trial on the basis of a soft-max function applied to the "reuse gains" of all source

policies (the reuse gain of a policy is equal to the average reward obtained in the new task). The newly learned policy is added to the policy library if its reuse gain is larger than a certain threshold with respect to those of existing policies. With respect to TERL, PRQ-learning has the advantage of computing an explicit metric of the similarity between tasks and the capacity to identify a core set of policies for a domain. However, at each trial it requires the selection of only one policy to use (the system is tested in a grid world) and evaluates the reuse gain of only that policy: this implies the need for an increasing number of trials to evaluate the source policies as their number increases. Instead, TERL evaluates the ability of experts to solve/evaluate the new task in parallel at each step. In addition, regarding brain modeling, PRQ-learning solves the interference problem by building a new data structure for each new task and so it could not be usable to model how the brain faces such problem.

Another class of TRL systems transfers knowledge encoded as "experience samples," in particular tuples $\langle s, a, s', r \rangle$ (i.e., the state and action at a certain time step and the state and reward at the following time step), rather than "compiling" it into the parameters of a function approximator (i.e., it follows a non-parametric approach). For example, the system proposed in [68] stores the experience samples during a "sampling phase" and then, during an offline learning phase, computes the approximation of the action-value function on the basis of the stored samples. The key idea of the system is to implement knowledge transfer by using samples from source tasks that are similar to the target task. The selection of the source tasks and of their tuples is done on the basis of the "compliance" (similarity) of such tuples with those experienced in the target task. The system has the advantages of basing transfer on the similarity between the dynamics and reward functions of different tasks and to be applicable to problems where both reward and transition functions change. However, its batch nature and the computation of compliance requires it to explicitly record experienced tuples, which is memory intensive, and to directly compare tuples, which implies a computational cost that increases with the number of experienced tuples. Moreover, as these mechanisms directly store tuples of experience they could not be used for modeling "compiled," semantic knowledge as done in biological neural networks. TERL, instead, tests the experts for selection in parallel and "compiles" the acquired information into the parameters of the experts (function approximators), so it requires computation time and memory resources that do not increase with the length of past experience.

C. MOSAIC Models

The systems most similar to TERL are those related to the MOSAIC model [69], so these are considered more in depth. TERL and some MOSAIC models have important similarities but also differences that make them suitable to face different problems. A first MOSAIC model similar to TERL is the *multiple model-based RL model* (MMRL) [70]. MMRL is composed of modules formed by a predictor of the world dynamics (taking as input the current state and planned action,

and returning as output the next state) and a RL controller. During each trial, the errors of the predictors determine the responsibility signals of the modules: for this purpose, each module predictor error is passed through a Gaussian function (measuring its “correctness”) and then all the Gaussian activations are passed through a soft-max function. The responsibility signals are used to both weight-average the contribution of the RL controllers to act, and to regulate the learning rate of both predictors and related RL controllers. In this way, different modules specialize in predicting and acting in subportions of the whole problem space based on the correctness of the predictors.

These ideas are further developed in the MOSAIC-RL model [71], the MOSAIC-family model most similar to TERL. MOSAIC-RL is formed by three sets of “modules”: “forward modules,” predicting the environment dynamics (i.e., models of the state transition function), “reward modules,” predicting the one-step reward (i.e., models of the reward function), and RL controllers (each based on a value-function approximator and a policy directly generated from the value approximator through control theory methods). With respect to previous MOSAIC models, in MOSAIC-RL the three types of modules are decoupled, so they segment the problem space differently on the basis of their respective errors: 1) the world-dynamics prediction errors; 2) the one-step reward prediction errors; and 3) the TD-errors. As in MMRL, these errors are used to compute the responsibility signals of the related type of modules on the basis of Gaussian and soft-max functions. However, now the three types of modules have distinct responsibility signals, so the RL controllers deciding the action are selected on the basis of their TD error, not on the basis of the world dynamics as in MMRL. This allows the model to face not only problems with nonstationary (i.e., hidden) dynamics, as MMRL, but also problems involving “a non-stationary reward function,” i.e., multiple reward functions (multiple tasks).

Given the similarities between TERL and MOSAIC-RL, in terms of the responsibility signals computed in parallel on the basis of a softmax of Gaussians of errors, and the segmentation of the whole problem into subproblems (“tasks”), it is important to clarify how they can be used to face different problems. The key difference between the two systems is that TERL uses the softmax of Gaussians to train the gating networks, not only to select the experts to act: the knowledge acquired with this learning process during the first trials in which a new task is faced (when TERL selects multiple experts for action as it still does not know which one is the best) is compiled into the parameters of the gating networks. Once acquired, this knowledge can be used to immediately select the best expert to use, based on the information regarding the identity of the task to solve, since the very first step of the trial. Instead, MOSAIC-RL uses the softmax of Gaussians to dynamically accumulate evidence on the experts to select step-by-step during the trial, or during multiple trials when these address the same task multiple times in sequence, even after the system has learned to solve the task. This implies an important feature of the problems that MOSAIC-RL can solve. In trial-based RL problems, relevant for TRL, reward

is often zero during the trial and high at the end of the trial when the task “goal” is accomplished, so the TD-error can be high only at the end of the trial. In these cases, the responsibility signals of controllers computed by MOSAIC-RL on the basis of their TD-errors are similar during the first part of the trial and differentiate only at the end of it, so they can start to have an effect on the selection of controllers only from the second trial onward, and only when the same task is experienced for more than one trial in sequence. In this respect, in commenting on the performance of MOSAIC-RL [71], it was reported that “after learning, the RL modules also successfully switched within a few (one or two) trials when the subenvironment changed” (note that the training and test of MOSAIC-RL was done in blocks each formed by 100 trials involving the same “subenvironment,” i.e., task).

This difference implies that MOSAIC-RL is not suitable to face problems as those used here to test TERL where the task is switched at each trial and its identity is known (interleaved condition tests, see Fig. 10). This situation is for example common in animals where motivation changes continuously after being “satisfied in one trial” and is known to the animal. After learning, for example, when an animal is hungry (hunger signals a first task identity) it is able to directly move to a food dispenser, and when it is thirsty (thirst signals a second task identity) it is able to directly move to a water dispenser: this without the need to sample each time the reward given by food or water. The same holds when animals pursue a goal: also in this case the task identity is known and so the animal can immediately recall the behavior to accomplish it (after this behavior has been acquired and associated to the goal). Note that this behavior of MOSAIC-RL is expected and *is not a drawback* of the system. Indeed, MOSAIC-RL has been designed to solve problems where the task identity is hidden and so has to be identified by repeated sampling: the mechanisms of the system are thus very good for facing these problems. Instead, for the TRL problems considered here, for which TERL has been designed, the information about the identity of the task to solve is clear and available before each trial. Since MOSAIC-RL addresses a different problem with respect to TERL, it does not use different responsibility signals for the evaluation function and the policy “experts,” nor does it use different responsibility signals for functioning and for learning as done by TERL. These features are very important, as shown here, when the system is used to address the skill-to-skill knowledge-transfer problem and the catastrophic-interference problem.

Overall, it can be said that MOSAIC models and TERL are best suited to solve different complementary problems: MOSAIC models are best suited to face the problem of segmentation of whole MDP problems into subtasks on the basis of the progressive online accumulation of evidence on the hidden (nonobservable) world features based on the errors of predictors of the world dynamics or the world reward. Instead, TERL can be used to solve the skill-to-skill transfer problem without incurring catastrophic interference on the basis of gating networks learning to map goals to responsibilities of experts. For these reasons, the two systems might be suitably integrated in the future.

D. Other Algorithms, Evolutionary Duplication, and Task-Policy Mapping

The generation of multiple neural copies by TERL links it to another class of models that implement neural duplication to solve different tasks or subparts of the same task. One of these models [42] is based on RL experts called “mots” that self-organize in a 2-D spatial grid in a way that is reminiscent of self-organizing neural networks. Mots are selected for action and for learning based on a previously proposed model called *selected expert reinforcement learner* [72]. In particular, each mot implements a RL policy—e.g., in the form of Q -learning—and at the same time learns to estimate the inaccuracy of that policy—e.g., in the form of error of the Q -learning values. At each time step, a “winning mot” is selected for action and for learning if it has the highest Q -value reduced by the Q -value estimated error. The mots close to the winning mot also learn and this leads to the emergence of spatially organized mots that tend to specialize on the same parts of the problem space, and thus to form “copies” of the same portion of behavior similar to TERL. Later [43], mots maps have also been developed to capture in space the temporal relations (sequential activation) of mots by training not only the winning mot and its neighbors but also the previous and/or following winner(s) and neighbors. Future work might aim to obtain a spatial organization of TERL experts by employing similar mechanisms as those used by mots, thus benefiting from the advantages that this carries in terms of smoothness, robustness, hierarchy by region, dimensionality reduction, and reuse (see [43]).

Another relevant class of models that implement neural duplication has been proposed in [44] and [73]. These models are developed within an evolutionary framework, so they can be used to model neural duplication based on genetic processes [73] whereas their use to model learning processes related to the single individuals, as in TERL, has been proposed only recently [74], [75]. In this respect, TERL represents a relevant hypothesis as to how neural duplication might take place within the single animal behavior and brain when trial-and-error learning mechanisms are involved [22].

Another class of systems learns to map tasks to solutions, in particular it uses the information regarding features of the goals, or reward functions, to best initialize the value functions or policies to solve new tasks [76], [77]. As in TERL, in these models the transition function is assumed to be constant. While solving different tasks, these systems learn a mapping between the goal features and the parameters of the value function or policy approximators. This mapping allows the systems to generate the initial policies and value functions to solve new tasks: later these policies and value functions can be refined with further training experience. These systems face a TRL problem that is different, in particular complementary, with respect to the one faced by TERL. Indeed, they can exploit information (features) on the goals of the new tasks *before* solving them, but then when they further refine the initial solutions they cannot further transfer knowledge from previously solved tasks. On the contrary, TERL cannot benefit from information (features) on the new tasks before starting to solve them, but on the other hand it can transfer information

from previously solved tasks *after* it starts to solve the new tasks based on the effectiveness of existing experts in solving them (rather than based on similarities between goals). To further clarify this difference, consider these examples. When the system faces a new task where the goal is very similar to the goal of a previously solved task, and indeed the two tasks require a very similar sensorimotor mapping to be solved, then the systems under discussion can formulate a good first guess on the solution to use. If instead the new task requires a very different sensorimotor mapping, and such sensorimotor mapping (or a similar one) was previously learned to accomplish a different goal, then the mechanisms of TERL would rapidly identify the expert encoding the latter mapping and use it to solve the new task. Similarly, the mechanisms of TERL would be useful when solving a new task where the goal does not give any useful information about the type of sensorimotor mapping to use. For these reasons, the mechanisms proposed by the systems discussed here and TERL might be integrated in future work (Section V sketches how TERL might do this).

V. CONCLUSION

A. Main Achievements

This paper has described and tested a RL architecture (TERL architecture) to learn multiple skills solving different related tasks. The architecture offers a solution to the RL source-task selection problem [6] requiring the identification of the skills acquired to solve previous tasks from which to transfer knowledge to best solve new tasks.

We sought a solution to this problem under two stringent conditions. First, TERL was not given any information regarding the similarity between the solved tasks and the new tasks and so it had to sample the actual performance of the possessed skills in the new task to infer their relevance for it. This condition was introduced to ensure the development of algorithms that are able to exploit the knowledge acquired while solving the new task without having any prior information on the similarity between the solved tasks and the new task. The solution proposed here is important because such knowledge on the effectiveness of previous solutions to solve the new task is always available in source-task selection problems. The idea is that we have developed efficient algorithms capable of fully exploiting such knowledge, we can develop and add to them other mechanisms able to exploit additional information, for example related to the similarity between tasks (see below). The second stringent condition was that the computational resources used by the system (here the experts) were constant. The use of the same resources to solve multiple tasks can support generalization and knowledge transfer but also introduces the well-known problem of catastrophic interference. The use of the same resources was dictated by our objective of building solutions that have biological plausibility.

Seeking the solution to the skill-to-skill transfer problem under the two conditions just described led to the development of two core mechanisms incorporated in TERL. The first mechanism accumulates evidence on the goodness of the skills, measured in terms of collected reward, during the

solution of the new task. This mechanism accumulates such evidence in parallel for all the experts, so continues to work well when the number of experts increases (it does not require the test of single experts separately one after the other). The second mechanism involves the decoupling of the responsibility signals used for the functioning and for the learning processes of the experts. This decoupling allows TERL to overcome the catastrophic interference problem.

The tests of TERL with a 2-D robotic dynamic arm showed that the system works both in conditions where tasks have to be learned and performed in sequence (i.e., with several blocks of trials for each task: this challenges the robustness of the system to catastrophic forgetting) or in random order (i.e., with one different task per each trial: this challenges the capacity of the system to rapidly decide which experts to use). Other tests showed how TERL scales up to larger numbers of experts and tasks with little additional computational costs. Also, in contrast to several non-parametric TRL systems that explicitly store experience in the form of state-action-state-reward tuples, the speed of functioning of TERL mechanisms is fully independent of the amount of knowledge already acquired. Preliminary experiments also showed that TERL can scale up to control a 3-D redundant dynamic robotic arm.

B. Future Work

The architecture of TERL could be further tested and improved in different ways in future work. In this paper, the experts were given only information about the task identity, but not on task features or task goals, to facilitate the development and study of the transfer algorithms tested here. For the same reasons, the gating network was informed only as to the task identity (with this the system only knows if two tasks are the same or different), and in particular it was not given a richer description of the task, or task goal, in terms of features (which might hint to the possible similarities between two tasks). In future work, it will be interesting: 1) to give the experts further information regarding the task to solve and 2) to give both experts and gating networks a rich description of the task/goal (e.g., as in the system proposed in [77] and [78], where a robot is informed on the position in space of a target to be hit with a dart or a ball). We speculate that this would have important consequences. Giving information about tasks to the experts might allow the system to encode knowledge of more than one task in the same expert, now not possible. In turn, this might allow the study of how the system encodes similar tasks in the same expert, and also allow the system to learn a number of tasks greater than the number of available experts. Giving the system a rich description of the task/goal might allow the gating networks to use such information to create useful predictions about the responsibility to assign to experts for each new task before experiencing it (here such prior signals were necessarily uniform as the system knew only the identity of tasks and had no information about their similarity). The mechanisms described here could then be used to accumulate evidence about the actual capacity of experts to solve the new task based on their test in the new task, so as to strengthen or weaken the responsibility signals predicted initially. The output of the gating networks before and after this

learning and accumulation of evidence could thus be interpreted as prior and posterior probability estimates of which experts are best to solve the new task similarly to what done in Bayesian approaches [79].

Future work might also further improve the mechanism that accumulates evidence on the goodness of experts for new tasks. In particular, here the mechanism searches only the average parameters of the Gaussian functions used to update the responsibility signals of experts, whereas their size (σ_A^2 and σ_C^2) is fixed. It might instead be possible to find a way to also estimate such parameters, similar to what is done in the mixture of Gaussian models [80], but taking into account the RL context considered here.

Future work should also further study and improve the mechanism used here to preserve the best experts against interference when facing some challenging conditions (e.g., when tasks are learned in sequence). Indeed, here we used a heuristic mechanism that lowers the learning rates of experts that achieve a high functioning responsibility in given tasks as this indicates that they can be reliably considered the best experts in those tasks. More principled mechanisms should hence be developed for this purpose.

Another aspect of the architecture to improve involves the treatment of noise. Here we adopted a simple solution as this was not the focus of this paper. Noise was kept low at the beginning of each trial during a fixed time interval that was sufficient to solve the task if the system had already learned it, and then noise was increased to let the system explore new solutions in case of failure to solve the task during such interval. Future work should find new solutions to this problem, in particular to regulate the level of noise depending on the system's performance suitably estimated on the basis of a meta-learning process (see [81], [82]).

Another aspect of TERL that might be developed in future work is that it currently assumes the existence of tasks to be solved, and that it receives information about the identity of the task to solve within each trial. These assumptions have important consequences on the system functioning, for example they give an episodic nature to the RL algorithms used by the system and this affects the algorithms to compute the global and local TD-errors, the exploration-exploitation noise, and the reset of some variables. Future work could make the system fully autonomous by endowing it with a component that is able to self-generate goals, tasks, and learning trials. For example, a recent work [82] has proposed a system that self-generates goals when the exploratory action of the simulated robot controlled by the model causes important effects in the environment, for example it turns on a spherical light by touching it. This system also self-generates "trial-terminations" either when it successfully accomplishes the currently pursued goal or when a time out elapses. TERL might be suitably integrated with similar mechanisms for goal (task) and trial self-generation.

Given the bio-inspired nature of the ingredients used to build it, TERL might also be used for investigating important behavioral and brain phenomena. One possible application is the study of the psychological processes of *assimilation and accommodation* postulated by Piaget [1], as started to do by [21] and [22]. Along this line, the system might

also be used to investigate the open-ended learning processes characterizing children [83] and leading them to progressively acquire a repertoire of increasingly sophisticated skills (as mentioned in Section I, this was one main motivation for designing TERL). In this respect, it would be interesting to endow the system with the autonomous capacities of self-generating goals and tasks, as mentioned above, and focusing on those that are learned with the highest learning rate [82], [84]–[87].

TERL can also be used to study psychophysical issues related to human motor learning, in particular problems related to how the motor system transfers knowledge between similar motor tasks. For example, the seminal work of [88] showed that learning and transfer of reaching movements are strongly facilitated when different tasks are learned in interleaved random trials rather than in sequential whole blocks of trials each focused on different tasks. Other studies have investigated how humans generalize a newly learned reaching skill, acquired to reach a given target point while compensating a disturbing force field, to other target points laying on a circumference centrad on the (fixed) starting point [89]–[91]. These studies show how transfer and generalization benefits only tasks involving target points that are very close in space. The architecture and transfer capabilities of TERL seem ideally suited to investigate these phenomena because they support transfer learning between similar tasks, and these transfer processes are parameterized under many respect, and because they can be linked to the architecture and functioning of the brain (see below; see also [92], for a review on experiments, models, and possible approaches to investigate these issues).

Some aspects of the architecture and functioning of TERL have been inspired by the brain organization and so they might be leveraged to model and investigate some open issues in neuroscience. In general, RL algorithms are suitable to study the organization and learning mechanisms of hierarchical architecture [93]. Within a biological perspective, TERL architecture might be used to model basal ganglia-cortical loops [101], a brain hierarchical system playing a key role in trial-and-error learning in organisms [48], [96]–[98]. For example, following previous works [18], [33], [99], the capacity of TERL gating networks to assign different tasks/goals to different experts might be used to model and study the mechanisms with which basal ganglia and cortex form separated loops and channels dedicated to different sensory inputs, actuator outputs, and input-output mappings [96], [100].

The capacity of TERL to form copies of experts on the basis of RL processes could be used to model the duplication of neural modules in the brain [41], a process that has been obtained through evolutionary algorithms mimicking DNA-based neural duplication [44], [73] and has also been previously observed in other systems [42], [70], [72]. In this respects, we are not aware of previous systems regulating the duplication process during learning on the basis of dedicated mechanisms such as the ranking and learning-responsibility mechanisms described here.

APPENDIX

Parameters Setting: The parameters of the model were set as indicated in Table II.

TABLE II
PARAMETERS OF THE MODEL USED TO CONTROL THE 2-D AND THE 3-D ROBOTIC ARM. THE TABLE REPORTS IN PARENTHESES THE PARAMETERS USED WITH THE 3-D ARM THAT DIFFER FROM THOSE USED WITH THE 2-D ARM

Parameter	Value
<i>State features</i>	
D	$21^2 (8^4)$
σ_f^2	1
<i>Responsibility accumulation process</i>	
κ	0.01
σ_{AG}^2	0.5 (0.3)
σ_{CG}^2	0.5 (0.3)
<i>Discount factor</i>	
γ	0.99
<i>Learning rates</i>	
η^{AG}	3.0
η^{CG}	1.0
η^A	1.2 (2.0)
η^C	0.012 (0.2)
<i>Noise generation</i>	
ϵ	20
τ	0.01
M	100
T'	600 (800)
ν	0.95
β	0.03

REFERENCES

- [1] J. Piaget, *The Origins of Intelligence in Children*. London, U.K.: Routledge and Kegan Paul, 1953.
- [2] N. E. Berthier and R. Keen, "Development of reaching in infancy," *Exp. Brain Res.*, vol. 169, no. 4, pp. 507–518, 2006.
- [3] H. Forssberg, A. C. Eliasson, H. Kinoshita, R. S. Johansson, and G. Westling, "Development of human precision grip I: Basic coordination of force," *Exp. Brain Res.*, vol. 85, no. 2, pp. 451–457, 1991.
- [4] E. Thelen, "Rhythmical stereotypies in normal human infants," *Animal Behav.*, vol. 27, pp. 699–715, Aug. 1979.
- [5] W. K. Geerts, C. Einspieler, J. Dibiasi, B. Garzarolli, and A. F. Bos, "Development of manipulative hand movements during the second year of life," *Early Human Develop.*, vol. 75, nos. 1–2, pp. 91–103, 2003.
- [6] M. E. Taylor and P. Stone, "Transfer learning for reinforcement learning domains: A survey," *J. Mach. Learn. Res.*, vol. 10, no. 1, pp. 1633–1685, 2009.
- [7] M. McCloskey and N. J. Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," *Psychol. Learn. Motiv.*, vol. 24, pp. 109–165, 1989.
- [8] R. M. French, "Catastrophic forgetting in connectionist networks," *Trends Cogn. Sci.*, vol. 3, no. 4, pp. 128–135, 1999.
- [9] J. Weng *et al.*, "Artificial intelligence. autonomous mental development by robots and animals," *Science*, vol. 291, no. 5504, pp. 599–600, 2001.
- [10] A. G. Barto, S. Singh, and N. Chentanez, "Intrinsically motivated learning of hierarchical collections of skills," in *Proc. Int. Conf. Develop. Learn. (ICDL)*, San Diego, CA, USA, 2004, pp. 112–119.
- [11] A. Stout, G. D. Konidaris, and A. G. Barto, "Intrinsically motivated reinforcement learning: A promising framework for developmental robot learning," in *Proc. AAAI Spring Symp. Develop. Robot.*, Stanford, CA, USA, 2005, pp. E1–E6.
- [12] P.-Y. Oudeyer, F. Kaplan, and V. V. Hafner, "Intrinsic motivation systems for autonomous mental development," *IEEE Trans. Evol. Comput.*, vol. 11, no. 2, pp. 265–286, Apr. 2007.
- [13] S. Thrun and T. M. Mitchell, "Lifelong robot learning," *Robot. Auton. syst.*, vol. 15, nos. 1–2, pp. 25–46, 1995.
- [14] G. Baldassarre and M. Mirolli, "What are the key open challenges for understanding the autonomous cumulative learning of skills?" *Newslett. Auton. Mental Develop. Tech. Committee*, vol. 7, no. 1, p. 11, 2010.
- [15] A. Lazaric, "Transfer in reinforcement learning: A framework and a survey," in *Reinforcement Learning: State of the Art*, M. Wiering and M. van Otterlo, Eds. Berlin, Germany: Springer, 2012, pp. 143–174.
- [16] J. Alcock, *Animal Behavior: An Evolutionary Approach*, 6th ed. Sunderland, MA, USA: Sinauer Assoc., 1998.
- [17] S. P. Singh, "Transfer of learning by composing solutions of elemental sequential tasks," *Mach. Learn.*, vol. 8, no. 3, pp. 323–339, 1992.

- [18] G. Baldassarre, "A modular neural-network model of the basal ganglia's role in learning and selecting motor behaviours," *Cogn. Syst. Res.*, vol. 3, no. 1, pp. 5–13, 2002.
- [19] G. Baldassarre, "Planning with neural networks and reinforcement learning," Ph.D. dissertation, Comput. Sci. Dept., Univ. Essex, Colchester, U.K., 2002.
- [20] D. Caligiore, M. Mirolli, D. Parisi, and G. Baldassarre, "A bioinspired hierarchical reinforcement learning architecture for modeling learning of multiple skills with continuous states and actions," in *Proc. 10th Int. Conf. Epigenetic Robot. (EpiRob)*, vol. 149. Lund, Sweden, 2010, pp. 27–34.
- [21] P. Tommasino, D. Caligiore, M. Mirolli, and G. Baldassarre, "Reinforcement learning algorithms that assimilate and accommodate skills with multiple tasks," in *Proc. IEEE Int. Conf. Develop. Learn. Epigenetic Robot. (ICDL-EpiRob)*, San Diego, CA, USA, 2012, pp. 1–8.
- [22] D. Caligiore, P. Tommasino, V. Sperati, and G. Baldassarre, "Modular and hierarchical brain organization to understand assimilation, accommodation and their relation to autism in reaching tasks: A developmental robotics hypothesis," *Adap. Behav.*, vol. 22, no. 5, pp. 304–329, 2014.
- [23] R. S. Sutton, "Generalization in reinforcement learning: Successful examples using sparse coarse coding," in *Proc. Adv. Neural Inf. Process. Syst.*, Denver, CO, USA, 1996, pp. 1038–1044.
- [24] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 1998.
- [25] J. Kober and J. Peters, "Reinforcement learning in robotics: A survey," in *Reinforcement Learning: State of the Art*, M. Wiering and M. van Otterlo, Eds. Berlin, Germany: Springer, 2012, pp. 579–610.
- [26] H. van Hasselt, "Reinforcement learning in continuous state and action spaces," in *Reinforcement Learning: State of the Art*, M. Wiering and M. van Otterlo, Eds. Berlin, Germany: Springer, 2012, pp. 207–251.
- [27] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *Proc. Adv. Neural Inf. Process. Syst.*, Denver, CO, USA, 2000, pp. 1057–1063.
- [28] R. C. O'Reilly, "Biologically plausible error-driven learning using local activation differences: The generalized recirculation algorithm," *Neural Comput.*, vol. 8, no. 5, pp. 895–938, Jul. 1996.
- [29] A. G. Barto, R. S. Sutton, and C. W. Anderson, "Neuronlike adaptive elements that can solve difficult learning control problems," *IEEE Trans. Syst., Man, Cybern., Syst.*, vol. SMC-13, no. 5, pp. 834–846, Sep./Oct. 1983.
- [30] A. G. Barto, "Adaptive critics and the basal ganglia," in *Models of Information Processing in the Basal Ganglia*, J. C. Houk, J. L. Davis, and D. G. Beiser, Eds. Cambridge, MA, USA: MIT Press, 1995, pp. 215–232.
- [31] K. Doya, "Complementary roles of basal ganglia and cerebellum in learning and motor control," *Current Opinions Neurobiol.*, vol. 10, no. 6, pp. 732–739, 2000.
- [32] D. Joel, Y. Niv, and E. Ruppin, "Actor-critic models of the basal ganglia: New anatomical and computational perspectives," *Neural Netw.*, vol. 15, nos. 4–6, pp. 535–547, 2002.
- [33] M. Khamassi, L. Lachèze, B. Girard, A. Berthoz, and A. Guillot, "Actor-critic models of reinforcement learning in the basal ganglia: From natural to artificial rats," *Adap. Behav.*, vol. 13, no. 2, pp. 131–148, 2005.
- [34] W. Schultz, P. Dayan, and P. R. Montague, "A neural substrate of prediction and reward," *Science*, vol. 275, no. 5306, pp. 1593–1599, 1997.
- [35] W. Schultz, "Getting formal with dopamine and reward," *Neuron*, vol. 36, no. 2, pp. 241–263, 2002.
- [36] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton, "Adaptive mixtures of local experts," *Neural Comput.*, vol. 3, no. 1, pp. 79–87, 1991.
- [37] R. A. Jacobs, M. I. Jordan, and A. G. Barto, "Task decomposition through competition in a modular connectionist architecture: The what and where vision tasks," *Cogn. Sci.*, vol. 15, no. 2, pp. 219–250, 1991.
- [38] D. McFarland, *Animal Behavior*, 2nd ed. Harlow, U.K.: Longman Group, 1993.
- [39] M. Mirolli, F. Mannella, and G. Baldassarre, "The roles of the amygdala in the affective regulation of body, brain, and behaviour," *Connect. Sci.*, vol. 22, no. 3, pp. 215–245, 2010.
- [40] R. E. Passingham and S. P. Wise, *The Neurobiology of the Prefrontal Cortex: Anatomy, Evolution, and the Origin of Insight*, vol. 50. Oxford, U.K.: Oxford Univ. Press, 2012.
- [41] D. Meunier, R. Lambiotte, and E. T. Bullmore, "Modular and hierarchically modular organization of brain networks," *Front. Neurosci.*, vol. 4, p. 200, Dec. 2010.
- [42] M. Ring, T. Schaul, and J. Schmidhuber, "The two-dimensional organization of behavior," in *Proc. Int. Conf. Develop. Learn. Epigenetic Robot. (ICDL)*, Trondheim, Norway, 2011, pp. 1–8.
- [43] M. Ring and T. Schaul, "The organization of behavior into temporal and spatial neighborhoods," in *Proc. IEEE Int. Conf. Develop. Learn. Epigenetic Robot. (ICDL)*, San Diego, CA, USA, 2012, pp. 1–6.
- [44] C. Fernando, K. K. Karishma, and E. Szathmáry, "Copying and evolution of neuronal topology," *PLoS One*, vol. 3, no. 11, p. e3775, 2008.
- [45] M. S. A. Graziano, *The Intelligent Movement Machine: An Ethological Perspective on the Primate Motor System*. Oxford, U.K.: Oxford Univ. Press, 2009.
- [46] T. N. Affalo and M. S. Graziano, "Organization of the macaque extrastriate visual cortex re-examined using the principle of spatial continuity of function," *J. Neurophysiol.*, vol. 105, no. 1, pp. 305–320, 2011.
- [47] K. Gurney, T. J. Prescott, and P. Redgrave, "A computational model of action selection in the basal ganglia. I. a new functional anatomy," *Biol. Cybern.*, vol. 84, no. 6, pp. 401–410, 2001.
- [48] A. M. Graybiel, "The basal ganglia and chunking of action repertoires," *Neurobiol. Learn. Memory*, vol. 70, nos. 1–2, pp. 119–136, 1998.
- [49] G. Rizzolatti and G. Luppino, "The cortical motor system," *Neuron*, vol. 31, no. 6, pp. 889–901, 2001.
- [50] M. S. Graziano, C. S. Taylor, T. Moore, and D. F. Cooke, "The cortical control of movement revisited," *Neuron*, vol. 36, no. 3, pp. 349–362, 2002.
- [51] S. Thill, D. Caligiore, A. M. Borghi, T. Ziemke, and G. Baldassarre, "Theories and computational models of affordance and mirror systems: An integrative review," *Neurosci. Biobehav. Rev.*, vol. 37, no. 3, pp. 491–521, 2013.
- [52] E. K. Miller and J. D. Cohen, "An integrative theory of prefrontal cortex function," *Annu. Rev. Neurosci.*, vol. 24, no. 1, pp. 167–202, 2001.
- [53] J. M. Fuster, "Upper processing stages of the perception-action cycle," *Trends Cogn. Sci.*, vol. 8, no. 4, pp. 143–145, 2004.
- [54] D. Caligiore, A. M. Borghi, D. Parisi, and G. Baldassarre, "TROICALS: A computational embodied neuroscience model of compatibility effects," *Psychol. Rev.*, vol. 117, no. 4, pp. 1188–1228, 2010.
- [55] K. Doya, "Reinforcement learning in continuous time and space," *Neural Comput.*, vol. 12, no. 1, pp. 219–245, 2000.
- [56] M. van Otterlo and M. Wiering, "Reinforcement learning and Markov decision processes," in *Reinforcement Learning: State of the Art*, M. Wiering and M. van Otterlo, Eds. Berlin, Germany: Springer, 2012, pp. 3–42.
- [57] Y. Uno, M. Kawato, and R. Suzuki, "Formation and control of optimal trajectory in human multijoint arm movement—Minimum torque-change model," *Biol. Cybern.*, vol. 61, no. 2, pp. 89–101, 1989.
- [58] L. Sciacivico and B. Siciliano, *Modelling and Control of Robot Manipulators*. Berlin, Germany: Springer, 2000.
- [59] D. Caligiore, E. Guglielmelli, A. M. Borghi, D. Parisi, and G. Baldassarre, "A reinforcement learning model of reaching integrating kinematic and dynamic control in a simulated arm robot," in *Proc. IEEE Int. Conf. Develop. Learn. (ICDL)*, Ann Arbor, MI, USA, 2010, pp. 211–218.
- [60] D. Caligiore, D. Parisi, and G. Baldassarre, "Integrating reinforcement learning, equilibrium points, and minimum variance to understand the development of reaching: A computational model," *Psychol. Rev.*, vol. 121, no. 3, pp. 389–421, 2014.
- [61] G. Sandini, G. Metta, and D. Vernon, "The iCub cognitive humanoid robot: An open-system research platform for enactive cognition," in *50 Years of Artificial Intelligence. Essays Dedicated to the 50th Anniversary of Artificial Intelligence*, M. Lungarella, F. Iida, J. Bongard, and R. Pfeifer, Eds., vol. 4850. Berlin, Germany: Springer, 2007.
- [62] A. G. Barto and S. Mahadevan, "Recent advances in hierarchical reinforcement learning," *Discr. Event Dyn. Syst.*, vol. 13, no. 4, pp. 341–379, 2003.
- [63] R. S. Sutton, D. Precup, and S. Singh, "Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning," *Artif. Intell.*, vol. 112, nos. 1–2, pp. 181–211, 1999.
- [64] G. Konidaris and A. Barto, "Building portable options: Skill transfer in reinforcement learning," in *Proc. 20th Int. Joint Conf. Artif. Intell.*, vol. 7. Hyderabad, India, 2007, pp. 895–900.

- [65] C. J. C. H. Watkins and P. Dayan, "Q-learning," *Mach. Learn.*, vol. 8, no. 3, pp. 279–292, 1992.
- [66] R. S. Sutton *et al.*, "Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction," in *Proc. 10th Int. Conf. Auton. Agents Multiagent Syst. (AAMAS)*, Taipei, Taiwan, May 2011, pp. 761–768.
- [67] F. Fernández and M. Veloso, "Probabilistic policy reuse in a reinforcement learning agent," in *Proc. 5th Int. Joint Conf. Auton. Agents Multiagent Syst.*, Hakodate, Japan, 2006, pp. 720–727.
- [68] A. Lazaric, M. Restelli, and A. Bonarini, "Transfer of samples in batch reinforcement learning," in *Proc. 25th Int. Conf. Mach. Learn.*, Helsinki, Finland, 2008, pp. 544–551.
- [69] D. M. Wolpert and M. Kawato, "Multiple paired forward and inverse models for motor control," *Neural Netw.*, vol. 11, nos. 7–8, pp. 1317–1329, 1998.
- [70] K. Doya, K. Samejima, K.-I. Katagiri, and M. Kawato, "Multiple model-based reinforcement learning," *Neural Comput.*, vol. 14, no. 6, pp. 1347–1369, 2002.
- [71] N. Sugimoto, M. Haruno, K. Doya, and M. Kawato, "Mosaic for multiple-reward environments," *Neural Comput.*, vol. 24, no. 3, pp. 577–606, Mar. 2012.
- [72] M. Ring and T. Schaul, "Q-error as a selection mechanism in modular reinforcement-learning systems," in *Proc. Int. Joint Conf. Artif. Intell. (IJCAI)*, vol. 22, Barcelona, Spain, 2011, pp. 1452–1458.
- [73] R. Calabretta, S. Nolfi, D. Parisi, and G. P. Wagner, "Duplication of modules facilitates the evolution of functional specialization," *Artif. Life*, vol. 6, no. 1, pp. 69–84, 2000.
- [74] C. Fernando, "Design for a darwinian brain: Part 1. Philosophy and neuroscience," *arXiv, arXiv preprint 1303.7200*, 2013.
- [75] C. Fernando and V. Vasas, "Design for a darwinian brain: Part 2. Cognitive architecture," *arXiv, arXiv preprint 1303.7201*, 2013.
- [76] N. Mehta, S. Natarajan, P. Tadepalli, and A. Fern, "Transfer in variable-reward hierarchical reinforcement learning," *Mach. Learn.*, vol. 73, no. 3, pp. 289–312, 2008.
- [77] B. C. da Silva, G. Konidaris, and A. G. Barto, "Learning parameterized skills," in *Proc. 29th Int. Conf. Mach. Learn. (ICML)*, Edinburgh, U.K., 2012, pp. 1679–1686.
- [78] B. C. da Silva, G. Baldassarre, G. Konidaris, and A. G. Barto, "Learning parameterized motor skills on a humanoid robot," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, Hong Kong, 2014, pp. 5239–5244.
- [79] A. Wilson, A. Fern, S. Ray, and P. Tadepalli, "Multi-task reinforcement learning: A hierarchical Bayesian approach," in *Proc. 24th Int. Conf. Mach. Learn.*, Corvallis, OR, USA, 2007, pp. 1015–1022.
- [80] R. Jacobs, "Mixture models," Dept. Brain Cogn. Sci., Univ. Rochester, Rochester, NY, USA, Tech. Rep., 2008.
- [81] V. G. Santucci, G. Baldassarre, and M. Mirolli, "Autonomous selection of the 'what' and the 'how' of learning: An intrinsically motivated system tested with a two armed robot," in *Proc. 4th Joint IEEE Int. Conf. Develop. Learn. Epigenetic Robot. (ICDL-EpiRob)*, Genoa, Italy, Oct. 2014, pp. 434–439.
- [82] V. G. Santucci, G. Baldassarre, and M. Mirolli, "GRAIL: A goal-discovering robotic architecture for intrinsically-motivated learning," *IEEE Trans. Cogn. Develop. Syst.*, vol. 8, no. 3, pp. 214–231, Sep. 2016.
- [83] G. Baldassarre and M. Mirolli, Eds., *Intrinsically Motivated Learning in Natural and Artificial Systems*. Berlin, Germany: Springer, 2013.
- [84] A. McGovern and A. Barto, "Automatic discovery of subgoals in reinforcement learning using diverse density," Dept. Comput. Sci., Faculty Publication Series, 2001.
- [85] G. Konidaris and A. G. Barto, "Skill discovery in continuous reinforcement learning domains using skill chaining," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, Vancouver, BC, Canada, 2009, pp. 1015–1023.
- [86] V. G. Santucci, G. Baldassarre, and M. Mirolli, "Intrinsic motivation mechanisms for competence acquisition," in *Proc. IEEE Int. Conf. Develop. Learn. (ICDL-EpiRob)*, San Diego, CA, USA, 2012, pp. 1–6.
- [87] V. G. Santucci, G. Baldassarre, and M. Mirolli, "Which is the best intrinsic motivation signal for learning multiple skills?" *Front. Neurobot.*, vol. 7, pp. 1–14, Nov. 2013.
- [88] J. B. Shea and R. L. Morgan, "Contextual interference effects on the acquisition, retention, and transfer of a motor skill," *J. Exp. Psychol. Human Learn. Memory*, vol. 5, no. 2, pp. 179–187, 1979.
- [89] R. Shadmehr and F. A. Mussa-Ivaldi, "Adaptive representation of dynamics during learning of a motor task," *J. Neurosci.*, vol. 14, no. 5, pp. 3208–3224, 1994.
- [90] R. Shadmehr and S. P. Wise, Eds., *The Computational Neurobiology of Reaching and Pointing*. Cambridge, MA, USA: MIT Press, 2005.
- [91] R. Shadmehr and S. Mussa-Ivaldi, *Biological Learning and Control: How the Brain Builds Representations, Predicts Events, and Makes Decisions*. Cambridge, MA, USA: MIT Press, 2012.
- [92] L. Lonini, C. Dimitrakakis, C. A. Rothkopf, and J. Triesch, "Generalization and interference in human motor control," in *Computational and Robotic Models of the Hierarchical Organization of Behavior*, G. Baldassarre and M. Mirolli, Eds. Berlin, Germany: Springer, 2013, pp. 155–176.
- [93] G. Baldassarre and M. Mirolli, Eds., *Computational and Robotic Models of the Hierarchical Organization of Behavior*. Berlin, Germany: Springer-Verlag, 2013.
- [94] M. Mirolli, V. G. Santucci, and G. Baldassarre, "Phasic dopamine as a prediction error of intrinsic and extrinsic reinforcements driving both action acquisition and reward maximization: A simulated robotic study," *Neural Netw.*, vol. 39, pp. 40–51, Mar. 2013.
- [95] V. G. Fiore *et al.*, "Keep focussing: Striatal dopamine multiple functions resolved in a single mechanism tested in a simulated humanoid robot," *Front. Psychol. Cogn. Sci.*, vol. 5, no. 124, pp. 1–17, 2014.
- [96] G. E. Alexander, M. R. DeLong, and P. L. Strick, "Parallel organization of functionally segregated circuits linking basal ganglia and cortex," *Annu. Rev. Neurosci.*, vol. 9, no. 1, pp. 357–381, 1986.
- [97] H. H. Yin and B. J. Knowlton, "The role of the basal ganglia in habit formation," *Nature Rev. Neurosci.*, vol. 7, no. 6, pp. 464–476, 2006.
- [98] M. M. Botvinick, Y. Niv, and A. C. Barto, "Hierarchically organized behavior and its neural foundations: A reinforcement learning perspective," *Cognition*, vol. 113, no. 3, pp. 262–280, 2009.
- [99] M. Khamassi, L.-E. Martinet, and A. Gullot, "Combining self-organizing maps with mixtures of experts: Application to an actor-critic model of reinforcement learning in the basal ganglia," in *From Animals to Animals 9*, S. Nolfi *et al.*, Eds. Berlin, Germany: Springer-Verlag, 2006, pp. 394–405.
- [100] J. W. Mink, "The basal ganglia: Focused selection and inhibition of competing motor programs," *Progr. Neurobiol.*, vol. 50, no. 4, pp. 381–425, 1996.
- [101] G. Baldassarre *et al.*, "Intrinsically motivated action-outcome learning and goal-based action recall: A system-level bio-constrained computational model," *Neural Netw.*, vol. 41, pp. 168–187, May 2013, doi: 10.1016/j.neunet.2012.09.015



Paolo Tommasino received the B.E. and M.E. degrees from the Università Campus Bio-Medico, Rome, Italy, in 2008 and 2011, respectively. He is currently working toward the Ph.D. degree with the Nanyang Technological University, Singapore.

From 2011 to 2012, he served as a Research Fellow with the Laboratory of Computational Embodied Neuroscience, Institute of Cognitive Science and Technologies, National Research Council, Rome, where he researched on reinforcement learning algorithms for the project funded by the European Commission "IM-CLeVer: Intrinsically Motivated Cumulative Learning Versatile Robots." In 2012, he joined the School of Mechanical and Aerospace Engineering, Nanyang Technological University. His current research interests include mechatronic technologies, novel robotic devices for rehabilitation and motor therapy after neurological injury, haptics, bio-inspired algorithms for controlling redundant manipulators, bio-inspired actuators, and computational embodied neuroscience.



Daniele Caligiore received the master's degree in electronics engineering from the University of Catania, Catania, Italy, in 2003, and the Ph.D. degree in biomedical engineering from the University Campus Bio-Medico di Roma, Rome, Italy, in 2011.

He was a Visiting Scholar with the Centre for Robotics and Neural Systems and the School of Psychology, University of Plymouth, Plymouth, U.K., and the Embodied Cognition Laboratory, University of Bologna, Bologna, Italy. Since 2004, he has been a Researcher with the Institute of

Cognitive Sciences and Technologies, National Research Council, Rome. He has participated in several European projects in the field of embodied cognition and developmental robotics, such as MindRACES—From Reactive to Anticipatory Cognitive Embodied Systems, ROSSI—Emergence of Communication in Robots Through Sensorimotor and Social Interaction, and IM-CLeVeR—Intrinsically Motivated Cumulative Learning Versatile Robots. He has authored/co-authored about 70 peer-reviewed publications appeared in international journals, books, and conference proceedings. His current research interests include developmental robotics, embodied cognition, system-level computational neuroscience, brain cortical and subcortical hierarchies, reinforcement learning.

Dr. Caligiore was a Guest-Editor for a consensus paper of the journal “Cerebellum” titled “Toward a systems-level view of cerebellar function: the interplay between cerebellum, basal ganglia and cortex.”



Marco Mirolli received the B.A. and M.A. degrees in philosophy and the Ph.D. degree in cognitive sciences from the University of Siena, Siena, Italy, in 2001 and 2006, respectively.

He is a Researcher with the Institute of Cognitive Sciences and Technologies, National Research Council, Rome, Italy. He has been (mostly) studying brain and behavior through computer simulations and robotic models. In particular, he has been researching on the evolution of communication and language, the role of language as a cognitive tool, the

concept of representations in cognitive science, intrinsic motivations, and the biological bases of conditioning, motivations, and emotions. He has co-edited four books, including the book entitled *Intrinsically Motivated Learning in Natural and Artificial Systems* (Springer) and published over 60 peer-reviewed papers. His current research interests include understanding the relationships between the body and the mind through theoretical analysis, computational modeling, and empirical research.



Gianluca Baldassarre received the B.A. and M.A. degrees in economics and the M.Sc. degree in cognitive psychology and neural networks with the Sapienza University of Rome, Rome, Italy, in 1998 and 1999, respectively, and the Ph.D. degree in computer science from the University of Essex, Colchester, U.K., in 2003, with a focus on planning with neural networks.

He was a Post-Doctoral Research Fellow with the Italian Institute of Cognitive Sciences and Technologies, National Research Council, Rome,

researching on swarm robotics until 2005. Since 2006, he has been a Researcher, and coordinates the Research Group called Laboratory of Computational Embodied Neuroscience, at the same institute. From 2006 to 2009, he was a Team Leader of the EU project “ICEA—Integrating Cognition Emotion and Autonomy.” From 2009 to 2013, he was the Coordinator of the European Integrated Project “IM-CLeVeR—Intrinsically-Motivated Cumulative-Learning Versatile Robots.” In 2016–2020, he is being the Coordinator of the European FET-OPEN project “GOAL-Robots—Goal-based Open-ended Autonomous Learning Robots.” He makes research with LOCEN by following two interdisciplinary research approaches: computational models aiming to understand the functioning of brain and behavior; machine-learning/autonomous-robotics approaches aiming to produce technologically useful robots. He has over 100 international peer-review publications. His current research interests include cumulative learning of multiple sensorimotor skills driven by extrinsic and intrinsic motivations, and higher-level cognition grounded on sensorimotor behavior with a particular focus on goal-directed behavior.