

# Imitation and mirror systems in robots through Deep Modality Blending Networks

M. Yunus Seker<sup>a,\*</sup>, Alper Ahmetoglu<sup>a</sup>, Yukie Nagai<sup>d</sup>, Minoru Asada<sup>b</sup>, Erhan Oztop<sup>b,c</sup>, Emre Ugur<sup>a</sup>

<sup>a</sup> Bogazici University, Bebek, Istanbul, 34342, Turkey

<sup>b</sup> Osaka University, Suita, Osaka, Japan

<sup>c</sup> Ozyegin University, Istanbul, Turkey

<sup>d</sup> The University of Tokyo, Bunkyo-ku, Tokyo, Japan

## ARTICLE INFO

### Article history:

Received 24 June 2021

Received in revised form 29 September 2021

Accepted 4 November 2021

Available online 16 November 2021

### Keywords:

Robot learning

Imitation learning

Representation learning

Multimodal learning

## ABSTRACT

Learning to interact with the environment not only empowers the agent with manipulation capability but also generates information to facilitate building of action understanding and imitation capabilities. This seems to be a strategy adopted by biological systems, in particular primates, as evidenced by the existence of mirror neurons that seem to be involved in multi-modal action understanding. How to benefit from the interaction experience of the robots to enable understanding actions and goals of other agents is still a challenging question. In this study, we propose a novel method, deep modality blending networks (DMBN), that creates a common latent space from multi-modal experience of a robot by blending multi-modal signals with a stochastic weighting mechanism. We show for the first time that deep learning, when combined with a novel modality blending scheme, can facilitate action recognition and produce structures to sustain anatomical and effect-based imitation capabilities. Our proposed system, which is based on conditional neural processes, can be conditioned on any desired sensory/motor value at any time step, and can generate a complete multi-modal trajectory consistent with the desired conditioning in one-shot by querying the network for all the sampled time points in parallel avoiding the accumulation of prediction errors. Based on simulation experiments with an arm-gripper robot and an RGB camera, we showed that DMBN could make accurate predictions about any missing modality (camera or joint angles) given the available ones outperforming recent multimodal variational autoencoder models in terms of long-horizon high-dimensional trajectory predictions. We further showed that given desired images from different perspectives, i.e. images generated by the observation of other robots placed on different sides of the table, our system could generate image and joint angle sequences that correspond to either anatomical or effect-based imitation behavior. To achieve this mirror-like behavior, our system does not perform a pixel-based template matching but rather benefits from and relies on the common latent space constructed by using both joint and image modalities, as shown by additional experiments. Moreover, we showed that mirror learning (in our system) does not only depend on visual experience and cannot be achieved without proprioceptive experience. Our experiments showed that out of ten training scenarios with different initial configurations, the proposed DMBN model could achieve mirror learning in all of the cases where the model that only uses visual information failed in half of them. Overall, the proposed DMBN architecture not only serves as a computational model for sustaining mirror neuron-like capabilities, but also stands as a powerful machine learning architecture for high-dimensional multi-modal temporal data with robust retrieval capabilities operating with partial information in one or multiple modalities.

© 2022 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## 1. Introduction

With appropriate and sufficient amount of data, a range of sensorimotor learning tasks encountered by robots and biological systems can be solved by deep learning. However, unlike the abundance of data for image recognition and language

\* Corresponding author.

E-mail address: [yunus.seker1@boun.edu.tr](mailto:yunus.seker1@boun.edu.tr) (M.Y. Seker).

modeling, robots and biological systems often need to harvest data themselves by either using self-exploration based learning or by observing the relevant behaviors of other agents. These two alternatives are studied in robotics and machine learning under the general titles of Reinforcement Learning (RL) (Sutton & Barto, 2018) and Learning from Demonstration (LfD) (Argall, Chernova, Veloso, & Browning, 2009). Although the use of self-observation during self-executed actions is common for forming a reward signal in RL, how to benefit the agent in a cognitive developmental sense is not well addressed. For example, for recognizing actions of others or forming a general imitation capacity. Learning to interact with the environment not only empowers the agent with manipulation capability but also generates information to facilitate the building of action understanding and imitation capabilities. This seems to be a strategy adopted by biological systems, in particular primates, as evidenced by the existence of mirror neurons (Di Pellegrino, Fadiga, Fogassi, Gallese, & Rizzolatti, 1992; Rizzolatti, Fadiga, Gallese, & Fogassi, 1996) in the ventral premotor cortex of those animals, which encode actions in a multi-modal fashion (Kohler, Keysers, Umiltà, Fogassi, Gallese, & Rizzolatti, 2002). For example, there are mirror neurons that become active when the animal breaks a peanut, observes an experimenter do the same act or hears the sound of peanut cracking (Keysers, Kohler, Umiltà, Nanetti, Fogassi, & Gallese, 2003). With such a system, sensed actions are mapped to one's own motor representation; and thus can bootstrap imitation, by for example, understanding the parts of an observed act in terms of the existing 'action vocabulary' of the animal, which can be reproduced in sequence yielding novel action imitation capability. Although, it is not clear whether mirror neurons play a role in imitation, as their exact function and mechanism are far from clear, computational modeling may help produce insights towards understanding them (Oztop, Kawato, & Arbib, 2013). Therefore, from a scientific and also technological point of view, it is desirable to develop a neural multi-modal action representation system that can learn/store actions and recall them from partial information that might be transformed as in the case of action observation from different perspectives. In fact, there exist a range of computational models related to mirror neurons and their function in the literature (Bonaiuto & Arbib, 2010; Bonaiuto, Rosta, & Arbib, 2007; Copete, Nagai, & Asada, 2016a; Demiris & Johnson, 2003; Oztop & Arbib, 2002; Tani, Ito, & Sugita, 2004) that have leveraged our understanding by creating hypotheses to be tested. Now, the time is ripe for a less constrained, end-to-end and more powerful multi-modal action representation mechanism for obtaining better insights. In particular, the existing multi-modal action representation schemes based on self-observation either fall short of providing robust recognition and imitation capability or rely on feature engineering.

In this study, we improve the state of the art in multi-modal action representation by showing for the first time that deep learning, when combined with a novel modality blending scheme, can facilitate feature-engineering-free action recognition and basic imitation capabilities under perspective changes with only partial information. Moreover, the modality blending scheme produces latent representations that can sustain both anatomical and effect-based imitation capabilities. We call the developed multi-modal action representation architecture as a Deep Modality Blending Network (DMBN).

DMBN connects multiple modalities by blending them as random mixtures of modality-specific latent representations to form a common latent representation for seamless transfer from one modality to another (see Fig. 1). The DMBN architecture follows an encoder–decoder structure where each modality is summarized by its corresponding encoder network, processing the sensorimotor data into a compact latent representation. While

learning, not only these latent representations are formed but they are blended together into a common representation through stochastic mixture weights. After learning, using the common representation, each decoder network can predict the corresponding modality for an arbitrary desired time step, effectively generating outcome predictions as temporal sequences for all the modalities. In this sense, the common latent layer in our network encodes representation of the complete multi-modal trajectories rather than encoding modalities in particular time steps. This feature sets our system apart from its competitors (Copete et al., 2016a; Zambelli, Cully, & Demiris, 2020) and give it a big advantage. To be concrete, our system can be conditioned on any desired sensory/motor value at any time step, and can generate a complete multi-modal trajectory consistent with the desired conditioning in one-shot by querying the network for all the sampled time points in parallel. DMBNs make temporal predictions independently for each query point in one-shot without requiring feeding back of the output as input. This one-shot full trajectory decoding ability makes our system very accurate as it does not suffer from the error accumulation faced by systems that need to chain next-state predictions in order to generate full trajectories.

To demonstrate the efficacy of the proposed DMBN architecture, we implemented it in a simulated manipulation setup. In this setup, an object was placed in the middle of a table, and an arm-gripper robot was set to execute grasp and push actions on it with different approach directions. The robot observed the consequences of its actions using an RGB camera from a fixed perspective, and learned the generated multi-modal sensory (visual and proprioceptive) signals as sensory trajectory distributions through the proposed DMBN architecture. After learning,

- Given desired images at any time point (such as images of objects lifted or pushed away), our system can find the joint trajectories that are required to generate changes in the environment to observe these images;
- Given joint angles at any time point(s), our system can generate the sequence of images that are expected to be observed during the execution of the action that is consistent with given angles;
- Given desired images from different perspectives, i.e. images generated by the observation of other robots placed on different sides of the table, our system can generate image and joint angle sequences that correspond to valid actions of the robot;
- Those valid actions, intriguingly correspond to either anatomical or effect based imitation behavior.

To clarify the last bullet above it would be useful to consider an example behavior observed in our simulations. Given an image that shows the snapshot of another robot on the other side of the table pulling the object to itself, our system can generate the sequence of images where its gripper pulls the object towards itself (anatomical imitation behavior) or pushes the object towards the other side of the table (effect based imitation or goal-emulation behavior) depending on the visual cues available to the robot. In our analysis, we show that the prediction capability of the proposed DMBN system does not simply perform a pixel-based template matching but rather benefits from and relies on the common latent space constructed by using both joint and image modalities. In addition to other interesting results, our experiments clearly show that our system outperforms a recent multimodal variational autoencoder model (Zambelli et al., 2020) in reconstructing long-horizon high-dimensional trajectories.

The outline of this paper is as follows: in Section 2, we review the related work, in particular, LfD systems as DMBN builds upon one such system and the competing multi-modal action representations. In Section 3, we describe our proposed method

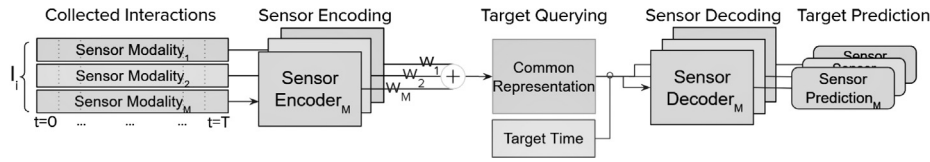


Fig. 1. General architecture of a Deep Modality Blending Network.

in detail. We explain our experiment setup in Section 4 and give experimental results in Section 5. Finally, we give a conclusion in Section 6.

## 2. Related work

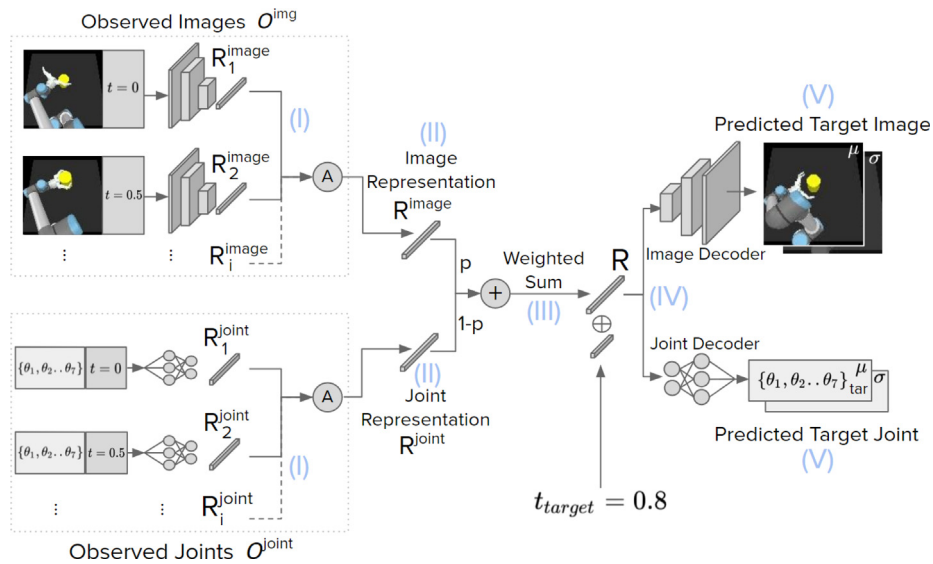
Imitation learning, or learning from demonstration (LfD) (Ar-gall et al., 2009), has been a popular research topic in robotic learning (Asfour, Azad, Gyrfas, & Dillmann, 2008; Ben Amor, Kroemer, Hillenbrand, Neumann, & Peters, 2012; Mühlig, Gienger, & Steil, 2012; Paraschos, Daniel, Peters, & Neumann, 2018; Pastor, Hoffmann, Asfour, & Schaal, 2009; Pastor, Righetti, Kalakrishnan, & Schaal, 2011). Various LfD methods have been proposed based on dynamic systems and statistical modeling (Calinon, 2016; Huang, Rozo, Silvério, & Caldwell, 2019; Schaal, 2006; Zhou & Asfour, 2017), where the parameters in the environment can be learned with Locally Weighted Regression (Atkeson, Moore, & Schaal, 1997; Kramberger, Gams, Nemeč, Chrysostomou, Mad-sen, & Ude, 2017; Ude, Gams, Asfour, & Morimoto, 2010) and Locally Weighted Projection Regression (Vijayakumar & Schaal, 2000). Gaussian Mixture Models (Calinon, Evrard, Gribovskaya, Billard, & Kheddar, 2009; Pervez & Lee, 2018) and Hidden Markov Models (Chu et al., 2013; Girgin & Ugur, 2018; Lee & Ott, 2011; Ugur & Girgin, 2020) are also frequently used to learn the motion distributions from multiple demonstrations. More recently, deep neural networks also started to be used in imitation learning to learn movement primitives from complex high-dimensional data (Droniou, Ivaldi, & Sigaud, 2015; Pahič, Gams, Ude, & Mo-rimoto, 2018; Pervez, Mao, & Lee, 2017; Xie, Chowdhury, De Pao-lis Kaluza, Zhao, Wong, & Yu, 2020). In our earlier work, we proposed Conditional Neural Movement Primitives (CNMPs) (Akbulut, Bozdoğan, Tekden, & Ugur, 2021; Akbulut, Oztop, Seker, Xue, Tekden, & Ugur, 2020; Seker, Imre, Piater, & Ugur, 2019) as an end-to-end deep LfD architecture that can learn temporal sensorimotor distributions of complex manipulation skills. Based on Conditional Neural Processes (Garnelo et al., 2018), CNMPs are deep learning from demonstration frameworks that use stochastic observation sampling and query prediction to learn complex temporal data. CNMPs are able to learn and generalize high-dimensional data due to their deep encoder–decoder architecture. The stochastic observation sampling used in the training process makes it possible to learn from a few examples. The most distinctive feature of CNMPs compared to the other approaches is the observation-query mechanism that allows the framework to collect observations and query predictions on any time-steps. Contrary to the other methods using recurrent models which are bound to their own outputs for the future predictions, CNMPs can make predictions for any time-steps independently before or after the given observations. The DMBN architecture developed in the current study builds upon CNMPs by introducing a novel mechanism for modality blending to learn a common latent representation that allows cross-modal temporal prediction with partial information.

Several works studied the emergence of the mirror neuron system (MNS) in the context of multi-modal sensor fusion. Nagai, Kawai, and Asada (2011) proposed a computational model for the early development of the MNS. In this model, the robot cannot make self-other discrimination in the early stages due to

the immature visual system. As the visual system develops, the robot starts to discriminate between itself and others, yet, still retains information regarding early experiences, producing the MNS as a by-product. Noda, Arie, Suga, and Ogata (2014) used time-delay neural networks (Waibel, Hanazawa, Hinton, Shikano, & Lang, 1989) as autoencoders to fuse multiple modalities and reconstruct the missing ones given others.

Copete, Nagai, and Asada (2016b) also used a similar autoen-coder architecture in a predictive learning context so to imagine the action of others. Jung, Matsumoto, and Tani (2019) proposed a top-down visual attention system to address the long-term visual prediction problem. In this system, the visual stream is divided into dorsal and ventral streams to decompose the difficulty of the problem into two sub-problems. These two streams are then merged for the visual prediction with the help of an external visuospatial memory which holds long-term visuospatial information. On the other hand, we provide a more holistic approach where there are only different submodules for different modalities. Our experiments show that DMBNs can output very accurate visual signals conditioned only on a single visual frame without any memory module. Among these studies, the learning problems considered in the work of Zambelli et al. (2020) is well-aligned with our study. They proposed a multimodal variational autoencoder (MVAE) (Suzuki, Nakayama, & Matsuo, 2016; Wu & Goodman, 2018) to fuse the sensorimotor information of an iCub humanoid robot for prediction and control. They showed that by training MVAE as a denoising autoencoder (Vincent, Larochelle, Bengio, & Manzagol, 2008), MVAE can predict the future sensorimotor states, reconstruct the missing modalities, and imitate based on human action observation. As MVAE is not a recurrent architecture, the temporal information should be explicitly stated in the input. To be concrete, in the training phase, the sensorimotor information at time  $t$  and  $t + 1$  were combined and given as input to the MVAE for reconstruction. Here, some sensorimotor information at time  $t + 1$  was randomly masked with  $-2$  (as in a denoising autoencoder) to train the network to reconstruct the future time step even if it was partially missing.

In the testing phase for future state predictions, states at  $t + 1$  were filled with mask values  $-2$ . Further steps could be predicted by feeding the output of the MVAE to the input. However, the error at one step cascades in the feedback loop as in RNNs. Therefore, the prediction power decreases as the trajectory horizon increases. This is not the case in our proposed model as DMBNs make temporal predictions in one-shot without requiring feeding back of the output as input. To concretely state, our work differs from the previous works in terms of modality fusion strategy and architecture: (1) we force the formation of a more common and robust representation space by taking stochastic mixtures of modalities during training, and (2) we learn individual modalities and their mixture as long range dependencies via CNPs (Garnelo et al., 2018), which allow arbitrary future and past temporal predictions. These key differences yield not only a more robust and better performing multi-modal action representation system, but also give rise to interesting generalization abilities as shown with the experiments presented in the Results section.



**Fig. 2.** Proposed framework for given *visual* and *joint* modalities. Image and joint observations are turned into their latent representations separately to be used to predict the image and joint positions given at another target time step.

### 3. Method

In this work, we propose Deep Modality Blending Networks<sup>1</sup> (DMBNs), that can learn and produce sensorimotor signals by forming and exploiting multi-modal representations acquired in a latent space. Assume  $M = \{\text{visual, proprioception, sound, haptic ...}\}$  corresponds to sensorimotor signals from multiple modalities collected by an agent through self-observation. The agent interacts with the environment using a variety of actions to leverage the information produced by the embodied interaction of the agent with the environment. In the current implementation, the action and action parameters are sampled from a predefined action repertoire. During every interaction, the sensorimotor values are recorded at each time step. The multi-sensorimotor interaction data set is defined as  $I$ , and the  $i$ th interaction is described as  $I_i = \{(t, S_t^M)\}_{t=0}^T$ , where  $t$  is time and  $S_t^M$  is the sensorimotor state collection for the given time step.  $S_t^M$  consists of multiple sensorimotor data,  $S_t^M = [S_t^{visual}, S_t^{joint}, S_t^{sound}, S_t^{haptic}, \dots]$ , where each member holds the corresponding state values of the sensorimotor modalities for the time step  $t$ . Fig. 2 shows the architecture of our model where the modalities in the system correspond to the *visual* and *proprioceptive* domains. These two domains are chosen specifically in order to show that our system can learn in an end-to-end fashion with both high (image) and low (joint) dimensional data and make more accurate target predictions on a long horizon compared to the sequential prediction models. In theory, all types of sensorimotor data can be included in the system with our formulation.

The aim of DMBN is to predict a conditional output distribution for a target query given a desired set of observation samples. At the beginning of each training iteration, an interaction  $I_d$  is selected randomly from the data set  $I$ . From this selected interaction,  $n$  data points of  $(t, S_t^M)$ , are randomly sampled as observations. Here,  $n$  is a changing number for each training iteration that is bounded by  $[1, obs_{max}]$  where  $obs_{max}$  is a hyper-parameter that decides the maximum number of sampled observations in the training. We define this sampled observation set as  $O^M = \{(t_i, S_{t_i}^M)\}_{i=1}^{obs_{max}}$  where  $(t_i, S_{t_i}^M) \in I_d$ . On the left side of Fig. 2.(1), example sampled observations  $O^{image}$  and  $O^{joint}$

are shown for the image and the joint domains. Besides  $O^M$ , a target tuple  $(t_{target}, S_{t_{target}}^M)$  is also sampled from the same selected interaction  $I_d$ . The purpose of a training iteration is to learn distributions on  $t_{target}$  for all modalities in the system, based on the observation set  $O^M$ .

Our aim is to merge the observations of all the modality signals in a single latent space to allow information sharing for a higher quality prediction. In order to achieve this, the observations of each modality,  $O^m$ , are first transformed into their latent representations  $R_i^m$ . For every modality  $m$  and every observation, latent states are calculated by the following equation:

$$R_i^m = E^m((t_i, S_{t_i}^m) | \theta^m) \quad (t_i, S_{t_i}^m) \in O^m, m \in M \quad (1)$$

where  $E^m$  is a deep encoder for the modality  $m$  with weights  $\theta^m$ , and  $R_i^m$  is the latent states of its  $i$ th observation. Fig. 2.(1) shows the encoded representations,  $R_i^{image}$  and  $R_i^{joint}$ , for each observation. After generating these representations, an averaged representation of each modality is calculated by:

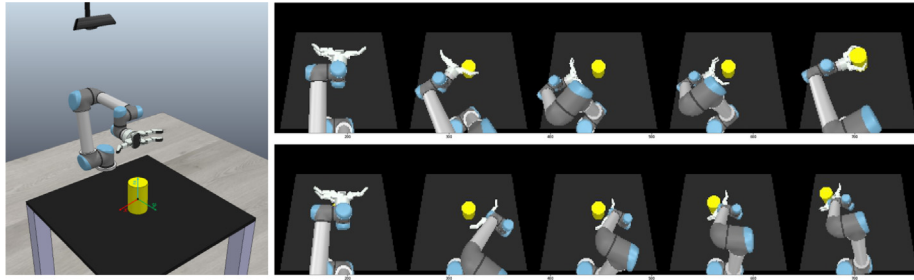
$$R^m = \frac{1}{n} \sum_i R_i^m \quad m \in M \quad (2)$$

where  $n$  is the size of the observations of this training iteration.  $R^{image}$  and  $R^{joint}$  in Fig. 2.(II) hold general knowledge about their modalities, and our aim is to use these representations in a shared latent space to allow information sharing between all modalities. To achieve this, a multi-modal general representation  $R$  that integrates all modalities is constructed by calculating a normalized weighted average:

$$R = \frac{\sum_m p^m R^m w^m}{\sum_m p^m w^m} \quad (3)$$

where  $w^M = [w^{image}, w^{joint}, w^m, \dots]$  is a vector representing the *weight* or *availability* of the individual modalities with  $0 \leq w^M \leq 1$  and  $w^M \neq 0$ , which could be used to model cases where one modality is more reliable than the other. On the other hand, modality blending during training is achieved through the random variables  $0 \leq p^m \leq 1$  that is sampled at every iteration, and obey the constraint  $\sum p^m = 1$ . Note that to avoid  $\sum p^m w^m$  ever becoming zero (See Eq. (3)), we may require  $p^m > 0$ ; but this is not an issue in practice. This follows the same intuition with dropout (Srivastava, Hinton, Krizhevsky, Sutskever, &

<sup>1</sup> Our implementation of DMBN is available at <https://github.com/myunusseker/Deep-Modality-Blending-Networks>.



**Fig. 3.** (Left) Experiment setup with vision sensor, UR5, and the object at the middle of the table. (Right) Example grasping and pushing actions recorded via the vision sensor.

Salakhutdinov, 2014); randomly dropping modalities forces the model to learn compact representations that can compensate for missing information. Fig. 2.(III) shows this process as a two-modality setup where  $w^{image} = w^{joint} = 0.5$  and  $p^{image} = p$  and  $p^{joint} = 1 - p$  where  $p$  is sampled uniformly from  $[0, 1]$ . Note that the dimension of each  $R^m$  should be the same in order to perform summation operation between vectors, so in the first place, all the encoders must be designed to produce the latent states with the same dimensions. Once all observations are merged into one general representation, this information can be used to infer target distributions on  $t_{target}$  for all the modalities as:

$$(\mu_{t_{target}}^m, \sigma_{t_{target}}^m) = Q^m((R, t_{target}) | \phi^m) m \in M \quad (4)$$

where  $Q^m$  is a deep decoder network with weights  $\phi^m$  that produces a distribution that consists of a mean  $\mu_{t_{target}}^m$  and variance  $\sigma_{t_{target}}^m$  for the modality  $m$ . Fig. 2.(IV) shows the decoders,  $Q^{image}$  and  $Q^{joint}$ , and predicted distributions,  $(\mu_{t_{target}}^{image}, \sigma_{t_{target}}^{image})$  and  $(\mu_{t_{target}}^{joint}, \sigma_{t_{target}}^{joint})$ , for two domains. The learning objective of our framework is to construct better distributions according to the given observations as in Garnelo et al. (2018) and Seker et al. (2019), so the loss term is defined as:

$$\mathcal{L} = - \sum_m \log P(S_{t_{target}}^m | \mu_{t_{target}}^m, \sigma_{t_{target}}^m) \quad (5)$$

where  $S_{t_{target}}^m \in S_{t_{target}}^M$  is the target sensorimotor value for modality  $m$  at time  $t_{target}$ .

After training, the system can be requested to make predictions for all the modalities and for all the time steps by fixing  $p^M = 1/M$  and assigning  $O^M$  as novel observations. By observing the sensorimotor state at any time step, any other time point before and after can be queried and predicted using our framework. According to the situation, if a sensorimotor modality does not seem to provide reliable signals, the weight given to that modality can be decreased by configuring availability vector  $w$ . Note that, the system can even predict missing modalities if the corresponding  $w^m$  is set to zero because of a lack of the modality. Our framework can use the shared latent space for multi-modal predictions. This also enables our framework to imitate other agents by observing their actions with, for example, vision and sound, and producing the agent own behavior by predicting the corresponding motor signals.

#### 4. Experiment setup

To demonstrate the capabilities of our system, we designed an experiment where the actions of the robot can be predicted from the visual and proprioceptual observations at the beginning of the movement execution. A simulated environment was built using CoppeliaSim (Rohmer, Singh, & Freese, 2013). The setup

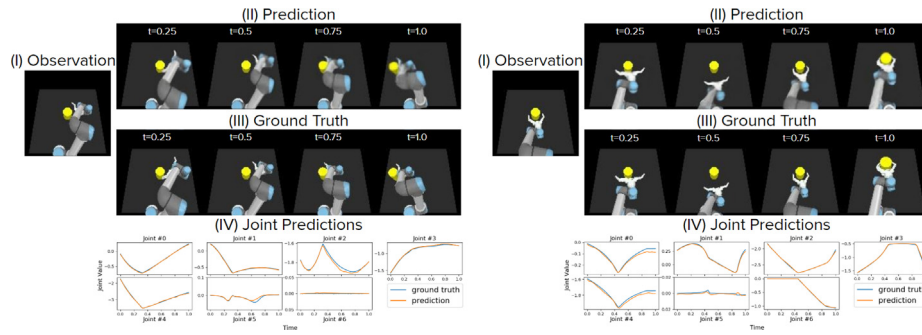
consisted of a UR5 robot equipped with a three-finger gripper, a vision sensor, and an object on a table to be manipulated by the robot (Fig. 3 left). The action repertoire of the robot was composed of parameterized push and grasp actions that allow reaching to the object from all directions, and the data collection protocol for each action execution (interaction) was as follows. At the beginning of each interaction, the robot initialized its wide-open hand at an initial position, and an object appeared in the middle of the table (Fig. 3 right). If the selected action was push, a random pushing angle was sampled and the robot pushed the object from this angle to a predetermined fixed distance of 30 cm while keeping the hand open. If the selected action was grasp, a random grasping angle was sampled and the robot started to close its hand while approaching to the object so as to grasp it and lift it to a fixed height over the table (30 cm). The collected data consisted of two modalities that are *proprioception* and *vision*. The proprioceptive signals were composed of seven joint angles of the robot (6 joints of the UR5 robot and 1 hand opening joint), whereas the visual signals were  $128 \times 128 \times 3$  RGB images. Visual signals were collected via the vision sensor that was placed to the point of view of the robot (see Fig. 3). In the end, 50 successful push and grasp interactions (100 in total) were collected using the simulator. The interactions were separated into train and test sets with 80% and 20% ratios respectively.

#### 5. Experimental results

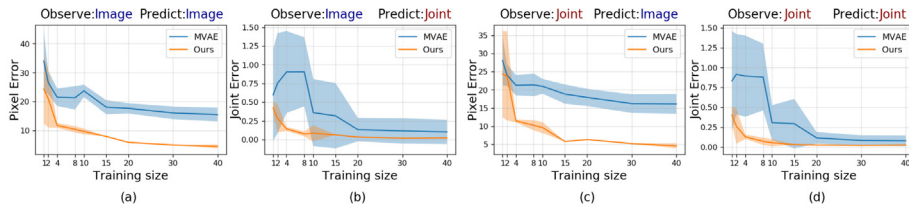
We conducted a set of experiments to test the capabilities of DMBN from different aspects. First, in Section 5.1, we verify the prediction capabilities of DMBN by generating complete image and joint trajectories conditioned only on single images. In Section 5.2, the performance of DMBN is compared with MVAE and multi-step errors made by these models are analyzed in Section 5.3. In Section 5.4, we show how the latent space of two modalities indeed blends with each other. In Section 5.5, we analyze the behavior of our model when conditioned with images from different perspectives and whether it can serve as a mirror neuron system in replicating observations from different agents. We analyze whether such generalization is due to the inductive bias of the model with two different ablation studies in Sections 5.6 and 5.7, together which lend support to the idea that mirror neuron formation can be mediated by self-observation and modality blending with DMBN. Lastly, we test the generalization of the model by conditioning on out-of-distribution samples and include the results in Appendices A and B.

##### 5.1. Long-term Prediction with Vision only

In this experiment, we verify whether our system can produce visual sequences and the corresponding joint values given a single image as input. Note that since we take the average of latent vectors for conditioned points, we might as well give multiple



**Fig. 4.** (I) Images that are used as observations. (II) DMBN visual predictions for the given time steps. (III) Ground truth images for the given time steps. (IV) DMBN 7D joint predictions for the whole action.



**Fig. 5.** The prediction errors on the test set for different modality input–output pairs with the increasing size of the training data (x-axis).

images, instead of a single one, to get a more accurate prediction (see Eq. (2)). Here, to demonstrate the capabilities of our system even in such a scenario where the information is minimum, the system is fed with a single visual observation which is obtained just before the robot interacts with the object. The availability vector is set to one for visual modality and to zero for proprioceptive modality since the observation includes only visual information. Then, the system is requested to produce visual and motor signals from the beginning to the end of the movement.

Fig. 4 shows two examples of pushing and grasping actions at the left and the right of the figure respectively. Fig. 4.(I) shows the obtained images that are used as observations from the test set. Fig. 4.(II–III) shows the predicted images together with the ground truth at the corresponding time steps. It can be seen that by exploiting the position, orientation, and hand state of the robot extracted from the observed image, our system could successfully predict the sequences of visual and proprioceptive signals from start to finish, which are highly accurate compared to the ground truth values. It is notable that even though there was no proprioceptive observation in these two examples, our model could make accurate joint predictions from start to the end of the movement by just having access to visual modality (see Fig. 4.(IV)). These results indicate that our model can use the representation encoded from an available modality to predict the signals of the other missing modalities. A more detailed quantitative analysis of cross-modality predictions is presented in the next section.

## 5.2. Missing Modality Prediction as a Function of Training Set Size

In this section, we test whether DMBN can indeed perform well when there are some missing modalities. In this experiment, we used the same network in the previous experiment which is trained by using either *visual* or *proprioceptive* modalities. During the test phase, we set the availability of one of the two modalities to zero. We tested whether our system can still predict missing modalities.

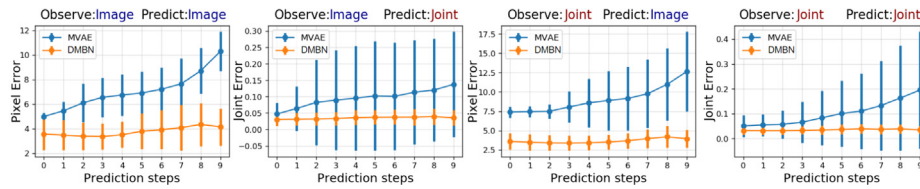
We compared our method with MVAE (Zambelli et al., 2020) as it can handle missing modalities. Moreover, we made several modifications to the original MVAE architecture to make a fair

comparison. First, we added convolutional layers for the visual input pipeline. All layers in the encoder and the decoder are exactly the same as in DMBN. Therefore, the number of parameters is the same except that MVAE uses an extra fully-connected layer to combine different encoder outputs. This extra layer is not needed in DMBN since the latent representation is shared and acquired via normalized weighted summation. Second, we remove the standard deviation prediction from the decoder as it gave better results in our preliminary experiments. We did not use the KL divergence term in the loss as in Zambelli et al. (2020). Third, we randomly mask the sensorimotor data at time  $t$  and predict the data at  $t + 1$ , in addition to other masking schemes reported in Zambelli et al. (2020). This additional masking scheme enables us to make full trajectory predictions (both forward and backward prediction) given the observation before contact. Our implementation<sup>2</sup> is based on Zambelli et al. (2020) and their code repository.<sup>3</sup>

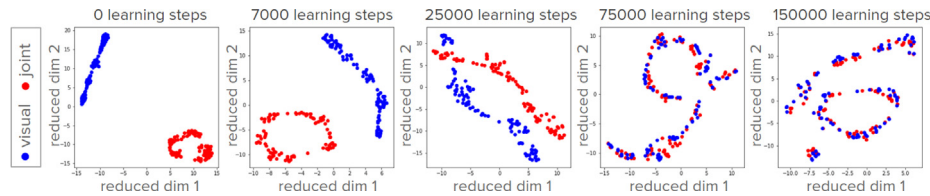
We report our results in Fig. 5 where the prediction accuracies with increasing number of training trajectories are shown. For the two modalities in our experimental setup, we tested four different combinations of modality masking: predicting visual states when either proprioceptive modality (Fig. 5.a) or the visual modality (Fig. 5.c) is missing, and predicting joint states when either proprioceptive modality (Fig. 5.b) or the visual modality (Fig. 5.d) is missing. We condition both DMBN and MVAE models with the observations taken from the same time step that is right before the robot interacts with the object. Both systems predict complete visual and joint trajectories starting from  $t = 0$  to  $t = T$ . Since DMBN is able to learn from few data, the error and its variation drop quickly even with the small training size, and it improves the accuracy while the data size is increased. For MVAE, the error slightly drops during the data size increase, yet, still far from DMBN. One reason for the error of MVAE is that it feeds the predictions back to itself as input, thus cascades the error in the long horizon. We investigate this phenomenon in the next section in detail.

<sup>2</sup> <https://github.com/alper111/multimodal-vae>

<sup>3</sup> [https://github.com/ImperialCollegeLondon/Zambelli2019\\_RAS\\_multimodal\\_VAE](https://github.com/ImperialCollegeLondon/Zambelli2019_RAS_multimodal_VAE)



**Fig. 6.** Multi-step prediction results. MVAE errors increase with the prediction steps due to error accumulation. However, our model preserves the error at the same level for increasing prediction steps since it predicts every time step independent from each other.



**Fig. 7.** t-SNE visualization of latent space during training. Blue points are visual encodings, and red points are joint encodings.

### 5.3. Analysis of long horizon predictions

In this section, we compared the capacity of DMBN on the long horizon predictions with the MVAE method. Both models are trained using the same two modalities in the same way as in the previous section.

In [Zambelli et al. \(2020\)](#), MVAE is used for one step ahead predictions to control the iCub humanoid robot in a closed loop. To make predictions about further time steps, the model can be fed with its output from the previous time step. They showed that when trained with sinusoidal data, the prediction accuracy remains the same for about 50 time steps, and then starts to degrade. In this experiment, we compared the two methods using the data that is collected during the self-exploration which is more complex and high-dimensional. In contrast to MVAE, DMBN does not need to feed its output back to itself as input to make further predictions since we can explicitly query any time step independently and make predictions on the long horizon directly.

We analyze the error versus the prediction step for two methods in [Fig. 6](#). The error of MVAE increases as the prediction step increases since the error is fed back in the input for future time step predictions. However, the error of DMBN remains around the same because the model does not have a feedback loop to connect an erroneously computed output to its input, and make predictions for every time step independently just by looking to the observations.

### 5.4. Multimodal latent space visualization

In this experiment, multi-modal latent space is visualized and analyzed. As mentioned in the previous section, we trained<sup>4</sup> the system with the two modalities that are *visual* and *proprioceptive*. For visualization purposes, the high-dimensional representation space (128 sized vector) is reduced to two dimensions using t-SNE ([van der Maaten & Hinton, 2008](#)) method at different stages of the training. [Fig. 7](#) shows the t-SNE visualization of the multimodal latent space at 0, 7k, 25k, 75k, and 150k learning steps from left to right. Blue and red points indicate the samples from the visual and proprioceptive modalities, respectively. [Fig. 7](#) shows that although the different modalities are clustered and separated from each other at the beginning of the training (0 and 7k learning steps), they start to share the representations

between each other after a while (25k learning steps), and turn into matching/overlapping representations in the later stages of the training (75k and 150k learning steps). Paired blue and red points in the overlapping representation space are analyzed and it is found that each paired blue-red point corresponds to two modalities recorded from the same state of the environment. These results suggest that our system can effectively learn multiple modalities in a common latent space in a way that every sensorimotor modality recorded from the same state of the environment ends up turning into the nearly same representation in the latent space. This allows our system to predict the missing modalities by using the representations produced by other available modalities, which was shown in [Section 5.2](#).

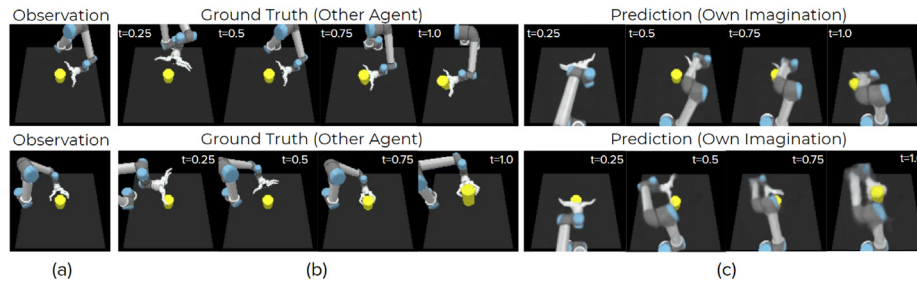
### 5.5. Imagining own actions by observing others: Emergence of mirror neuron system behavior

In this experiment, we tested our system to see if it can generate its own sensorimotor data by observing another agent perform an action. In order to do that, an agent was placed on the different sides of the table and their performed actions are observed via our agent's visual sensor. Note that in the training data, interactions were only performed and recorded just by our agent, so observing other agents in the test time is a novel information that is completely outside of our training set. Since we were using only visual data as the observation, the availability vector is set to one for visual modality and to zero for proprioceptive modality. Because the observations are on another agent but the predictions are made for our agent, this prediction process can be considered as forming a visual representation of the action of another agent for the self.

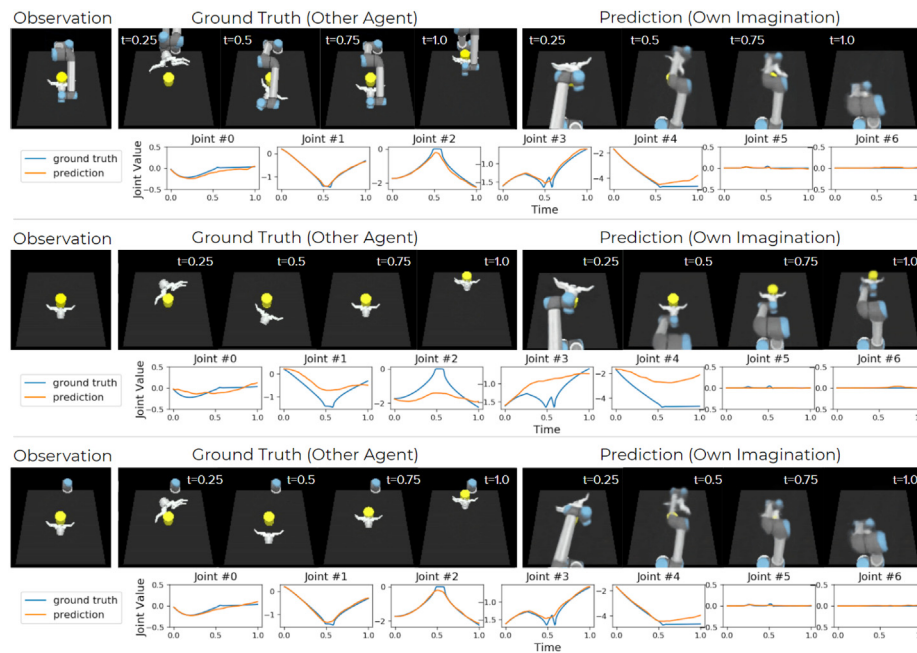
[Fig. 8](#) shows the prediction results of our model in two different pushing and grasping scenarios where the observations are shown in [Fig. 8\(a\)](#). In the first scenario, the other agent was placed on the opposite side of the table, and in the second scenario, the other agent was placed on the left side of the table. [Fig. 8\(b\)](#) shows the visual signals during the other agent performed its action, and [Fig. 8\(c\)](#) shows the full trajectory prediction of our system as it imagines the visual signals for itself. As it can be seen in the predictions, our agent is able to generate visual trajectories from its own perspective that matches the approaching angle and the action type in the observation, hence, imagining an action that would be an effect-based imitation of the observed action.

However, when we further analyzed our model, we saw that DMBN behaves differently in some specific scenarios. Surprisingly, when the other agent pulls the object towards itself, our

<sup>4</sup> Training details about the network can be found in [Appendix C](#).



**Fig. 8.** Examples of DMBN effect imitation behavior. First row: Observing the other agent just before it pushes away the object. Second row: Observing the other agent just before it grasps the object.



**Fig. 9.** Examples of DMBN egocentric imitation behavior. First row: Emergence of mirror neuron behavior where the agent observes the other agent pull the object towards itself. Second row: The agent observes a hand without the body. Third row: The agent observes a hand and a base without the arm.

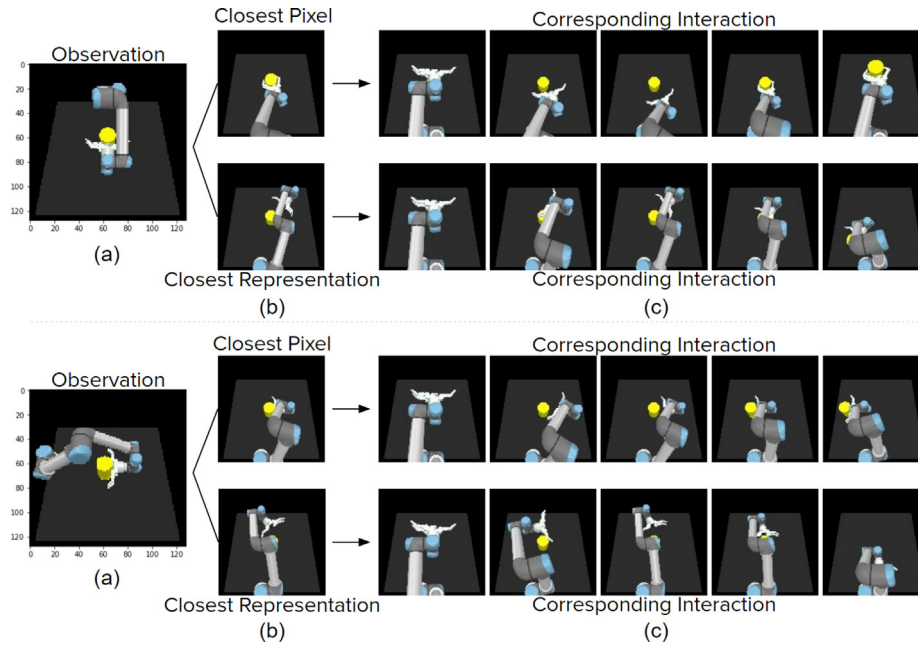
agent imagines an action that egocentrically imitates the observed action (Fig. 9, first row) and generates motor signals that would also pull the object towards itself rather than creating the effect on the object as shown in Fig. 8. We can say that, in this particular action observation case, an emergent mirror neuron property was exhibited by our DMBN. Interestingly this behavior ‘switches’ so that the action imagined corresponds to effect-based imitation (i.e. *emulation*) of the observed action when the body of the robot is removed from the interaction (Fig. 9, second row). Finally, when the robot is partially revealed by disclosing the base of the robot, the system starts to understand the observed action again as bringing the object toward one’s self (Fig. 9, third row), thereby showing a mirror neuron response as in the first row of Fig. 9.

These results show that when conditioned with the visual signals of other agents, DMBN has the potential to produce output signals similar to that of a mirror neuron system. However, signals generated can correspond to either effect-based or egocentric imitation depending on the specific visual signals available from the other agent and the environment. Therefore, it is viable to modulate the behavior of DMBNs via other cognitive mechanisms, e.g. attention, to purposefully control the operation of the model.

### 5.6. Template matching according to pixel and latent space distances

In this experiment, we aimed to see whether the mirror neuron emergence in our system was due to the rich representations constructed in the latent space during the learning, or it could be simply explained by a straightforward image-based template matching. For this, first, two test cases in which true mirror response (i.e. action representation that would yield egocentric imitation) was observed were selected. Fig. 10.(a) shows these two observation cases where the agents are placed on the opposite and the left side of the table, respectively. For each test case, the closest image in the training set that gives the minimum average pixel error, and the corresponding image of the closest representation that gives the minimum MSE error in the representation latent space are found and compared. Fig. 10.(b) shows the corresponding closest pixel and representation images for each case respectively. The closest pixel image is found by comparing the observed image with all of the images in the training dataset and selecting the image that gives the minimum pixelwise MSE error, where the closest representation is found by comparing the encoding of the observed image with all of the encodings in the training dataset and selecting the image that gives the minimum MSE encoding error. Fig. 10.(c) also shows





**Fig. 10.** Closest Pixel: Pixel Space Distance; Close Representation: Latent Space Distance.

**Table 1**

Test results of the two models in ten different training sessions with two action observation cases (see Fig. 10) of demonstrating agent positioned across (Case 1) and the left side of the agent (Case 2). Success: The model produced a signal output that corresponds to true mirror response; i.e. the execution of the action based on those signals would yield egocentric imitation. Fail: the model produced disturbed image signals.

True Mirror Response		Success	Fail
Image + Joint Model	Case 1	10	0
	Case 2	10	0
Only Image Model	Case 1	6	4
	Case 2	4	6

the corresponding full trajectory interactions of the found results from the training set.

Results of the both examples show that the corresponding interactions of the closest pixel images do not exhibit true mirror response (i.e. the predicted signals would not yield an egocentric imitation when executed on the robot). On the other hand, the corresponding interactions of the closest latent space representations show true mirror response. These results suggest that the output signals that DMBN produces are not based on a simple image-based error minimization but on rich representations that are learned during the multi-modal training with modality blending. The contribution of the deep modality blending to the mirror neuron emergence is further inspected in detail in the next section.

### 5.7. Analysis of the contribution of multimodal learning to mirror neuron emergence

In this experiment, we tested if deep modality blending contributes to the mirror neuron emergence in our system. To do that, a model that only uses visual modality was trained next to our model which was trained by using both visual and proprioceptive modalities. In order to prevent the training biases that can occur because of the initial network weights or sampling seeds, both models were trained 10 times with different random initializations. After the training, both of the models were tested with two test cases and checked whether the networks produce

output signals that correspond to mirror neuron emergence. The two test cases used in this experiment were the same examples as in the Experiment 5.6 where the demonstrating agents were placed at the opposite and the left side of the table (see Fig. 10).

Table 1 shows the results of the two models in ten different training initializations with two test cases. Results indicate that the model that uses deep modality blending (the model with Image + Joint) produces coherent images that correspond to egocentric imitation in every test case where the model that uses only one modality (Only Image Model) produces disturbed images on the ten test cases out of twenty. Fig. 11 shows some example fail cases for the only image model where the image is disturbed or the arm of the robot is disappeared. These results suggest that using deep modality blending with visual and proprioceptive modalities contribute to the emergence of mirror neuron behavior.

## 6. Conclusion

In this work, we proposed Deep Multi-modal Blending Network (DMBN) as a multi-modal action representation system that learns the sensorimotor signals corresponding to the actions, in a robust latent representation allowing temporal cross-modal predictions with limited information. DMBNs can generate complete signal trajectories in any desired modality even with zero information on the desired modality by using other available modalities. The performance of the network surpasses the available multi-modal learning systems due to long-range one-shot prediction capability and its novel stochastic modality blending mechanism.

DMBNs build powerful internal representations that allow surprisingly dynamic extrapolation properties, making it a strong contender as a feature-engineering-free Mirror Neuron System model. To be specific, after learning proprioception and visual signals based on self action observations, when tested with different perspective action observations, it successfully generates valid signals that represent its own actions. Depending on the visual setting, the network either acts a true mirror system matching an observed act to its own repertoire in an egocentric way, or acts as an effect-based action matching system. Thus, the network has

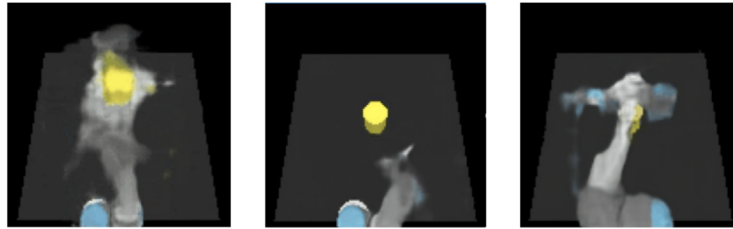


Fig. 11. Example failing scenarios for the only image model. The images are disturbed and the robot arm is disappearing.

potential to sustain egocentric and effect-based action recognition and imitation capabilities when envisioned in the cognitive system of an artificial or biological agent.

In this vein, future work should focus on developing biologically plausible and developmentally realistic end-to-end mirror neuron systems that learn along with sensorimotor skill acquisition. In the current study, we used a fixed action repertoire to systematically study the properties of DMBNs; yet in a developing artificial or biological cognitive agent, mirror neuron formation and action learning should go in parallel creating potentially non-trivial interactions worth studying. Another direction that should be pursued is to use the basic imitation capacity acquired by the model, to construct *novel imitation* capacity, where the parts of an observed novel act can be understood and matched to the existing action repertoire of the agent with the help of DMBN implementing the developing mirror neuron system. We believe that work around these directions will not only stimulate the computational study of mirror neurons as a full end-to-end system but also form a framework for lifelong sensorimotor learning for social robots. As a final point, investigating our model with Spiking Neural Networks (SNNs) is another direction to be pursued to develop biologically plausible systems. Recent studies show that state-of-the-art SNNs are real-time efficient (Yang, Wang, deng, Rahimi Azghadi, & Linares-Barranco, 2021b), scalable (Yang et al., 2019), and biologically plausible (Yang et al., 2020) systems that can be used in the applications such as: motor control with supervised learning (Yang, Wang, Zhang, deng, Pang, & Rahimi Azghadi, 2021c), real-world robotics (Lobov, Mikhaylov, Shamshin, Makarov, & Kazantsev, 2020), and object recognition with neurobotic control (Yang, Gao, Wang, Deng, Lansdell, & Linares-Barranco, 2021a). We believe that improvements towards this direction, such as converting our network to an SNN (Rueckauer, Lungu, Hu, Pfeiffer, & Liu, 2017), not only makes our model developmentally more realistic but also drastically reduces the power and memory consumption during the training phase which is caused by the large scale of layers and parameters in the model.

#### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

#### Acknowledgments

This research has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement no. 731761, IMAGINE; was partially supported by Japan Science and Technology Agency, Japan CREST “Cognitive Mirroring” under Grant No. JPMJCR16E2, by the International Joint Research Promotion Program of Osaka University, Japan under the project “Developmentally and biologically realistic modeling of perspective invariant action understanding” and by the Turkish Directorate of Strategy and Budget under the TAM Project

number 2007K12-873. The numerical calculations reported in this paper were partially performed at TUBITAK ULAKBIM, High Performance and Grid Computing Center (TRUBA resources).

#### Appendix A. Generalization of the system to the novel gripper configurations

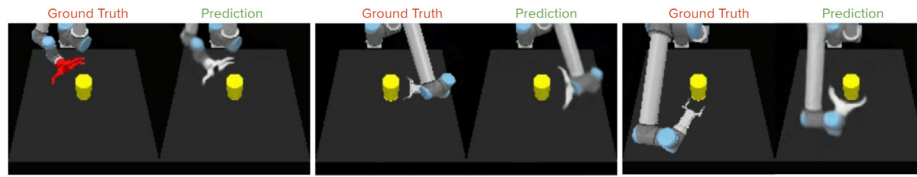
In this experiment, we tested our system with different novel gripper configurations that have different properties than in the training set. Fig. A.12 shows the generalization performances of the three different configurations. The left side of the figure shows a scenario in which the color of the gripper was changed to red. The middle and the right side of the figure shows two sample scenarios where a novel gripper is used to perform push and grasp actions, respectively. Although there are no samples with different colors or different types of grippers in the training set, our system could successfully predict the correct actions in all scenarios. It can be seen that in the results the gripper is predicted as in its own configuration in the training set even the test configurations of the gripper are different. These results suggest that when encountered with novel configurations, our system could transfer the knowledge extracted from these configurations to make predictions for its own form that it was trained for.

#### Appendix B. Generalization of the system to the novel environmental configurations

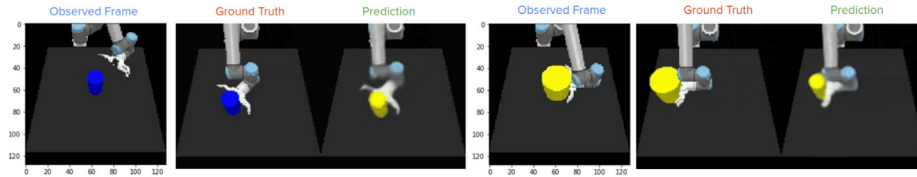
In this experiment, we tested our system with different novel environmental configurations that have different properties than the training set. Fig. B.13 shows the generalization performances of the two different configurations. Left side of the figure shows a scenario in which the color of the object was different from the object in the training data, and the right side shows a configuration where the size of the object was changed. Despite not seeing a big or blue object in the training, our system could successfully predict the correct approaching angle and the action using the observed image in both configurations. It can be seen that the color and the size of the objects are predicted as in the configuration in the training images. This is expected since the only configuration for the object in the training scene was yellow and small. Even though the object in the observed image was not the same with the training object, our system could use the knowledge that is learned in the training data to predict a correct output in its own configurations that satisfies the given observation.

#### Appendix C. Network architecture and training details of DMBN

In this section, the network architecture and training configurations of DMBN are shown. Tables C.2 and C.3 show the image and joint encoder architectures respectively. Tables C.4 and C.5 show the image and joint decoder architectures respectively. DMBN is trained with Adam optimizer (Kingma & Ba, 2014) for one million iterations with a batch size of one and a learning rate of 0.0001. We set  $obs_{max}$  to 5.



**Fig. A.12.** Generalization performance of the proposed system in three different configurations. Left: The color of the gripper is red. Middle: The robot is performing a push action with a novel gripper. Right: The robot is performing a grasp action with a novel gripper.



**Fig. B.13.** Generalization performance of the proposed system in two different configurations. Left side: the color of the object is blue. Right side: the size of the object is bigger than the original one.

**Table C.2**

DMBN Image Encoder.

Layer	Input size	Output size
Conv3 × 3 + ReLU + MaxPool2 × 2	(4, 128, 128)	(32, 64, 64)
Conv3 × 3 + ReLU + MaxPool2 × 2	(32, 64, 64)	(64, 32, 32)
Conv3 × 3 + ReLU + MaxPool2 × 2	(64, 32, 32)	(64, 16, 16)
Conv3 × 3 + ReLU + MaxPool2 × 2	(64, 16, 16)	(128, 8, 8)
Conv3 × 3 + ReLU + MaxPool2 × 2	(128, 8, 8)	(128, 4, 4)
Conv3 × 3 + ReLU + MaxPool2 × 2	(128, 4, 4)	(256, 2, 2)
Flatten	(256,2,2)	1024
Dense	1024	128
Multiply (Image Coefficient)	128 * 128	128

**Table C.3**

DMBN Joint Encoder.

Layer	Input size	Output size
Dense + ReLU	8	32
Dense + ReLU	32	64
Dense + ReLU	64	64
Dense + ReLU	64	128
Dense + ReLU	128	128
Dense + ReLU	128	256
Dense + ReLU	256	128
Multiply (Joint Coefficient)	128 * 128	128

**Table C.4**

DMBN Image Decoder.

Layer	Input size	Output size
Add (Image + Joint Representations)	128 + 128	128
Concatenate (Target Layer)	128	129
Dense + ReLU	129	1024
Reshape	1024	(256, 2, 2)
Conv3 × 3 + ReLU + UpSample2 × 2	(256, 2, 2)	(256, 4, 4)
Conv3 × 3 + ReLU + UpSample2 × 2	(256, 4, 4)	(128, 8, 8)
Conv3 × 3 + ReLU + UpSample2 × 2	(128, 8, 8)	(128, 16, 16)
Conv3 × 3 + ReLU + UpSample2 × 2	(128, 16, 16)	(64, 32, 32)
Conv3 × 3 + ReLU + UpSample2 × 2	(64, 32, 32)	(64, 64, 64)
Conv3 × 3 + ReLU + UpSample2 × 2	(64, 64, 64)	(32, 128, 128)
Conv3 × 3 + ReLU	(32, 128, 128)	(16, 128, 128)
Conv3 × 3 + ReLU	(16, 128, 128)	(8, 128, 128)
Conv3 × 3 + Sigmoid	(8, 128, 128)	(3, 128, 128)

**Table C.5**

DMBN Joint Decoder.

Layer	Input size	Output size
Add (Image + Joint Representations)	128 + 128	128
Concatenate (Target Layer)	128	129
Dense + ReLU	129	1024
Dense + ReLU	1024	512
Dense + ReLU	512	216
Dense + ReLU	216	128
Dense + ReLU	128	32
Dense	32	14

**Table D.6**

MVAE Image encoder.

Layer	Input size	Output size
Conv3 × 3 + ReLU + MaxPool2 × 2	(6, 128, 128)	(32, 64, 64)
Conv3 × 3 + ReLU + MaxPool2 × 2	(32, 64, 64)	(64, 32, 32)
Conv3 × 3 + ReLU + MaxPool2 × 2	(64, 32, 32)	(64, 16, 16)
Conv3 × 3 + ReLU + MaxPool2 × 2	(64, 16, 16)	(128, 8, 8)
Conv3 × 3 + ReLU + MaxPool2 × 2	(128, 8, 8)	(128, 4, 4)
Conv3 × 3 + ReLU + MaxPool2 × 2	(128, 4, 4)	(256, 2, 2)
Flatten	(256, 2, 2)	1024
Dense + ReLU	1024	128

**Table D.7**

MVAE Joint encoder.

Layer	Input units	Output units
Dense+ReLU	14	32
Dense+ReLU	32	64
Dense+ReLU	64	64
Dense+ReLU	64	128
Dense+ReLU	128	128
Dense+ReLU	128	256
Dense+ReLU	256	128

**Table D.8**

MVAE shared encoder–decoder. The activation after the first decoder layer is sliced into two, and each slice is given to a different decoder.

Layer	Input units	Output units
<b>Encoder</b>		
Concatenate (Image+Joint)	128, 128	256
Dense + Tanh	256	128 mean, 128 std
<b>Decoder</b>		
Dense+ReLU	128	256
Slice (for image and joint dec.)	256	128, 128

## Appendix D. Network architecture and training details of MVAE

In this section, the network architecture and training configurations of MVAE are shown. Tables D.6 and D.7 show the image and the joint encoder architectures respectively. Table D.8

**Table D.9**

MVAE Image Decoder. The last activation is sliced into two (6, 128, 128) shaped tensors for mean and std. See the original implementation (Zambelli et al., 2020) for further details.

Layer	Input size	Output size
Dense + ReLU	128	1024
Reshape	1024	(256, 2, 2)
Conv3 × 3 + ReLU + UpSample2 × 2	(256, 2, 2)	(256, 4, 4)
Conv3 × 3 + ReLU + UpSample2 × 2	(256, 4, 4)	(128, 8, 8)
Conv3 × 3 + ReLU + UpSample2 × 2	(128, 8, 8)	(128, 16, 16)
Conv3 × 3 + ReLU + UpSample2 × 2	(128, 16, 16)	(64, 32, 32)
Conv3 × 3 + ReLU + UpSample2 × 2	(64, 32, 32)	(64, 64, 64)
Conv3 × 3 + ReLU + UpSample2 × 2	(64, 64, 64)	(32, 128, 128)
Conv3 × 3 + ReLU	(32, 128, 128)	(16, 128, 128)
Conv3 × 3 + ReLU	(16, 128, 128)	(12, 128, 128)
Conv3 × 3	(12, 128, 128)	(12, 128, 128)

shows the shared encoder–decoder architecture. Tables D.9 and D.10 show the image and joint decoder architectures respectively. MVAE is trained with Adam optimizer (Kingma & Ba, 2014) for 200 epochs with a batch size of 128 and a learning rate of 0.001.

**Table D.10**

MVAE Joint Decoder. The last activation is sliced into two for mean and std.

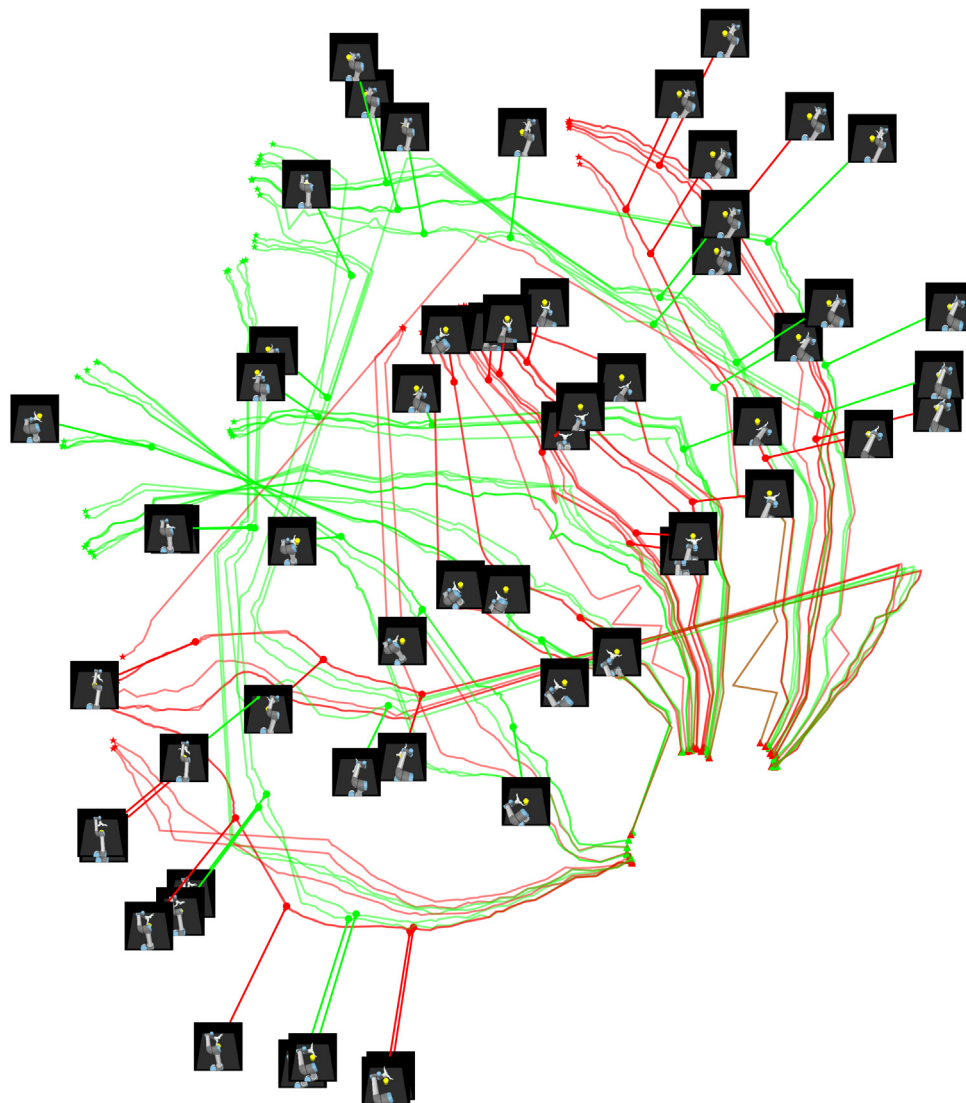
Layer	Input units	Output units
Dense+ReLU	128	256
Dense+ReLU	256	128
Dense+ReLU	128	128
Dense+ReLU	128	64
Dense+ReLU	64	64
Dense+ReLU	64	32
Dense	32	28

**Appendix E. t-SNE visualization of the latent space**

In this section, the detailed version of the latent space is investigated. Fig. E.14 shows the encodings of all of the training trajectories in the latent space.

**Appendix F. Supplementary data**

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.neunet.2021.11.004>.



**Fig. E.14.** t-SNE (van der Maaten & Hinton, 2008) visualization of the encoder output. Here, green and red represent ‘move’ and ‘grasp’ actions, respectively. The initial and the final point of a trajectory is represented with a triangle and a star, respectively.

## References

- Akbulut, M., Bozdogan, U., Tekden, A., & Ugur, E. (2021). Reward conditioned neural movement primitives for population based variational policy optimization. In *International conference on robotics and automation (ICRA)*.
- Akbulut, M. T., Oztop, E., Seker, M. Y., Xue, H., Tekden, A. E., & Ugur, E. (2020). ACNMP: Skill transfer and task extrapolation through learning from demonstration and reinforcement learning via representation sharing. In *4th Conference on Robot Learning (CoRL 2020)*.
- Argall, B. D., Chernova, S., Veloso, M., & Browning, B. (2009). A survey of robot learning from demonstration. *Robotics and Autonomous Systems*, 57(5), 469–483.
- Asfour, T., Azad, P., Gyafas, F., & Dillmann, R. (2008). Imitation learning of dual-arm manipulation tasks in humanoid robots. *International Journal of Humanoid Robotics*, 5, 183–202. <http://dx.doi.org/10.1109/ICHR.2006.321361>.
- Atkeson, C. G., Moore, A. W., & Schaal, S. (1997). Locally weighted learning for control. In *Lazy learning* (pp. 75–113). Springer.
- Ben Amor, H., Kroemer, O., Hillenbrand, U., Neumann, G., & Peters, J. (2012). Generalization of human grasping for multi-fingered robot hands. In *IROS*.
- Bonaiuto, J., & Arbib, M. A. (2010). Extending the mirror neuron system model, II: what did I just do? A new role for mirror neurons. *Biological Cybernetics*, 102(4), 341–359.
- Bonaiuto, J., Rosta, E., & Arbib, M. (2007). Extending the mirror neuron system model, I - Audible actions and invisible grasps. *Biological Cybernetics*, 96(1), 9–38.
- Calinon, S. (2016). A tutorial on task-parameterized movement learning and retrieval. *Intelligent Service Robotics*, 9(2016), <http://dx.doi.org/10.1007/s11370-015-0187-9>.
- Calinon, S., Evrard, P., Gribovskaya, E., Billard, A., & Kheddar, A. (2009). Learning collaborative manipulation tasks by demonstration using a haptic interface. In *Advanced robotics* (pp. 1–6).
- Chu, V., McMahon, I., Riano, L., McDonald, C. G., He, Q., Martinez Perez-Tejada, J., et al. (2013). Using robot exploratory procedures to learn the meaning of haptic adjectives. In *ICRA* (pp. 3048–3055).
- Copete, J. L., Nagai, Y., & Asada, M. (2016a). Motor development facilitates the prediction of others' actions through sensorimotor predictive learning. In *2016 joint IEEE international conference on development and learning and epigenetic robotics (ICDL-epirob)* (pp. 223–229).
- Copete, J. L., Nagai, Y., & Asada, M. (2016b). Motor development facilitates the prediction of others' actions through sensorimotor predictive learning. In *2016 joint IEEE international conference on development and learning and epigenetic robotics (ICDL-EpiRob)* (pp. 223–229). <http://dx.doi.org/10.1109/DEVLRN.2016.7846823>.
- Demiris, Y., & Johnson, M. (2003). Distributed, predictive perception of actions: a biologically inspired robotics architecture for imitation and learning. *Connection Science*, 15(4), 231–243. <http://dx.doi.org/10.1080/09540090310001655129>.
- Di Pellegrino, G., Fadiga, L., Fogassi, L., Gallese, V., & Rizzolatti, G. (1992). Understanding motor events: A neurophysiological study. *Experimental Brain Research*, 91, 176–180. <http://dx.doi.org/10.1007/BF00230027>.
- Droniou, A., Ivaldi, S., & Sigaud, O. (2015). Deep unsupervised network for multimodal perception, representation and classification. *Robotics and Autonomous Systems*, 71, 83–98.
- Garnelo, M., Rosenbaum, D., Maddison, C., Ramalho, T., Saxton, D., Shanahan, M., et al. (2018). Conditional neural processes. In *ICML* (pp. 1704–1713).
- Girgin, H., & Ugur, E. (2018). Associative skill memory models. In *IROS* (pp. 6043–6048).
- Huang, Y., Rozo, L., Silvério, J., & Caldwell, D. (2019). Kernelized movement primitives. *International Journal of Robotics Research*, 38, 833–852. <http://dx.doi.org/10.1177/0278364919846363>.
- Jung, M., Matsumoto, T., & Tani, J. (2019). Goal-directed behavior under variational predictive coding: Dynamic organization of visual attention and working memory. <http://arxiv.org/abs/1903.04932>.
- Keysers, C., Kohler, E., Umiltà, M. A., Nanetti, L., Fogassi, L., & Gallese, V. (2003). Audiovisual mirror neurons and action recognition. *Experimental Brain Research*, 153(4), 628–636.
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980).
- Kohler, E., Keysers, C., Umiltà, M. A., Fogassi, L., Gallese, V., & Rizzolatti, G. (2002). Hearing sounds, understanding actions: action representation in mirror neurons. *Science*, 297(5582), 846–848.
- Kramberger, A., Gams, A., Nemeč, B., Chrysostomou, D., Madsen, O., & Ude, A. (2017). Generalization of orientation trajectories and force-torque profiles for robotic assembly. *Robotics and Autonomous Systems*, [ISSN: 0921-8890] 98, 333–346. <http://dx.doi.org/10.1016/j.robot.2017.09.019>.
- Lee, D., & Ott, C. (2011). Incremental kinesthetic teaching of motion primitives using the motion refinement tube. *Autonomous Robots*, 31(2–3), 115–131.
- Lobov, S. A., Mikhaylov, A. N., Shamshin, M., Makarov, V. A., & Kazantsev, V. B. (2020). Spatial properties of STDP in a self-learning spiking neural network enable controlling a mobile robot. *Frontiers in Neuroscience*, 14, 88.
- Mühlig, M., Gienger, M., & Steil, J. (2012). Interactive imitation learning of object movement skills. *Autonomous Robots*, 32, 97–114.
- Nagai, Y., Kawai, Y., & Asada, M. (2011). Emergence of mirror neuron system: Immature vision leads to self-other correspondence. Vol. 2, In *2011 IEEE international conference on development and learning (ICDL)* (pp. 1–6). <http://dx.doi.org/10.1109/DEVLRN.2011.6037335>.
- Noda, K., Arie, H., Suga, Y., & Ogata, T. (2014). Multimodal integration learning of robot behavior using deep neural networks. *Robotics and Autonomous Systems*, [ISSN: 0921-8890] 62(6), 721–736.
- Oztop, E., & Arbib, M. A. (2002). Schema design and implementation of the grasp-related mirror neuron system. *Biological Cybernetics*, 87, 116–140.
- Oztop, E., Kawato, M., & Arbib, M. A. (2013). Mirror neurons: Functions, mechanisms and models. *Neuroscience Letters*, 540, 43–55.
- Pahič, R., Gams, A., Ude, A., & Morimoto, J. (2018). Deep encoder-decoder networks for mapping raw images to dynamic movement primitives. In *2018 IEEE international conference on robotics and automation (ICRA)* (pp. 5863–5868). <http://dx.doi.org/10.1109/ICRA.2018.8460954>.
- Paraschos, A., Daniel, C., Peters, J., & Neumann, G. (2018). Using probabilistic movement primitives in robotics. *Autonomous Robots*, 42(2018), <http://dx.doi.org/10.1007/s10514-017-9648-7>.
- Pastor, P., Hoffmann, H., Asfour, T., & Schaal, S. (2009). Learning and generalization of motor skills by learning from demonstration. In *ICRA* (pp. 763–768).
- Pastor, P., Righetti, L., Kalakrishnan, M., & Schaal, S. (2011). Online movement adaptation on previous sensor experiences. In *IROS*.
- Pervez, A., & Lee, D. (2018). Learning task parameterized dynamic movement primitives using mixture of GMMs. *Intelligent Service Robotics*, 11, 61–78.
- Pervez, A., Mao, Y., & Lee, D. (2017). Learning deep movement primitives using convolutional neural networks. In *2017 IEEE-RAS 17th international conference on humanoid robotics (humanoids)* (pp. 191–197). IEEE.
- Rizzolatti, G., Fadiga, L., Gallese, V., & Fogassi, L. (1996). Premotor cortex and the recognition of motor actions. *Cognitive Brain Research*, 3(2), 131–141.
- Rohmer, E., Singh, S. P. N., & Freese, M. (2013). Coppeliastik (formerly V-REP): a versatile and scalable robot simulation framework. In *Proc. of the international conference on intelligent robots and systems (IROS)*. [www.coppeliarobotics.com](http://www.coppeliarobotics.com).
- Rueckauer, B., Lungu, I.-A., Hu, Y., Pfeiffer, M., & Liu, S.-C. (2017). Conversion of continuous-valued deep networks to efficient event-driven networks for image classification. *Frontiers in Neuroscience*, [ISSN: 1662-453X] 11, 682. <http://dx.doi.org/10.3389/fnins.2017.00682>.
- Schaal, S. (2006). Dynamic movement primitives—a framework for motor control in humans and humanoid robotics. In *Adaptive motion of animals and machines* (pp. 261–280). Springer.
- Seker, M. Y., Imre, M., Piater, J., & Ugur, E. (2019). Conditional neural movement primitives. In *Proceedings of robotics: science and systems*. Freiburg/Breisgau, Germany: <http://dx.doi.org/10.15607/RSS.2019.XV.071>.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1), 1929–1958.
- Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction*. MIT Press.
- Suzuki, M., Nakayama, K., & Matsuo, Y. (2016). Joint multimodal learning with deep generative models. arXiv preprint [arXiv:1611.01891](https://arxiv.org/abs/1611.01891).
- Tani, J., Ito, M., & Sugita, Y. (2004). Self-organization of distributedly represented multiple behavior schemata in a mirror system: reviews of robot experiments using RNNPB. *Neural Networks*, 17(8–9), 1273–1289.
- Ude, A., Gams, A., Asfour, T., & Morimoto, J. (2010). Task-specific generalization of discrete and periodic dynamic movement primitives. *IEEE Transactions on Robotics*, 26(5), 800–815.
- Ugur, E., & Girgin, H. (2020). Compliant parametric dynamic movement primitives. *Robotica*, 38(3), 457–474. <http://dx.doi.org/10.1017/S026357471900078X>.
- van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(86), 2579–2605.
- Vijayakumar, S., & Schaal, S. (2000). Locally weighted projection regression: Incr. real time learning in high dimensional space. In *ICML* (pp. 1079–1086).
- Vincent, P., Larochelle, H., Bengio, Y., & Manzagol, P.-A. (2008). Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on machine learning* (pp. 1096–1103).
- Waibel, A., Hanazawa, T., Hinton, G., Shikano, K., & Lang, K. J. (1989). Phoneme recognition using time-delay neural networks. *IEEE Transactions on Acoustics, Speech and Signal Processing*, 37(3), 328–339.
- Wu, M., & Goodman, N. (2018). Multimodal generative models for scalable weakly-supervised learning. In *Advances in neural information processing systems* (pp. 5575–5585).
- Xie, F., Chowdhury, A., De Paolis Kaluza, M. C., Zhao, L., Wong, L., & Yu, R. (2020). Deep imitation learning for bimanual robotic manipulation. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, & H. Lin (Eds.), Vol. 33, *Advances in neural information processing systems* (pp. 2327–2337). Curran Associates, Inc..

- Yang, S., Deng, B., Wang, J., Li, H., Lu, M., Che, Y., et al. (2020). Scalable digital neuromorphic architecture for large-scale biophysically meaningful neural network with multi-compartment neurons. *IEEE Transactions on Neural Networks and Learning Systems*, 31(1), 148–162. <http://dx.doi.org/10.1109/TNNLS.2019.2899936>.
- Yang, S., Gao, T., Wang, J., Deng, B., Lansdell, B., & Linares-Barranco, B. (2021a). Efficient spike-driven learning with dendritic event-based processing. *Frontiers in Neuroscience*, [ISSN: 1662-453X] 15, 97.
- Yang, S., Wang, J., Deng, B., Liu, C., Li, H., Fietkiewicz, C., et al. (2019). Real-time neuromorphic system for large-scale conductance-based spiking neural networks. *IEEE Transactions on Cybernetics*, 49(7), 2490–2503. <http://dx.doi.org/10.1109/TCYB.2018.2823730>.
- Yang, S., Wang, J., deng, B., Rahimi Azghadi, M., & Linares-Barranco, B. (2021b). Neuromorphic context-dependent learning framework with fault-tolerant spike routing. *IEEE Transactions on Neural Networks and Learning Systems*, PP, 1–15. <http://dx.doi.org/10.1109/TNNLS.2021.3084250>.
- Yang, S., Wang, J., Zhang, N., deng, B., Pang, Y., & Rahimi Azghadi, M. (2021c). Cerebellumorphic: Large-scale neuromorphic model and architecture for supervised motor learning. *IEEE Transactions on Neural Networks and Learning Systems*, PP, 1–15. <http://dx.doi.org/10.1109/TNNLS.2021.3057070>.
- Zambelli, M., Cully, A., & Demiris, Y. (2020). Multimodal representation models for prediction and control from partial information. *Robotics and Autonomous Systems*, 123, Article 103312.
- Zhou, Y., & Asfour, T. (2017). Task-oriented generalization of dynamic movement primitive. In *IROS* (pp. 3202–3209).