# Robotic self-representation improves manipulation skills and transfer learning

Phuong D.H. Nguyen[1], Manfred Eppe[1], Stefan Wermter[1]

*Abstract*— Cognitive science suggests that the self-representation is critical for learning and problem-solving. However, there is a lack of computational methods that relate this claim to cognitively plausible robots and reinforcement learning. In this paper, we bridge this gap by developing a model that learns bidirectional action-effect associations to encode the representations of body schema and the peripersonal space from multisensory information, which is named multimodal BidAL. Through three different robotic experiments, we demonstrate that this approach significantly stabilizes the learning-based problem-solving under noisy conditions and that it improves transfer learning of robotic manipulation skills.

## I. INTRODUCTION

Humans and other biological agents depend on appropriate representations of their body and their surroundings to facilitate their activities in safety and comfort. These representations, known as body schema and peripersonal space (PPS) (see Fig. 1), result from the integration of different sensorimotor modalities that are involved when physically interacting with the environment [1]–[4].

Consequently, the body schema and peripersonal space are not innate but develop incrementally. This developmental process starts in infants through self-exploration and motor babbling, and continues later through goal-directed physical and social interactions [5]–[7]. Acquiring a body schema and a peripersonal space representation involves associative learning, a mechanism enabling infants to detect the sensorimotor contingencies in their environment. However, associative learning alone is not sufficient to explain the ability of humans to generate goal-directed actions. The conversion of the learned contingencies into goal-directed actions is related to another central ability that is only little investigated in computational methods: The ability to distinguish between self-caused body-schematic sensory effects and externally caused sensory effects [8], [9]. But how can we model these bidirectional body-schematic action-effect associations computationally, such that an agent can distinguish between its own body, the peripersonal space and the external world? And how does this affect the learning and problem-solving performance of artificial agents?

In this work, we argue that artificial agents/robots require two vital elements to develop the body-schema and peripersonal space representations: (1) a multimodal sensory
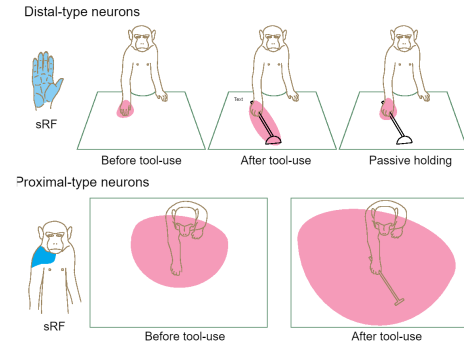
[1]Phuong D.H. Nguyen, Manfred Eppe and Stefan Wermter are with Department of Informatics, University of Hamburg, Hamburg, Germany {pnguyen, eppe, wermter}@informatik.uni-hamburg.de



Fig. 1: The body schema and peripersonal space representations during tool-use (from [10] ).

integration model, and (2) a mechanism to learn bidirectional associations between actions and effects. The core hypotheses of this paper are that (i) intrinsic rewards created by predictability of an agent's self-caused multimodal sensory effects enable the learning of bidirectional action-effect associations, (ii) these associations implement body schema and peripersonal space representations that make the robotic action-selection robust to noise, and (iii) the abstract body-schematic representations improve the transfer learning performance.

We address our hypotheses by integrating multisensory association with a forward and inverse mode and a reinforcement learning policy. Our proposed framework (Fig. 2) is based on the bidirectional associative model by Pathak *et al.* [11]. The model combines the learning of an inverse model and a forward model to emulate the bidirectional action-effect associations. The key property of the model is that it automatically eliminates information from the sensor signals over which the agent has no control, including noise and non-deterministic effects. This is achieved by training simultaneously an abstraction function $\phi(\cdot)$, a forward model $f(\cdot)$ and an inverse model $g(\cdot)$, such that the forward and inverse model operate in the learned abstract state. The abstraction function learns to ignore sensory information that are not directly related to the agent's own actions because the forward and inverse model operate in the abstract space, ignoring all sensorimotor contigencies that are not self-caused [11]. Therefore, the model implicitly learns to distinguish between the agent's own body and the external world. Moreover, the model implies the representations of the agent's body schema and its peripersonal dynamics. We exploit this property here to investigate: (1) whether the learned representations benefit from multimodal sensor signals, (2) to what extent they are invariant to noise and (3)

whether they are transferable to learn a new skill.

For this investigation, we extend Pathak *et al.*'s approach as follows: (i) We use multimodal sensory input instead of only visual input for our model; (ii) we build on intrinsic motivation to minimize the prediction error, especially in noisy conditions; and (iii) we investigate the transfer learning performance. In addition, we extend the discrete action space used by Pathak *et al.* towards a continuous action space. To realize these extensions, we build on several mechanisms and concepts related to multisensory representation learning, forward models and intrinsic motivation.

## II. BACKGROUND AND RELATED WORK

Our proposed model in this paper relates to a body of literature in following main topics:

**Multisensory representation learning of the body schema** enables humans to perform pose estimation of their body parts, and coordinate transformation between sensory sources, which, ultimately, enables action [12], [13]. Robotic and computational models of body-schematic representations mostly focus on exploiting sensory information from proprioceptive and tactile sensing [14]–[16], or proprioceptive and visual sensing [17]–[19] and cast the representation learning as calibration, pose estimation or visuomotor mapping.

**The peripersonal space representation** serves as an interface between the agent's body and the environment [3]. Existing robotic and computational models construct the PPS representation from sensory data including vision, audio, touch and proprioception [20]–[28]. Most of the approaches base on the random movements of joints inspired by infants' motor babbling to generate the training data.

**Forward models** are computational models that map the current state of the system to the next state through actions. Some approaches utilize this forward model to learn the imitating actions from multisensory input [29]–[32]; others employ the forward model to learn the single sensory embedding for control, e.g. [11], [33]–[36]. Differently, in some neurorobotics models, forward model plays the core role in high-level cognitive functions such as self/other distinction, sense of agency or body ownership [37]–[39].

**Intrinsic motivation** is an internal system that drives human to "engage in particular activities for their own sake, rather than a step towards solving practical problem" [40]. Some recent models revisit this concept with the computational models of neural networks (cf. [41] for an overview). For example, Dilokthanakul *et al.* [42] implement the intrinsic motivation as the changes in the image features of two consecutive frames, allowing gaming agents to learn by maximizing the change of the visual representations. Pathak *et al.* [11] propose to use prediction error of an additional forward model as intrinsic motivation to drive gaming agents to explore the space. Pathak *et al.* [43] further extend the idea by using an ensemble of forward functions and exploiting the disagreement among prediction errors in the ensemble as the intrinsic motivation. Röder *et al.* [44] adopt the ideas of the intrinsic motivation from [11] but with only proprioceptive input instead of visual input.
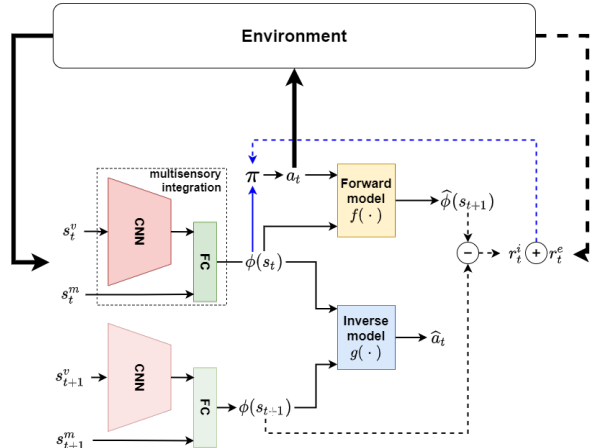


Fig. 2: Overview of the learning framework.

## III. METHODOLOGY

The objective of our research is to enable a robot to learn to solve a given task in a general goal-directed way, by generalizing the learnt knowledge to any goal $g$ in the set of possible goals $\mathcal{G}$. For example, considering the toy example of a planar reaching task (see Fig. 3a), we aim for a framework that enables a 2DoF-robot to learn to move its end-effector to every position in the plane within the robot's reachable region. More advanced applications involve self-locomotion and object manipulation, where an agent should be able to generalize over goal locations of itself or of external objects.

### A. Overview of the framework

To achieve the desired generalization and transfer learning abilities, we train a reinforcement learning policy in an abstract generalized state representation imposed by the agent's body schema and peripersonal space model. This model is learned with the bio-inspired learning framework presented in Fig. 2. The forward model $f(\cdot)$ predicts a sensory effect $\hat{\phi}(s_{t+1})$ from a currently conducted action $a_t$ and the currently perceived sensory state representation $\phi(s_t)$. The *policy* $\pi$ generates motor actions $a_t$ under constraints exerted by the environment and under consideration of the prediction error $e_{t+1}$ of the forward model. Both the forward model and the policy operate in the latent space of the multimodal sensory input, which is compressed by the *multisensory integration* process. We specify the operation of these modules as follows:

Multisensory representations:
$$\phi(s_t) = \phi^{PPS}(s_t^{e\cup i}) \bigcup \phi^{body}(s_t^i) \qquad (1)$$

Forward model:
$$\hat{\phi}(s_{t+1}) = f\big(\phi(s_t), a_t; \boldsymbol{\theta}_F\big) \qquad (2)$$

Prediction error:
$$e_{t+1} = \frac{1}{2} \left\| \hat{\phi}(s_{t+1}) - \phi(s_{t+1}) \right\|_2^2 \qquad (3)$$

where $\phi^{PPS}(s_t^{e\cup i})$ denotes the representation of the PPS, $\phi^{body}(s_t^i)$ denotes the body schema representation, and $\boldsymbol{\theta}_F$

denotes parameters of a two-layer fully-connected neural network approximating the function $f(\cdot)$.

Instead of using the prediction error as an intrinsic motivation to drive agents in seeking for the new information, which is suitable for the navigation task [11], we propose to use the prediction error for a different purpose in this work. For our agents, the lower the prediction error means the better their ability to anticipate their own action, which is represented by the intrinsic reward $r_t^i$ as follows:

$$r_t^i = -\frac{\lambda}{2} \left\| \hat{\phi}(s_{t+1}) - \phi(s_{t+1}) \right\|_2^2 \tag{4}$$

where $\lambda$ is a hyperparameter to decide the importance of the intrinsic reward within the whole reward the agents receive.

In the following sections III-B and III-C we describe how the intrinsic reward is combined with the extrinsic reward to train an actor-critic reinforcement learning method (cf. Eq. 8). Herein, all modules learn simultaneously through the agent's interactive experience in the environment.

### B. Multisensory integration with neural network

We aim to enable robots to exploit the relevant information from different sources of their sensory input, to construct a representation of the environment and finally to facilitate the task learning. Therefore, we first implement a neural network for visuo-motor integration based on Nguyen *et al.* [19], and construct our input preprocessing network with two branches, one for vision and one for proprioception. The former branch consists of four convolution layers with *Elu* activation and a fully-connected layer, while the latter branch contains one layer of fully-connected units. The visual and proprioceptive features are concatenated and combined by another fully-connected layer to produce a compressed latent feature $\phi(s_t)$ of the environment (including the robot itself). The construction of multisensory network is presented in the left side of our general architecture in Fig. 2.

Unlike Nguyen *et al.* [19] and other methods that learn this representation separately, e.g. Watter *et al.* [33], Zambelli *et al.* [29], we adopt the approach by Pathak *et al.* [11] for state representation learning within the reinforcement framework. While Pathak *et al.* [11] consider only visual input, we extend it to multisensory input. In addition, we also employ the inverse model $g(\cdot)$ for representation learning from multiple inputs by minimizing the inverse loss as follows:

$$\mathcal{L}_I\left(\hat{a}_t, a_t\right) = \frac{1}{2} \left\| \hat{a}_t - a_t \right\|_2^2, \tag{5}$$

where $\hat{a}_t = g\left(s_t, s_{t+1}; \boldsymbol{\theta}_I\right)$ is the estimated action of $a_t$ through the inverse mapping function $g(\cdot)$, approximated by a two-layer fully-connected neural network with parameters $\boldsymbol{\theta}_I$. By minimizing the difference between the estimated and real action, the learnt inverse model plays as an encoder of the relevant information for the task from multiple sensory input [11].

### C. Reinforcement learning policy

We consider a reinforcement learning setting, where an agent interacts with an environment and tries to maximize the long-term expected reward. The interacting environment can be defined as a set of state (or observation) $\mathcal{S}$, an action set $\mathcal{A}$, a reward function $r : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$, transition probability $\rho(s_{t+1}|s_t, a_t) : \mathcal{S} \times \mathcal{S} \times \mathcal{A}$, and a discount rate $\gamma \in [0, 1]$.

In our setting, we employ the policy gradient approach that allows agent to select actions directly through a parameterized policy instead of consulting the value function. This means, at time $t$, the agent takes an action $a_t$ drawn from the policy $\pi(a|s, \boldsymbol{\theta})$, a probability distribution over the action space given the current state $s_t$ with parameter vector $\boldsymbol{\theta}$.

Any policy gradient-based reinforcement learning algorithm can be applied to construct the *policy* in our framework. In this work, we employ and validate the Deep deterministic policy gradient algorithm (DDPG) [45] with continuous actions in combination with our BidAL method. The previous implementation by Pathak *et al.* [11] uses AC3 [46] with a discrete action space. Moreover, we do not manually design any task-specific rewards, but employ the spare-reward scheme for the external reward $r_t^e$ in our proposal. This means the robots receive 0 reward for successfully completing the desired task and -1 for failing the task, as defined in Eq. 6:

$$r_t^e = -\left[ |g - s_t| > \epsilon \right] \tag{6}$$

where $\epsilon$ is a reasonable threshold value to determine whether an achieved state $s_t$ is close enough to an desired goal $g$ to consider the goal achieved. We implement our goal-directed method using deep deterministic policy gradient (DDPG), universal value function approximators (UVFA)[47] and hindsight experience replay (HER)[48], as described below:

*1) Deep deterministic policy gradient (DDPG):* While policy gradient methods in general refer to a parameterized, stochastic policy, deterministic policy gradient methods aim to learn parameters for a deterministic policy, $\mu_{\boldsymbol{\theta}} : \mathcal{S} \mapsto \mathcal{A}$. Silver *et al.* [49] shows that the deterministic policy gradient (DPG) exists as a special case of stochastic policy, but can be estimated more efficiently. The efficient exploration of a deterministic policy can be guaranteed with an off-policy algorithm in which actions are chosen according to a stochastic behaviour policy $\beta(a|s)$, but to learn about a deterministic target policy [49]. The off-policy deterministic policy gradient is written as:

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} J_\beta(\mu_{\boldsymbol{\theta}}) &\approx \int_{\mathcal{S}} \rho^\beta(s) \nabla_{\boldsymbol{\theta}} \mu_{\boldsymbol{\theta}} Q^\mu(s, a) ds \\ &= \mathbb{E}_{s \sim \rho^\beta} \left[ \nabla_{\boldsymbol{\theta}} \mu_{\boldsymbol{\theta}}(s) \nabla_a Q^\mu(s, a)|_{a=\mu_{\boldsymbol{\theta}}(s)} \right] \end{aligned} \tag{7}$$

where $\rho^\beta$ denotes the state distribution of $\beta(a|s)$.

DDPG [45] extends the actor-critic approach of DPG with neural network function approximators. The actor network of $\mu_{\boldsymbol{\theta}}$ is updated with deterministic policy gradient as Eq. 7, while the critic network estimates the action-value function $Q(s, a)$ by minimizing the loss in Eq. 8.

$$\mathcal{L} = \mathbb{E}_{s_t \sim \rho^\beta, a_t \sim \beta} \left[ \left( Q(s_t, a_t | \boldsymbol{\theta}^Q) - y_t \right)^2 \right]$$

with target: $y_t = r_t + \gamma Q(s_{t+1}, \mu s_{t+1} | \boldsymbol{\theta}^Q)$ $\qquad$ (8)

and reward: $r_t = r_t^i + r_t^e$

Both actor and critic neural networks are updated by sampling from a finite sized replay buffer where tuples of exploration $(s_t, a_t, r_t, s_{t+1})$ have been stored.

*2) Universal value function approximation (UVFA):* In order to generalize the learnt policy to both the state space and the goal space, [47] proposes the UVFA approach to represent a set of optimal action-value functions by using a unified function approximator. This is achieved by extending the Q-function to depend not only on the state-action $(s, a)$ pair but also the goal $g$. In terms of function approximation with neural networks, we can concatenate $(s, g)$ or their embeddings $(\psi(s), \eta(g))$ for the actor network, and $(s, a, g)$ or $(\psi(s), a, \eta(g))$ for the critic network. Our networks follow the embeddings structure. This technique is important for our setup as it allow robots to learn a general goal-directed policy: Instead of achieving a specific task, e.g. reaching a certain position, they learn to complete a more general one, i.e. reaching every position within the reachable space. This also makes our approach different from the work by Pathak *et al.* [11].

*3) Hindsight experience replay (HER):* We employ this technique from Andrychowicz *et al.* [48] to enrich the replay buffer by assuming that some random unsuccessful achieved state is the goal state.

### D. Evaluation environments



(a) Planar reaching      (b) UR5 reaching

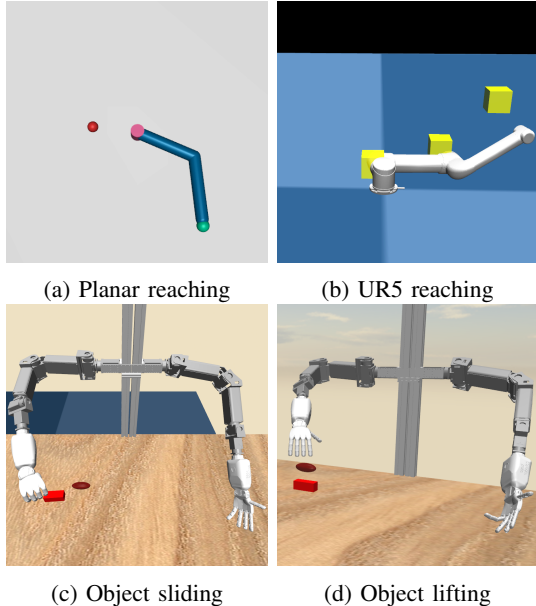(c) Object sliding      (d) Object lifting

Fig. 3: Experimental environments

We evaluate our method by conducting experiments in simulated environments based on the Mujoco software [50]. Depending on the specific environment, the setup can be varied but generally all robots have the common properties:

- They have access to cameras which provide rendered RGB frames of the environment, including parts of the robot itself. The rendered frames are then converted to gray format and scaled down to smaller dimension. This source of input is denoted as $\mathbf{s}_t^v$ in this work;

| Env. | $\mathbf{s}_t^v$ | $\mathbf{s}_t^m$ | $\mathbf{a}_t$ |
|---|---|---|---|
| Planar reaching | $\left[\mathbf{I}\right]_{(42,42,1)}$ | $\left[\mathbf{q}, \dot{\mathbf{q}}\right]_{(0,4)}$ | [0,2] |
| UR5 reaching | | $\left[\mathbf{q}, \dot{\mathbf{q}}\right]_{(0,6)}$ | [0,3] |
| Object sliding | | $\left[\mathbf{x}, \mathbf{q}, \dot{\mathbf{x}}, \dot{\mathbf{q}}\right]_{(0,40)}$ | [0,7]* |
| Object lifting | | | |

TABLE I: Summary of experimental environments. For object manipulation tasks[(*)], i.e. sliding and lifting, the agent can control the position of the arm in the Cartesian space, roll and tilt its hand, and move all fingers and its thumb abduction.

- Robots have access to encoder measurements of all joints composed of their body. Additionally, the end-effector and object information are available for object manipulation environment. We use the term "proprio-ception" for this input in an extended meaning, and denote it as $\mathbf{s}_t^m$ in this work;
- Agents generate action in only the joint space (in reaching tasks) or mixes of the joint and the Cartesian space (in object manipulation environments). The size of generated actions also varies in different environments.

We illustrate the specific environments that we use for our experiments in Fig. 3:

*1) Planar reaching:* is a robot with two revolute joints and the end-effector depicted with the green dot. Its task is to reach a target on a 2D plane (depicted with the red dot) without a priori knowledge of its kinematics model.

*2) UR5 reaching:* is composed of the three first joints of a UR5 robot[1]. The goal of the robot in this environment is to explore the 3D space and reach the desired joint angles (shown by the three yellow cubes).

*3) NICOL manipulator:* is a humanoid torso composed of two OpenManipulator-P arms[2] and two R8H hands[3]. The robot is tasked to learn in varieties of object manipulating scenarios: (i) Reach and move a cube to a desired position (marked by the dark red ellipsoid) on the table–object sliding, or (ii) reach and lift an object up to the air–object lifting. Note that the lifting task implicitly requires the robot to learn a robust grasping skill to complete.

## IV. EXPERIMENTS & RESULTS

We evaluate our proposed method in three experiments. We first train the agent with combined intrinsic and extrinsic rewards, then we investigate the robustness to noise, and, finally, we investigate the transfer learning performance.

### A. Learning performance with combined reward

To investigate how combined intrinsic and extrinsic rewards affect the learning performance, the robot receives two types of reward for its actions: A sparse extrinsic reward

---

[1] https://www.universal-robots.com/de/produkte/ur5-roboter/

[2] http://www.robotis.us/openmanipulator-p/

[3] https://www.seedrobotics.com/rh8d-adult-robot-hand

(a) Planar reaching

(b) UR5 reaching

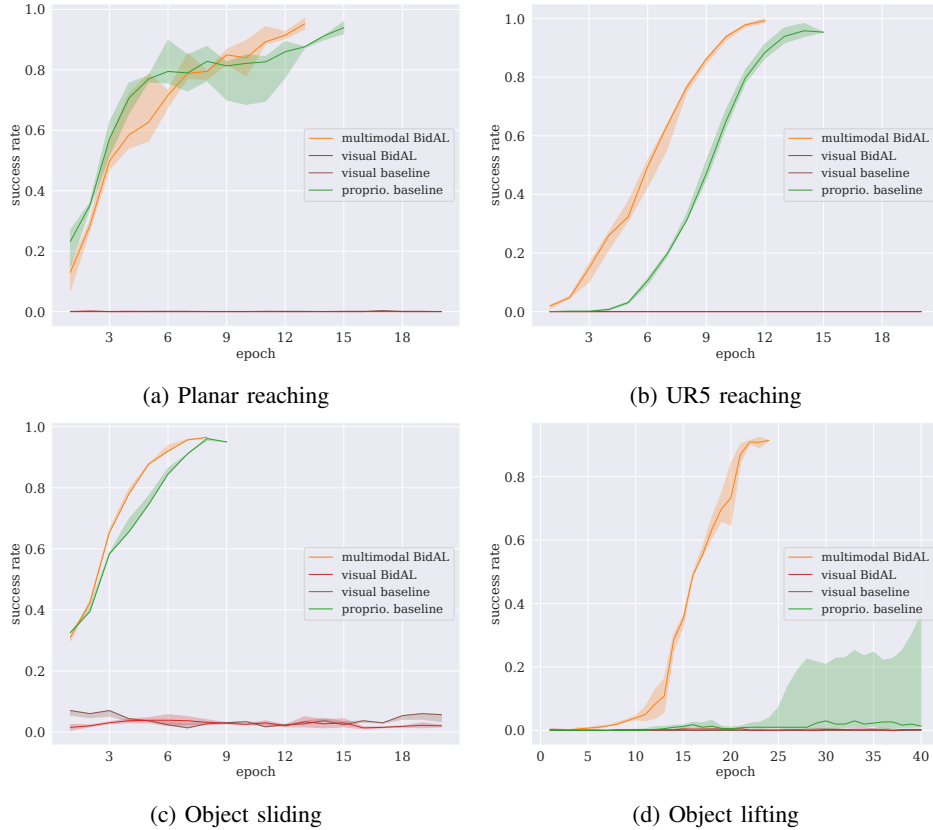(c) Object sliding

(d) Object lifting

Fig. 4: Learning performance in different environment. Note that we use the term "proprioception" for both joint and Cartesian measurements. See Section III-D for more details
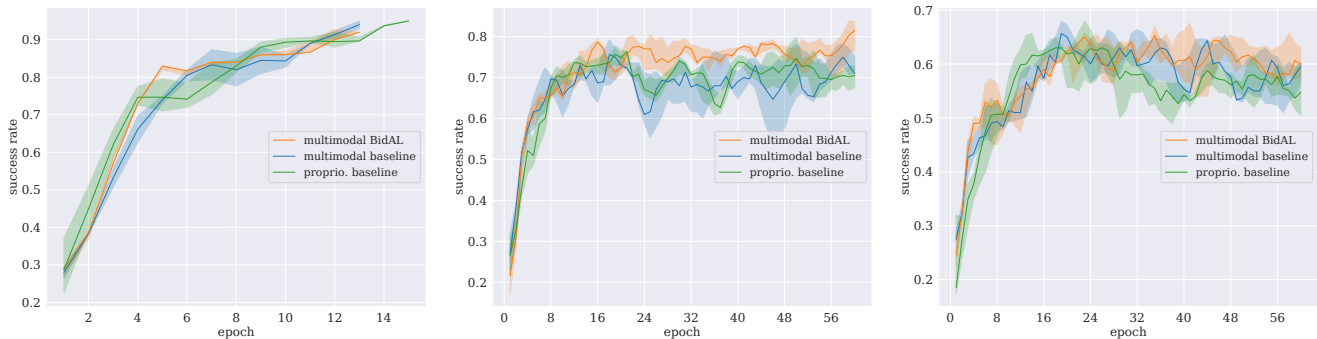


Fig. 5: Learning performance in different noise conditions for the object sliding tasks: 5%–Left, 10%–Center and 15%–Right

from the environment for achieving the task and an intrinsic reward computed from the prediction error.

Herein, we investigate four configurations: As baselines, we use DDPG+HER with proprioceptive input (proprio. baseline) and visual input (visual baseline). We compare these baselines with our proposed algorithm, where we combine the bidirectional associative learning (BidAL) method with the visual input (visual BidAL) and multisensory integration (multimodal BidAL). For these configurations, we perform training in all simulated environments presented in Section III-D.

The results in Fig. 4 show that our proposed algorithm learns all tasks efficiently, while the baselines with only visual input fail to learn any task. The baselines with

proprioceptive input perform significantly better than those with visual input. Overall, our approach outperforms all baselines, especially in the object lifting task. This result provides direct evidence for our core hypothesis that the bidirectional associations improve the learning performance. The results suggest further that minimizing the prediction loss is beneficial for learning the multisensory representation, which is a prerequisite for the robots to learn the main desired task.

### B. Robustness to noise

We further investigate the role of proprioception in the multimodal setting by adding observational noise to the end-effector pose and the object position. This experiment

focuses on the task of object sliding (see Fig. 3c), and aims to simulate that the robot's end-effector and object pose cannot be measured directly in the real world. Normally, these information is obtained through additional estimation processes, mostly from visual input. These estimations may contribute to noise or inaccuracy in the general observation. We realize the noisy observations as follows:

$$s_{noisy} = s + \kappa \cdot n \tag{9}$$

where $\kappa$ denotes the noise coefficient, $n$ denotes noise and $n \sim \mathcal{N}(0, \sigma)$, $\sigma$ denotes the range between $75^{th}$ and $25^{th}$ percentile of the continuous history of $s$.

We perform the training with the combined intrinsic and extrinsic rewards at three different values of the noise coefficient, namely 5%, 10% and 15%. Fig. 5 illustrates that a high noise coefficient affects the performance significantly in all experiments. However, our BidAL approach is more robust to noise, performing significantly more stable than the baseline as we compare the mean and standard deviation of the success rate over the last $N$ epochs[4] over 5 training runs (see Table II).

| Noise ($\kappa$) | multi. BidAL | multi. baseline | proprio. baseline |
|---|---|---|---|
| 5% | $0.874 \pm 0.0127$ | $0.886 \pm 0.0143$ | $0.882 \pm 0.0191$ |
| 10% | $0.764 \pm 0.0222$ | $0.696 \pm 0.0542$ | $0.699 \pm 0.0519$ |
| 15% | $0.615 \pm 0.0320$ | $0.590 \pm 0.0410$ | $0.554 \pm 0.0358$ |

TABLE II: Evaluation of robustness to noise.

For the mean success rates, we see that a drop in the success rate when comparing 5% noise with 15% noise: by a factor of 0.70 and 0,66 for the multimodal BidAL and multimodal baseline, respectively. The higher factor of the BidAL approach illustrates a slight increase in the robustness to noise.

To investigate the stability of the learning under noisy conditions we consider the mean standard deviation of the success rate. With 10% noise, the mean standard deviation of the multimodal BidAL approach is 0.022, compared to a significantly larger value of 0.054 for the multimodal baseline. Hence, the mean standard deviation with the BidAL approach is less then half (factor of 0.4) of the multimodal baseline. In the case of 15% noise, these numbers are 0.032 and 0.041 respectively: BidAL mean standard deviation is around 3/4 (factor of 0.78) of the baseline.

Furthermore, both metrics favor our multimodal BidAL over the proprioceptive baseline.

### C. Transfer learning from a simple to a complex skill

This experiment addresses our hypothesis that previously learned peripersonal space and body schema representations, encoded in the forward and inverse model, foster the learning of new tasks. We propose to pre-train these representations first with a simple skill, and then re-use the representations later to learn a novel skill. Therefore, we first train the NICOL robot to slide objects (see Fig. 3c), and then continue
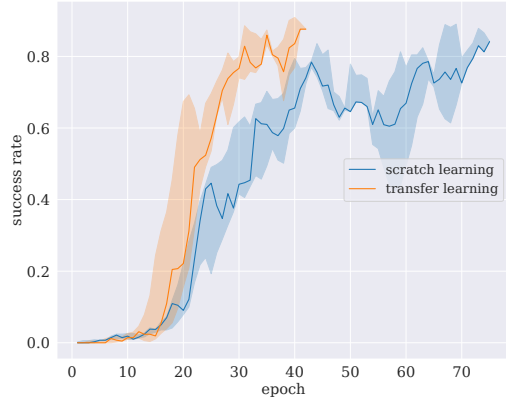


Fig. 6: The robot learns the task of object lifting from multimodal input, with and without transfer learning

learning the more complex task of lifting an object (see Fig. 3d).

To obtain more meaningful results, we increase the difficulty of the lifting task by reducing proprioceptive inputs of the robot that are important for grasping. Specifically, we remove the object orientation, the object velocity, and the relative distance between object and end-effector. This reduced input explains the poorer performance of lifting task in this section than in Section IV-A when the robot learns from scratch. The results in Fig. 6 clearly show that through pre-training the robot with the sliding task, it can more quickly learn to complete the lifting task.

### V. CONCLUSIONS

We have introduced a new goal-directed reinforcement learning framework that builds on multisensory self-representations, realized as bidirectional action-effect associations. Our approach aims to allow agents to exploit relevant sensory input and interact with the environment in a continuous manner. We evaluate our approach in different continuous environments. Our three experiments address our three main hypotheses. We show that intrinsic predictability-driven rewards enable the learning of bidirectional action-effect associations, and we show that these associations implement body schema and peripersonal space representations that make the robotic action-selection robust to noise. Finally, we demonstrate that our approach is beneficial for transfer learning.

Future work involves optimization of hyperparamters and the combination with a planning method. The planning will be based on the forward model that we have learned. We hypothesize that integrating planning with a reinforcement learning policy will further improve the learning performance, as demonstrated by related approaches [51]–[53].

### ACKNOWLEDGMENT

---

[4]$N = 30$ for 10-15% noise and $N = 5$ for 5% noise.

## References

[1] H. Head and G. Holmes, "Sensory disturbances from cerebral lesions," *Brain*, vol. 34, no. 2-3, pp. 102–254, Nov. 1911.

[2] F. de Vignemont, "Body schema and body image-Pros and cons," *Neuropsychologia*, 2010.

[3] J. Cléry *et al.*, "Neuronal bases of peripersonal and extrapersonal spaces, their plasticity and their dynamics: Knowns and unknowns," *Neuropsychologia*, vol. 70, pp. 313–326, Apr. 2015.

[4] A. Serino, "Peripersonal space (PPS) as a multisensory interface between the individual and the environment, defining the space of the self," *Neuroscience and Biobehavioral Reviews*, vol. 99, no. August 2018, pp. 138–159, 2019.

[5] P. Rochat and R. Morgan, "Spatial Determinants in the Perception of Self-Produced Leg Movements by 3- to 5-Month-Old Infants," *Developmental Psychology*, vol. 31, no. 4, pp. 626–636, 1995.

[6] A. J. Bremner *et al.*, *Multisensory Development*, eng. Oxford: Oxford University Press, 2012, p. 392.

[7] D. Corbetta *et al.*, "The Embodied Origins of Infant Reaching: Implications for the Emergence of Eye-Hand Coordination," *Kinesiology Review*, vol. 7, no. 1, pp. 10–17, 2018.

[8] B. Elsner and B. Hommel, "Effect anticipation and action control," *Journal of Experimental Psychology: Human Perception and Performance*, vol. 27, no. 1, pp. 229–240, 2001.

[9] S. A. Verschoor and B. Hommel, "Self-by-doing: The role of action for self-acquisition," *Social Cognition*, vol. 35, no. 2, pp. 127–145, 2017.

[10] A. Maravita and A. Iriki, "Tools for the body (schema)," *Trends in Cognitive Sciences*, vol. 8, no. 2, pp. 79–86, 2004.

[11] D. Pathak *et al.*, "Curiosity-driven Exploration by Self-supervised Prediction," in *International Conference on Machine Learning (ICML)*, May 2017.

[12] M. Hoffmann *et al.*, "Body Schema in Robotics: A Review," *IEEE Transactions on Autonomous Mental Development*, vol. 2, no. 4, pp. 304–324, Dec. 2010.

[13] V. Gallese and C. Sinigaglia, "The bodily self as power for action," *Neuropsychologia*, vol. 48, no. 3, pp. 746–755, 2010.

[14] A. Roncone *et al.*, "Automatic kinematic chain calibration using artificial skin: Self-touch in the iCub humanoid robot," in *IEEE International Conference on Robotics and Automation*, 2014, pp. 2305–2312.

[15] P. Vicente *et al.*, "Online body schema adaptation based on internal mental simulation and multisensory feedback," *Frontiers Robotics AI*, vol. 3, no. MAR, 2016.

[16] M. Hoffmann *et al.*, "Robotic homunculus: Learning of artificial skin representation in a humanoid robot motivated by primary somatosensory cortex," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 10, no. 2, pp. 163–176, 2018.

[17] G. Schillaci *et al.*, "Online learning of visuo-motor coordination in a humanoid robot. A biologically inspired model," in *ICDL-EPIROB 2014*, IEEE, 2014, pp. 130–136.

[18] L. P. Wijesinghe *et al.*, "Robot end effector tracking using predictive multisensory integration," *Frontiers in Neurorobotics*, vol. 12, no. October, pp. 1–16, 2018.

[19] P. D. Nguyen *et al.*, "Transferring Visuomotor Learning from Simulation to the Real World for Robotics Manipulation Tasks," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, IEEE, Oct. 2018, pp. 6667–6674.

[20] E. Magosso *et al.*, "Visuotactile representation of peripersonal space: A neural network study," *Neural computation*, vol. 22, no. 1, pp. 190–243, 2010.

[21] A. Serino *et al.*, "Extending peripersonal space representation without tool-use: Evidence from a combined behavioral-computational approach," *Frontiers in Behavioral Neuroscience*, vol. 9, Feb. 2015.

[22] Z. Straka and M. Hoffmann, "Learning a peripersonal space representation as a visuo-tactile prediction task," in *International Conference on Artificial Neural Networks*, Springer, 2017, pp. 101–109.

[23] E. Chinellato *et al.*, "Implicit sensorimotor mapping of the peripersonal space by gazing and reaching," *IEEE Transactions on Autonomous Mental Development*, vol. 3, no. 1, pp. 43–53, 2011.

[24] J. Juett and B. Kuipers, "Learning to grasp by extending the peri-personal space graph," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems*, IEEE, 2018, pp. 8695–8700.

[25] A. Roncone *et al.*, "Peripersonal space and margin of safety around the body: Learning visuo-tactile associations in a humanoid robot with artificial skin," *PLoS ONE*, vol. 11, no. 10, pp. 1–32, 2016.

[26] D. Nguyen *et al.*, "Compact Real-time Avoidance on a Humanoid Robot for Human-robot Interaction," in *ACM/IEEE International Conference on Human-Robot Interaction*, 2018.

[27] P. D. H. Nguyen *et al.*, "Reaching development through visuo-proprioceptive-tactile integration on a humanoid robot - a deep learning approach," in *ICDL-EpiRob 2019*, Aug. 2019, pp. 163–170.

[28] G. Pugach *et al.*, "Brain-inspired coding of robot body schema through visuo-motor integration of touched events," *Frontiers in Neurorobotics*, vol. 13, no. March, 2019.

[29] M. Zambelli *et al.*, "Multimodal representation models for prediction and control from partial information," *Robotics and Autonomous Systems*, vol. 123, 2020.

[30] J. L. Copete *et al.*, "Motor development facilitates the prediction of others' actions through sensorimotor predictive learning," in *ICDL-EpiRob 2016*, IEEE, 2017, pp. 223–229.

[31] J. Hwang *et al.*, "Dealing with Large-Scale Spatio-Temporal Patterns in Imitative Interaction between a Robot and a Human by Using the Predictive Coding Framework," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 50, no. 5, pp. 1918–1931, 2020.

[32] C. Lang *et al.*, "A deep convolutional neural network model for sense of agency and object permanence in robots," in *ICDL-EpiRob 2018*, IEEE, 2018, pp. 257–262.

[33] M. Watter *et al.*, "Embed to control: A locally linear latent dynamics model for control from raw images," in *Advances in neural information processing systems*, 2015, pp. 2746–2754.

[34] A. Byravan *et al.*, "SE3-Pose-Nets: Structured deep dynamics models for visuomotor control," in *IEEE International Conference on Robotics and Automation*, 2018, pp. 3339–3346.

[35] P. Agrawal *et al.*, "Learning to poke by poking: Experiential learning of intuitive physics," in *Advances in Neural Information Processing Systems*, 2016, pp. 5074–5082.

[36] J. C. Park *et al.*, "Learning for Goal-Directed Actions Using RNNPB: Developmental Change of 'What to Imitate'," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 10, no. 3, pp. 545–556, 2018.

[37] P. Lanillos *et al.*, "Yielding Self-Perception in Robots Through Sensorimotor Contingencies," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 9, no. 2, pp. 100–112, 2017.

[38] N. A. Hinz *et al.*, "Drifting perceptual patterns suggest prediction errors fusion rather than hypothesis selection: Replicating the rubber-hand illusion on a robot," in *ICDL-EpiRob 2018*, IEEE, 2018, pp. 125–132.

[39] P. Lanillos *et al.*, "Robot self/other distinction: active inference meets neural networks learning in a mirror," in *24th European Conference on Artificial Intelligence (ECAI 2020)*, 2020.

[40] P. Oudeyer *et al.*, "Intrinsic Motivation Systems for Autonomous Mental Development," *IEEE Transactions on Evolutionary Computation*, vol. 11, no. 2, pp. 265–286, Apr. 2007.

[41] Y. Burda *et al.*, "Large-Scale Study of Curiosity-Driven Learning," in *International Conference on Learning Representations (ICLR)*, 2019.

[42] N. Dilokthanakul *et al.*, "Feature Control as Intrinsic Motivation for Hierarchical Reinforcement Learning," *IEEE Transactions on Neural Networks and Learning Systems*, 2019.

[43] D. Pathak *et al.*, "Self-Supervised Exploration via Disagreement," in *International Conference on Machine Learning (ICML)*, 2019.

[44] F. Röder *et al.*, "Curious Hierarchical Actor-Critic Reinforcement Learning," in *29th International Conference on Artificial Neural Networks (ICANN2020)*, 2020.

[45] T. P. Lillicrap *et al.*, "Continuous control with deep reinforcement learning," in *International Conference on Learning Representations (ICLR)*, 2015.

[46] V. Mnih *et al.*, "Asynchronous methods for deep reinforcement learning," in *International Conference on Machine Learning (ICML)*, vol. 4, PMLR, 2016, pp. 2850–2869.

[47] T. Schaul *et al.*, "Universal Value Function Approximators," in *International Conference on Machine Learning (ICML)*, 2015, pp. 1312–1320.

[48] M. Andrychowicz *et al.*, "Hindsight Experience Replay," in *Conference on Neural Information Processing Systems (NIPS)*, 2017, pp. 5048–5058.

[49] D. Silver *et al.*, "Deterministic Policy Gradient Algorithms," in *International Conference on Machine Learning (ICML)*, 2014, pp. 1–9.

[50] E. Todorov *et al.*, "MuJoCo: A physics engine for model-based control," in *IEEE International Conference on Intelligent Robots and Systems*, 2012, pp. 5026–5033.

[51] A. Srinivas *et al.*, "Universal planning networks," in *International Conference on Machine Learning (ICML)*, vol. 11, 2018.

[52] A. H. Qureshi *et al.*, "Motion planning networks," in *IEEE International Conference on Robotics and Automation*, IEEE, 2019, pp. 2118–2124.

[53] M. Eppe *et al.*, "From Semantics to Execution: Integrating Action Planning With Reinforcement Learning for Robotic Causal Problem-Solving," *Frontiers in Robotics and AI*, vol. 6, p. 123, 2019.