

Bio-inspired Visuomotor Control of Robot to Reach 3D Objects

Xiaodan Chen, Alexandre Pitti

University of Cergy-Pontoise

1 Introduction

As the active development of cooperation between human and robot, the demand for robots to accomplish tasks independently arises either. This kind of robot has the characteristic of being high adaptive to the changing environment in the real-time and the ability of self-correct when it comes to the calibration problems. Studies in human infants and chimpanzees infant give us a thread for conceiving a self-accomplishment of subconscious reaching and grasping task of these cognitive robots.

A test on the 7 infants between 6 and 25 weeks of age reveals that human infants rely on proprioceptive information, which is the emergence of the interaction between eye and hand, to control and direct their arm toward a specific location in 3-dimensional space[1][2]. Thus, correlating the sensory signals with variations occurring in the motor space is essential for the development of complex motor abilities. Inside our brain, more precisely in the posterior parietal cortex(PPC), there are neuron networks who encode them in order to control our hands to get close to an object and finally grasp it. [3] explains it by taking a cup as our hand target as an example. The spatial location of this cup is first defined within visual coordinates and then neurons will convert this cup representations with respect to the eye and to the hand. Within the PPC, these multi-modalities significantly modulate the encoding and processing of one another [4]. Records of these neurons' activities show that targets are encoded in reference frames that span the continuum between such intuitive cases[5]. [6][7] assume that during this early stage of planning, the location of the goal is remapped into egocentric coordinates [8][9], which is then associated with hand position information to form a simplified hand-centered plan of the intended movement trajectory, according to vectorial planning hypotheses. We propose hence to apply a modulated network on the manipulation of the robot for the task of reaching by means of correlating the coordination of the robot's arm with the information of the target.

It is important to discriminate grasping activity from simply reaching movement for the reason that firstly, supplemental information of target is required, e.g. size of the object, orientation of the target, etc. Besides, the neurons relevant to grasping behaviour lie mostly in primary somatosensory cortex, instead of PPC. In addition, when arms move in same trajectory but with different orientation of the palm such as pronation and supination, the activities of neurons in the same region vary from one to another, indicating the existence of preferred orientation between them. In [10], Iberall and Arbib proposed

a reinforcement-type model named Infant Learning to Grasp Model with a conception known as "Affordance". However, in practice, we will not adapt completely this model for some particular reasons, which will be explained later in the next session.

Nevertheless, how the neurons encode remains unknown. Based on micro-electrode penetrations, Georgopoulos et al. show that preferred directions range throughout the 3D directional continuum and are multiply represented in the arm area of the motor cortex. They assume that a motor cortical cell might relate to weighted combinations of muscles[11]. Besides, [12][13][14] posit that the way it is done is by multiplicative interaction across the different sensorimotor signals to convolve conditionally variables from each other. Additionally, Michael et al. analyze in [15] that the units in the output layer exploit complex, gaze-centered visual receptive fields, rifully exist vectorial muscle commands, and "gain-field"-type eye position sensitivities. This leads to subtle adjustments of each output unit in the corresponding motor contributions, so that sum of vector of each involved neuron, say, definitive vector, rotates into the preferred direction, assuring that the relative muscles are proportionally activated and then the effector moves accordingly.

Inspired from it, we can then obtain visuo-motor coordinations by manipulating the product of sensory variables for learning transformations through multiplicative gain-field network. This kind of mathematical method has some advantages such as a better discrimination than deep networks affine transformations[16], applied extensively by Memisevic to the learning of optical flow. We name this network that undertakes the transformation from visual reference frame to extrinsic egocentric reference frame a retina-centered gain-field network. On the other hand, we define this model that guides the robot arm toward a visual target by involving mathematically a computed vectorial transformation from direction in visual space to direction in motor space a forward model.

Inversely, within the robot, it is always possible to get the information of motor, such as joint representations as well as the effect of the command in 3 dimension space. Knowing the effect of the predictive command and correlating it with the current sensory, we can verify if the calculated command should be executed autonomously. As a result, it becomes a loop if we combine this inverse model with the forward model mentioned above, the same as [17][18][19].

There are also other choices for this task. For example, [20][21] proposed to employ multi-layered perceptrons

whereas [22][23][24][25][26][27] used other techniques based on Bayesian networks and Gaussian mapping. Our idea of using gated network is based on the fruits of others researchers in the domain of Neuroscience as we mentioned previously. Until now, gain-field network has been applied in several experiments. [28][29] developed it for categorizing and retrieving motor sequences with the ICub and at the same time, it was applied for audio-visual and visuomotor integration and adaptation as during tool-use [30][31][32]. However, very few teams practices it for a 3-dimension reaching and grasping on an open source robot Reachy.

In this paper, we begin by introducing what the Gated Network is. Moreover, we will come up with some explications before starting the experiments. In particular we show, through computational modeling, how a gain-field network realize with or without preferred direction or preferred orientation a basic behavior (reaching) and a more complex behavior (grasping) through interactive goal-directed trial.

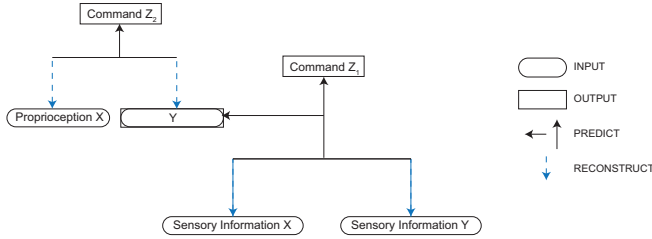


Figure 1: Left Top: Hand-centered Gain Field Network. Right Bottom: Eye-centered Gain Field Network. Gain-field network for reaching and grasping task, an integration of sensorial, motor and proprioceptive information. The Gain field network (on the left) learns commands Z integrating the proprioceptive information that refers to the encoding of parameters between each joint and visual direction which is also the output of another Gain Field network. The Gain field network on the right learns the relative hand-centered reference frame, the visual direction, from mapping of the Z and the visual input X and Y . The results are hand-centered receptive field for grasping tasks of nearby targets.

2 Method and Experimental Setup

2.1 Gain-Field Network

Gain-field network, different from bi-partite networks, is a model involving three-way multiplicative interactions that each unit in the model can gate connections between visible units, in which the latent variables act as dynamic mapping units by encoding the relationship between other variables, with the dependency of the value of each variable which is thus no longer the weighted sum but the product of the other variables [33].

In [34], Paul Smolensky named this kind of product a Sigma-Pi unit, as shown in Figure 2](a). Being multiplied before entering into Sigma function, these multiplicative pairs possess the characteristic that if one of the pairs whose value is zero, other pairs will have no effect on the output, while with its value equals one, the value of output will depend only on other pairs, without any probability distribution calculations. These two special pairs indicate exactly how one conjunct can

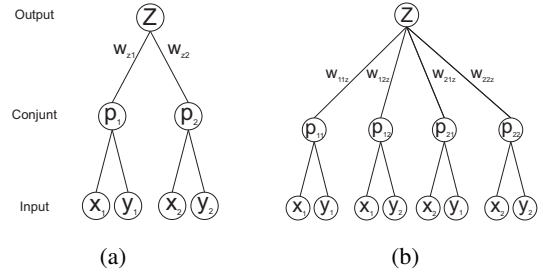


Figure 2: (a) Sigma-Pi Units. The input to unit P is the conjunct X_1 and Y_1 and X_2 and Y_2 . The value of unit Z equals to the weighted sum of units P which is the products of $X_1 Y_1$ and $X_2 Y_2$. (b) Sigma-Pi units in Gain-field network, where input pairs are the combination of matrix X and Y , which contribute proportionally to the unit P .

modulate other conjuncts through the Pi connection. Take p as conjunct, x and y as inputs in visible layer and suppose that $I \times 1$ (resp. $J \times 1$) is the number of elements in input vector X ($x \in X$) (resp. Y ($y \in Y$)), we can get:

$$Z_k = \sum_{k=1}^K W_{nk} \cdot P_n \quad (1)$$

$$P_n = \prod_{t=1}^N a_{nt} = x \cdot y \quad (2)$$

Gated-network has developed Sigma-Pi units, of which the conjuncts alter into the combination of matrix X and Y , among which have different weights for the contribution to the unit P as shown in Figure 2](b). Mathematically, by adding synaptic weights to each input and appending the subscript i,j to X,Y respectively, here comes the output Z in the hidden layer of Gain field network:

$$\begin{aligned} Z_k &= \sum_{i,j,k=1}^{I,J,K} W_{ijk} \cdot P_{ij} \\ &= \sum_{i,j,k=1}^{I,J,K} W_{ijk} \cdot X_i \cdot Y_j \end{aligned} \quad (3)$$

where adaptive parameters have 3 subscripts ijk denote two inputs and the output itself separately. The 3(b) exhibits it in a more concrete way which reveals a notable problem with this network that the gating parameter W is a cubic-like matrix (only when $I = J = K$, it turns out to be an exact cubic matrix), resulting in an exponential relation between the number of training samples and a better generalisation of the parameter.

For this reason, we introduce the factorisation of modulation connections for the reason that by allowing to sharing features, the dimension of this matrix will reduce and at the same time the network will improve its performance by requiring less priors, as shown in [35]. We obtain then the intermediate equa-

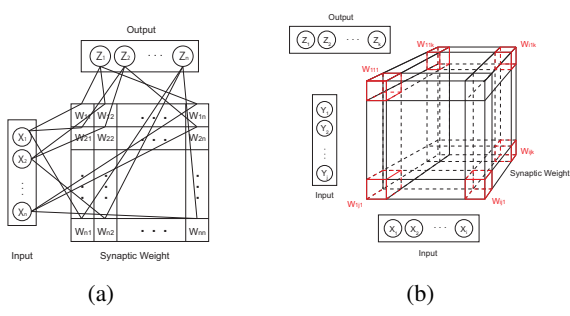


Figure 3: Comparison between two networks. (a) Standard Feature Learning Model. (b) Gated Feature Learning Model Before Factorisation

tion as below:

$$\sum_{i,j,k=1}^{I,J,K} W_{ijk} \cdot X_i \cdot Y_j \cdot Z_k$$

$$= \sum_{h=1}^H \left(\sum_{i=1}^I W_{ih} X_i \right) \cdot \left(\sum_{j=1}^J W_{jh} Y_j \right) \left(\sum_{k=1}^K W_{kh} Z_k \right) \quad (4)$$

where H is the number of hidden factors that should be determined by hand or by cross-validation[33], and W_{ih} , W_{jh} , W_{kh} are respectively the projection of X , Y , Z onto latent factors, as Figure 4(b). Combining Equation 3 and Equation 4 with respect to distributive property:

$$Z_k = \sum_{h=1}^H \left(\sum_{i=1}^I W_{ih} X_i \right) \cdot \left(\sum_{j=1}^J W_{jh} Y_j \right) W_{kh} \quad (5)$$

where the dimensionality of adaptive parameter reduces to $O(I \times H + J \times H + K \times H)$, approximately $O(N^2)$. For the sake of simplification, we define $\sum_{i=1}^I W_{ih} X_i$, $\sum_{j=1}^J W_{jh} Y_j$, $\sum_{k=1}^K W_{kh} Z_k$ as G_x , G_y , G_z respectively. Furthermore, since terms X and Y are all linear terms as Z , so the Equation 5 satisfies commutativity:

$$X_i = \sum_{h=1}^H W_{ih} \cdot G_y \cdot G_z \quad (6)$$

$$Y_j = \sum_{h=1}^H W_{jh} \cdot G_x \cdot G_z \quad (7)$$

$$Z_k = \sum_{h=1}^H W_{kh} \cdot G_x \cdot G_y \quad (8)$$

where Z_k is the output of the k th high-order neuron unit. Equation 6 to 8 show that z (resp. x , y) comes from x (resp. y , z) and y (resp. z , x), and that Sigma-Pi neural network is a generative model combining bottom-up and top-down interactions, which corresponds to the necessity of information process inside a developmental robot[29]. Therefore, the boundary between latent layer and visible layer on a holistic scale becomes blurred, for the reason that the output calculated from the input in an initial model will be used as input in its mirror model.

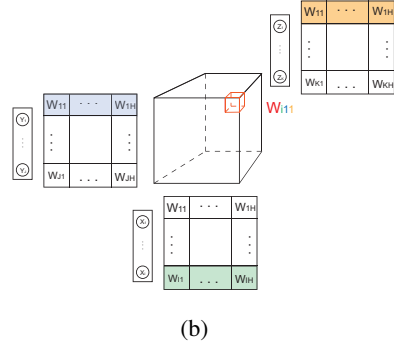
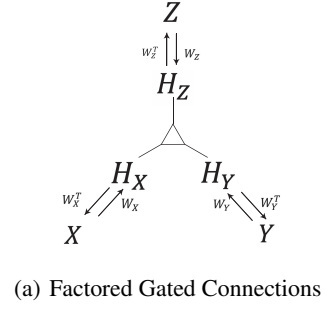


Figure 4: (a) Before being projected on the output layer Z , the projections of two given variables X and Y are multiplied in the intermediate factors layer. (b) Here is an instance of factorizing W_{I11} . After factorisation, the dimension of W_X , W_Y and W_Z are $H \times I$, $H \times J$, $H \times K$, corresponding to each form around the cubic.

2.2 Reachy Robot

As put forward in the introduction, Reachy robot, that will be employed and integrated in our experiment, is an open source human-like robot. With its totally 7-DoF human-like robotic arm, among which 3-DoF at shoulder level, 2-DoF for the wrist and 1-DoF of forearm and for upper arm respectively, Reachy's arm measures 60cm, similar to an adult's arm, allowing complex motions and postures in the 3D space. Note that the fifth motor operates forearm pronation-supination. Since Reachy have no lower limb so its peripersonal space depends only on the total length of his arm. For the sake of simplification, in the rest part, we consider that the shoulder has 3-DoF, elbow and wrist are both of 2-DoF as shown in Figure 5(b).

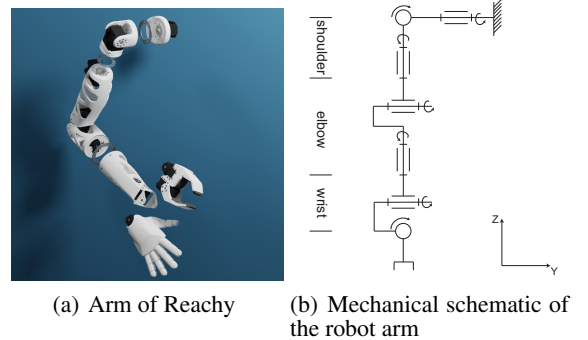


Figure 5: Architecture of Reachy

With regard to the equipped media device, Reachy provides also the high-level motor information about the joints's angular positions as well as mechanical parameters thanks to Robotis Dynamixel servomotors and embedded sensors. Inversely, Reachy can be then controlled through a serial port with a software interface called Pypot, a software written in Python which works out the communication with Dynamixel servomotors to drive the robot, e.g., sending motor commands, retrieving data from embedded sensors, which makes it possible to program a trajectory and manipulate the robot's movement through programming language, meeting the requirement of our experiment of reaching to a visually located target. Last but not least, with its two cameras installed inside its head as shown on the left side of Figure [6], one for observing the surrounding environment and another camera to focus on the task of manipulating, making it possible to get directly sensory information.

Attracted by these features, recently, some researchers have carried out several projects, getting the most out of the features of Reachy such as the inverse Kinematics problem on the exercise of minimisation of cost function, etc. However, these projects remain basic but they did inspire us to experiment our network on Reachy to make full use of its advantages like open source and extensive customization.

2.3 Task of Grasping in Gain-field Network

In [36], James Bonaiuto et al. came up with an Integrated Learning of Grasps and Affordances model when concerning reach-to-grasp movements of robots. However, since Reachy is equipped with a force sensor inside the gripper, so it has the ability to detect and adjust automatically the grip to maintain enough force to hold the object while not forcing too much and overheat the motor [37]. Besides, the graspable object of our experiment remains the same, a cylinder, hence there is no need for the virtual finger layer for our model.

Nevertheless, unlike the simple approaching to an object, the orientation of palms should be taken into account. As shown on the left side of Figure [6], with a preset identical displacement vector, the postures of robot may be different, namely, the object can be grasped and touched from many angles. Given different preferred orientations to grasping an object, as on the right of Figure 6, the parameters between each joint are different, so are the commands of motor. In [38], the author proves that when the wrist rotates, the preferred visual direction is discharged proportionally in order to compensate the muscles displacement. Apparently, hence the need for the additional variable fields, hence the need for additional parameters inside the model.

In order to better explain how robot reacts during the grasping movement, we introduce a structure based on the ILGM in [39], of which the dissimilarity is that there is no need for a Virtual Finger Layer. Note that certain relevant variables are the input of our network.

Direction Layer

Taking position s as the center of a target during reaching, the arm must reach a via-point around the object from which

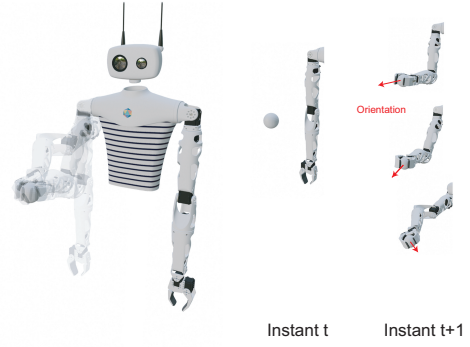


Figure 6: Same Visual Direction but Different Joint Angles and thus Different Motor Commands. The white ball is the object which stands also for the target point. Red arrows represent different orientations.

it moves toward the target. We define initial position by the vector sum $s + h$, among which the offset vector h , is the result of D layer. The parameter h is consisted of azimuth, elevation and radius components, has respectively a preferred value: azimuth, elevation, radius. The MG module associates the vector s with the offset h and then guides the arm to a via-point at $s+h$ and progressively to the target center point s .

Orientation Layer

The Orientation layer calculates the vector, r , which informs the fifth motor of the extent of the forearm movements, namely hand supination/pronation, to be made during the reach. Similar to the D layer, each unit has preferred hand supination/pronation values. The net movement parameter is calculated using a population-coding scheme as in the D layer. The relevant motor uses the output to rotate the hand via a linear interpolation between the current forearm angles and the value during reaching.

In practice, however, the vector h which is the output of the eye-centered gain-field network, as well as the vector r provided directly for the experiment.

Interestingly, if these grasping features are passed on to the gaze-centered gain-field network which learns to combine the information into representation of reaching for grasping the object thereby the relative layers specifies the hand approach direction relative to the object [39], and the representations are subsequently used to decide motor parameters for executing reaching and grasping activities, it enables the network to compute the optimal parameters and finally manipulate the grasping movement.

2.4 Experimental Setup

The experiments are composed of 3 different gain-field networks, corresponding to 3 different experiments. We will first illustrate the similarities between each experiment.

The experiment of reaching and grasping a located target through a neuronal network, of which the training is in terms

of a sequence of kinematic and dynamic activities, is actually an inverse kinetics problem against the backdrop of neuron network: instead of giving directly motor command in a mechanic way, the neuronal model will compute the coordinate transformations and return commands to the motor thereupon then it can be put in motion cognitively.

First of all, it is important to understand that the reference frame used for representing the target corresponds to which is applied for realizing reaching task because the brain determines appropriate axes to represent in space the configuration of the effector and the position of the object to bring one to the other. In practice, in these experiments, we chose hand-centered coordinates instead of should- head-centered coordinates for the reason that we do not take rotation of hand-eye into account and the saccade vector in hand- or shoulder-centered coordinates remains identical.

The eye-centered gain-field network takes advantage of the input furnished by two cameras installed inside the head of Reachy that play the role of eyes. An image of the hand of the robot and another photo for reachable target are two retina input in gaze-centered reference. The result of this gated network is the motor command in visual-centered frame. Since these input are in retina-centered reference, they can therefore be interpreted into eye-centered hand location HE and eye-centered target location TE , as shown below.

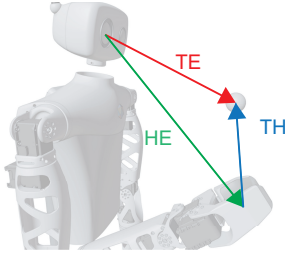


Figure 7: Vector TH : Hand-centered Target location, Vector HE : Eye-centered Hand location and vector TE : Eye-centered Target location.

Here, TH is hand-centered object location, which is regarded as desire movement vector, obtaining from mapping of motor command as well as HE and TE in the previous network. The association this hand-centered representation of target position $V(t)$ with the proprioceptive information $P(t)$ which consists of joints' parameters provided directly by Reachy is the input of the hand-centered gain-field network. The learned model ends up by returning motor command $A(t)$, as presented in Figure [8].

The combination of these two gain-field networks allows sensory stimuli to guide an effector, realising visuo-motor reaching in 3D space. The learning processes is done by using some exact but stochastic motor exploration which allows combining visual direction or preferred orientation with a specific posture provided by sensor. As soon as the different crucial and representative positions have been learned, the synaptic weights will be fixed, and the gain-field network can therefore easily interpolate for intermediate positions allowing the reaching of the visual position.

As described previously, Reachy provides the joints' pa-

rameters with its embedded sensors monitoring angular speed and position, serving as input X (or $P(t)$) of the hand-centered network, which are in the form of initial posture ($\varphi(s)$, $\varphi(e)$, $\varphi(w)$), corresponding to the joint angle of shoulder, elbow and wrist respectively. Note that shoulder joint rotates along three axes separately whereas elbow and wrist rotate only along two axes: axe Y and Z for elbow and axe Y and X for wrist. We adapt therefore a 1×7 matrix where each φ represents the degree of yaw, pitch and roll.

As for the output Z , the desired motor command $A(t)$ is also a vector C consists on 27 units corresponding to $3^3 = 27$ different motor synergies of the shoulder-elbow- wrist motion triplet $\{\Delta\varphi(s), \Delta\varphi(e), \Delta\varphi(w)\}$ whose discrete values are comprised in the small repertoire of three speeds $\{-1, 0, +1\}$; i.e., moving clockwise, release or moving anticlockwise.

2.4.1 Reaching with Given PD

The first experiment is to test whether learned network can produce appropriate trajectory with a given preferred direction. Input Y (or $V(t)$) in this network is a vector corresponding to a displacement vector in Cartesian hand-centered coordinates. The input layer is the output of the H layer mentioned previously, a 3 dimension vector corresponding to azimuth, elevation and radius value. In practice, we use 20 units of which their topology forms a three-dimensional mesh. The relation between joints' parameters and the degree of azimuth and elevation is shown below.

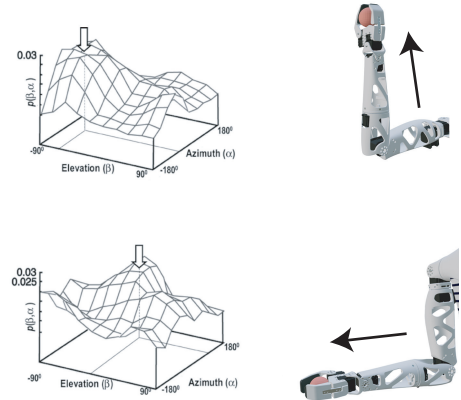


Figure 8: The correspondence between the region with the highest peak and the approach direction [39]

As a result, another input $P(t)$ is the initial posture, namely the parameters of each joint, which is also a vector, as discussed before. With the input as well as posture $P(t)$ mentioned earlier, the model will compute and return the motor command Z in the form of angular displacement.

2.4.2 Grasping with Given PO

The difference between this experiment and the previous one is that the second input is no longer the preferred direction but the desire palm orientation. However, since the network has

trained and learned from the critical posture, the gated parameters contain enough feature to generate an appropriate command.

Similarly, the input corresponds to the result of r layer. In our experiment, we use the effector as shown in Figure [6], with its embedded sensor, the wrist movement is composed of hand supination/pronation. Therefore, this input is a one dimension vector, or simply a scalar.

2.4.3 Reaching and Grasping with Given PD and PO

Apparently, different from the last two experiments, the dimensionality of input Y in this experiment has induced. Taken visual direction as well as the desire palm orientation into account, the $V(t)$ turns out to be a vector of 4 dimension.

3 Experiment

Acknowledgements

References

- [1] R. K. Clifton, D. W. Muir, D. H. Ashmead, and M. G. Clarkson, "Is visually guided reaching in early infancy a myth?," *Child development*, vol. 64, pp. 1099–1110, aug 1993.
- [2] D. Corbetta, S. L. Thurman, R. F. Wiener, Y. Guan, and J. L. Williams, "Mapping the feel of the arm with the sight of the object: on the embodied origins of infant reaching.," *Frontiers in psychology*, vol. 5, p. 576, 2014.
- [3] J. F. X. DeSouza, G. P. Keith, X. Yan, G. Blohm, H. Wang, and J. D. Crawford, "Intrinsic reference frames of superior colliculus visuomotor receptive fields during head-unrestrained gaze shifts.," *The Journal of neuroscience : the official journal of the Society for Neuroscience*, vol. 31, pp. 18313–18326, dec 2011.
- [4] S. Chivukula, M. Jafari, T. Aflalo, N. A. Yong, and N. Pouratian, "Cognition in Sensorimotor Control: Interfacing With the Posterior Parietal Cortex," *Frontiers in Neuroscience*, vol. 13, p. 140, 2019.
- [5] M. Flanders, S. I. H. Tillery, and J. F. Soechting, "Early stages in a sensorimotor transformation," *Behavioral and Brain Sciences*, vol. 15, no. 2, pp. 309–320, 1992.
- [6] M. F. Ghilardi, J. Gordon, and C. Ghez, "Learning a visuomotor transformation in a local area of work space produces directional biases in other areas.," *Journal of neurophysiology*, vol. 73, pp. 2535–2539, jun 1995.
- [7] P. Vindras and P. Viviani, "Frames of reference and control parameters in visuomanual pointing.," 1998.
- [8] J. McIntyre, F. Stratta, and F. Lacquaniti, "Viewer-centered frame of reference for pointing to memorized targets in three-dimensional space.," *Journal of neurophysiology*, vol. 78, pp. 1601–1618, sep 1997.
- [9] M. Carrozzo, J. McIntyre, M. Zago, and F. Lacquaniti, "Viewer-centered and body-centered frames of reference in direct visuomotor transformations.," *Experimental brain research*, vol. 129, pp. 201–210, nov 1999.
- [10] E. Ugur, Y. Nagai, H. Celikkanat, and E. Oztop, "Parental scaffolding as a bootstrapping mechanism for learning grasp affordances and imitation skills," *Robotica*, vol. 33, no. 5, pp. 1163–1180, 2015.
- [11] A. P. Georgopoulos, H. Merchant, T. Naselaris, and B. Amirikian, "Mapping of the preferred direction in the motor cortex," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 104, no. 26, pp. 11068–11072, 2007.
- [12] E. Salinas and P. Thier, "Gain modulation: a major computational principle of the central nervous system.," jul 2000.
- [13] A. Pouget and T. J. Sejnowski, "Spatial transformations in the parietal cortex using basis functions.," *Journal of cognitive neuroscience*, vol. 9, pp. 222–237, mar 1997.
- [14] S. Deneve and A. Pouget, "Bayesian multisensory integration and cross-modal spatial links.," *Journal of physiology, Paris*, vol. 98, no. 1-3, pp. 249–258, 2004.
- [15] M. A. Smith and J. D. Crawford, "Distributed population mechanism for the 3-D oculomotor reference frame transformation," *Journal of Neurophysiology*, vol. 93, no. 3, pp. 1742–1761, 2005.
- [16] J. Abrossimoff, A. Pitti, and P. Gaussier, "Visuo-Motor Control Using Body Representation of a Robotic Arm with Gated Auto-Encoders," 2019.
- [17] M. I. Jordan and D. E. Rumelhart, "Forward models: Supervised learning with a distal teacher," *Cognitive Science*, vol. 16, no. 3, pp. 307–354, 1992.
- [18] D. Bullock, S. Grossberg, and F. H. Guenther, "A Self-Organizing Neural Model of Motor Equivalent Reaching and Tool Use by a Multijoint Arm," *Journal of Cognitive Neuroscience*, vol. 5, no. 4, pp. 408–435, 1993.
- [19] D. M. Wolpert, K. Doya, and M. Kawato, "A unifying computational framework for motor control and social interaction.," *Philosophical transactions of the Royal Society of London. Series B, Biological sciences*, vol. 358, pp. 593–602, mar 2003.
- [20] G. Schillaci, B. Lara, and V. V. Hafner, "Internal simulations for behaviour selection and recognition," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 7559 LNCS, no. October, pp. 148–160, 2012.

- [21] G. Schillaci, V. V. Hafner, and B. Lara, "Online learning of visuo-motor coordination in a humanoid robot. A biologically inspired model," *IEEE ICDL-EPIROB 2014 - 4th Joint IEEE International Conference on Development and Learning and on Epigenetic Robotics*, pp. 130–136, 2014.
- [22] B. Moldovan, P. Moreno, M. Van Otterlo, J. Santos-Victor, and L. De Raedt, "Learning relational affordance models for robots in multi-object manipulation tasks," *Proceedings - IEEE International Conference on Robotics and Automation*, pp. 4373–4378, 2012.
- [23] J. M. Galeazzi, B. M. Mender, M. Paredes, J. M. Tromans, B. D. Evans, L. Minini, and S. M. Stringer, "A Self-Organizing Model of the Visual Development of Hand-Centred Representations," *PLoS ONE*, vol. 8, no. 6, 2013.
- [24] E. Ugur and J. Piater, "Emergent structuring of interdependent affordance learning tasks," *IEEE ICDL-EPIROB 2014 - 4th Joint IEEE International Conference on Development and Learning and on Epigenetic Robotics*, pp. 489–494, 2014.
- [25] A. Roncone, M. Hoffmann, U. Pattacini, and G. Metta, "Automatic kinematic chain calibration using artificial skin: Self-touch in the iCub humanoid robot," *Proceedings - IEEE International Conference on Robotics and Automation*, pp. 2305–2312, 2014.
- [26] J. Born, J. M. Galeazzi, and S. M. Stringer, "Hebbian learning of hand-centred representations in a hierarchical neural network model of the primate visual system," *PLOS ONE*, vol. 12, no. 5, pp. 1–35, 2017.
- [27] P. Lanillos, E. Dean-Leon, and G. Cheng, "Yielding Self-Perception in Robots Through Sensorimotor Contingencies," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 9, pp. 100–112, jun 2017.
- [28] A. Droniou, S. Ivaldi, and O. Sigaud, "Deep unsupervised network for multimodal perception, representation and classification," *Robotics and Autonomous Systems*, vol. 71, pp. 83–98, 2015.
- [29] O. Sigaud and A. Droniou, "Towards Deep Developmental Learning," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 8, no. 2, pp. 99–114, 2015.
- [30] A. Pitti, A. Blanchard, M. Cardinaux, and P. Gaussier, "Gain-field modulation mechanism in multimodal networks for spatial perception," *IEEE-RAS International Conference on Humanoid Robots*, no. August 2015, pp. 297–302, 2012.
- [31] S. Mahé, R. Braud, P. Gaussier, M. Quoy, and A. Pitti, "Exploiting the gain-modulation mechanism in parieto-motor neurons: Application to visuomotor transformations and embodied simulation," *Neural Networks*, vol. 62, pp. 102–111, 2015.
- [32] R. Braud, A. Pitti, and P. Gaussier, "A Modular Dynamic Sensorimotor Model for Affordances Learning, Sequences Planning, and Tool-Use," *IEEE Transactions on Cognitive and Developmental Systems*, vol. 10, pp. 72–87, mar 2018.
- [33] R. Memisevic, "Learning to relate images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1829–1846, 2013.
- [34] P. Smolensky, "Tensor product variable binding and the representation of symbolic structures in connectionist systems," *Artificial Intelligence*, vol. 46, no. 1, pp. 159–216, 1990.
- [35] R. Memisevic and G. E. Hinton, "Learning to represent spatial transformations with factored higher-Order boltzmann machines," *Neural Computation*, vol. 22, no. 6, pp. 1473–1492, 2010.
- [36] J. Abrossimoff, A. Pitti, P. Gaussier, J. Abrossimoff, A. Pitti, P. Gaussier, V. Learning, J. Abrossimoff, A. Pitti, and P. Gaussier, "Visual Learning for Reaching and Body-Schema with Gain-Field Networks To cite this version : HAL Id : hal-01976669 Visual Learning for Reaching and Body-Schema with Gain-Field Networks," 2019.
- [37] S. Mick, M. Lapeyre, P. Rouanet, C. Halgand, J. Benois-Pineau, F. Paclet, D. Cattaert, P. Y. Oudeyer, and A. De Rugy, "Reachy, a 3D-printed human-like robotic arm as a testbed for human-robot control strategies," *Frontiers in Neurorobotics*, vol. 13, no. August, pp. 1–12, 2019.
- [38] J. Bonaiuto and M. Arbib, "Learning to grasp and extract affordances: the Integrated Learning of Grasps and Affordances (ILGA) model," *Biological Cybernetics*, vol. 109, 2015.
- [39] E. Oztop, N. S. Bradley, and M. A. Arbib, "Infant grasp learning: A computational model," *Experimental Brain Research*, vol. 158, no. 4, pp. 480–503, 2004.