

# Cognitive Architecture for Mutual Modelling

Alexis Jacq<sup>1,2</sup>, Wafa Johal<sup>1</sup>, Pierre Dillenbourg<sup>1</sup>, Ana Paiva<sup>2</sup>

<sup>1</sup>CHILI Lab, École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

<sup>2</sup>INESC-ID & Instituto Superior Técnico, University of Lisbon, Portugal

**Abstract**—In social robotics, robots need to be able to be understood by humans. Especially in collaborative tasks where they have to share mutual knowledge. For instance, in an educative scenario, learners share their knowledge and they must adapt their behaviour in order to make sure they are understood by others. Learners display behaviours in order to show their understanding and teachers adapt in order to make sure that the learners' knowledge is the required one. This ability requires a model of their own mental states perceived by others: “*has the human understood that I(robot) need this object for the task or should I explain it once again ?*” In this paper, we discuss the importance of a cognitive architecture enabling second-order Mutual Modelling for Human-Robot Interaction in educative contexts.

## I. INTRODUCTION

A social robot is brought to interact with humans. The quality of this interaction depends on its ability to behave in an acceptable and understandable manner by the user. Hence the importance for a robot to take care of his image: how much it is perceived as an automatic and repetitive agent, or contrarily as a surprising and intelligent character. If the robot is able to detect this perception of itself, it can adapt its behaviour in order to be understood: “you think I am sad while I am happy, I want you to understand that I am happy”.

In a collaborative context, where knowledge must be shared, agents must exhibit that they are acquiring the shared information with an immediate behaviour: “I look at what you are showing me, do you see that I am looking at it, do you think I am paying attention to your explanation ?”; “I have understood your idea, do you understand that I have understood ?”. As humans, we have different strategies to exhibit understanding or to resolve a misunderstanding. As an example, if someone is talking about a visual object, we alternatively gaze between the object and the person to make sure he saw that we gazed at the object. Or if we detect that the other person has not understood a gesture (e.g. pointing at an object) we would probably exaggerate the gesture.

Developed by Baron-Cohen and Leslie [1], the Theory of Mind (ToM) describes the ability to attribute mental states and knowledge to others. In interaction, humans are permanently collecting and analysing huge quantity of information to stay aware of emotions, goals and understandings of their fellows. In this work, we focus on a generalization of this notion: Mutual Modelling characterizes the effort of one agent to model the mental state of another one [2].

Until now, the work conducted by the Human-Robot Interaction (HRI) community to develop mutual modelling abilities in robots was limited to a first level of modelling (see

related work in section II). Higher levels require the ability to recursively attribute a theory of mind to other agents (*I think that you think that ...*) and their application to HRI remains unexplored. However, a knowledge of oneself perceived by others is necessary to adapt a behaviour to keep mutual understandings.

An important challenge of social robotics is to provide assistance in education. The ability of robots to support adaptive and repetitive tasks can be valuable in a learning interaction. The CoWriter Project [3], [4] introduces a new approach to help children with difficulties in learning handwriting. Based on the *learning by teaching* paradigm, the goal of the project is not only to help children with their handwriting, but mainly to improve their self-confidence and motivation in practising such exercise.

*Learning by teaching* engages students to conduct the activity in the role of the teachers in order to support their learning process. This paradigm is known to produce motivational, meta-cognitive and educational benefits in a range of disciplines [5]. The CoWriter project is the first application of the learning by teaching approach to handwriting.

The effectiveness of this learning by teaching activity is built on the “protégé effect”: the teacher feels responsible for his student, commits to the student’s success and possibly experiences student’s failure as his own failure to teach. The main idea is to promote the child’s extrinsic motivation to write letters (he does it in order to help his “protégé” robot) and to reinforce the self-esteem of the child (he plays the teacher and the robot actually progresses).

In that context, the robot needs to pretend enough difficulties to motivate the child to help it. This ability of the robot to pretend strongly depends on the perception of the robot by the child: the displayed behaviours (gestures, gazes and sounds) by the robot, the initial level and learning speed of the robot must match with what the child imagines of a “robot in difficulty”. In order to adapt to the child, the robot needs then to have a model of how it is perceived by the child. On the other side, the child builds also a model of the robot’s difficulties and attitude. This mutual-modelling is primordial in order to have mutual understanding and fluid interaction between learner and teacher.

## II. RELATED WORKS

A large amount of fields have introduced frameworks to describe mutual modelling ability [6]. In developmental psychology Flavell [7] denotes two different levels of perspective taking: the *cognitive connection* (I see, I hear, I want, I

like...) and *mental representation* (what other agents feel, hear, want...).

From a computational perspective, Epistemic logic describes knowledges and beliefs shared by agents. This framework enables consideration of infinite-level of mutual modelling. It defines a *shared-knowledge* (all the agents of a group know  $\mathbf{X}$ ) and a *common-knowledge* (all the agents of a group know  $\mathbf{X}$ , and know that all the agent know  $\mathbf{X}$ , and know that all the agents know that all the agents know  $\mathbf{X}$ , ...) [8].

Mutual modelling has also been studied through educational contexts. Roschelle and Teasley [9] suggested that collaborative learning requires a *shared understanding* of the task and of the shared information to solve it. The term “mutual modelling” was introduced in Computer-Supported Collaborative Learning (CSCL) by Dillenbourg [2]. It focused on knowledge states of agents. Dillenbourg developed in [10] a computational framework to represent mutual modelling situations.

However, HRI research has not, until now, explored the whole potential of mutual modelling. In [11], Scassellati supported the importance of Leslie’s and Baron-Cohen’s theory of mind to be implemented as an ability for robots. He focused his work on attention and perceptual processes (face detection or colour saliency detection). Thereafter, some works (including Breazeal [12], Trafton [13], Ros [14] and Lemaignan [15]) were conducted to implement Flavell’s first level of perspective taking [16] (“*I see (you do not see the book)*”), ability that is still limited to visual perception.

Breazeal [17] and Warnier [18] reproduced the Sally and Anne’s test of Wimmer [19] with robots able to perform visual perspective taking. The robot was able to infer the knowledge of a human given the history of his visual experience.

In [20], Lemaignan implemented a system that computes the visual field of agents and estimates which objects are looked at in real-time. This time, the robot is not just aware of what *can be seen* by agents, but it perceives what *is currently being looked at*. Lemaignan used this system to measure Sharma’s *with-me-ness* [21], visual commitment based on expected focus of attention in an activity.

### III. MM-BASED REASONING

A first intuition for mutual modelling is to assume that all agents have the same basic architecture. In [12], Breazeal show a MM-based reasoning where the robot uses its own architecture to model other agents. We can imagine a second level of modelling where the robot recursively attribute to other agents the mutual modelling ability. But it would create an infinite recursive loop: the agent then models the robot that models the agent etc. Another reason to avoid a recursive approach is that different agent must have different behaviours: in similar situations, they will not necessary take similar decisions.

We propose a different approach of modelling, where we define two orders of agents: the *first-order-agents* deal with direct representations of agents by the robot (for example the child), while the *second-order-agents* deal with the representation

of agents by agents (for example the *robot-perceived-by-the-child*). Modelling *second-order-agent* like the *robot-perceived-by-the-child* will help to model how the child perceives the robot, e.g. to make sure the child understands that the robot is learning from his demonstrations. We can also define  $n^{th}$ -order agents with a higher level of theory of mind. But taking into account high levels of mutual modelling would be difficult to process in real time. Unlike the epistemic logic, our proposed framework will not take into account infinite regress [22] of mutual modelling.

All sensors (cameras, micros, motor positions etc. and in the case of CoWriter the tablet’s inputs) are used to perceive information about the physical behaviour of agents. We call all the measurable quantities or qualities that provide information like position in space, the direction of the gaze, speech, movement and facial expressions etc. as *perceived variables*. Each agent’s model is associated with a set of perceived variables that describes his physical behaviour.

Emotional states of agents cannot be directly measured directly from sensors. We call *abstract variables* all the quantities or qualities that describe the mental state of an agent. Abstract variables are deduced from the dynamic of perceived variables. As an example, if the robot points at an object with its arm, it expects the child to look at the object. If then the child looks at the hand of the robot, the robot can deduce that the child has not understood the meaning of its gesture. The perceived variables are the robot’s gesture and the gaze direction of the child. The deduced abstract variable is the understanding of the gesture by the child.

A model of an agent is the set of all the values of the perceived or abstract variables associated with this model. Since the values of variables are likely to change with the time, the models must be dynamic.

In order to deduce the values of abstract variables (that can’t be obtained from direct perception), we propose to build a Bayesian model based on the knowledge from perceived variables. The choice of a probabilist approach instead of a symbolic approach comes from the errors in the perception of the robot: knowledge and other mental states of agents can not be directly perceived through the behaviours. They must be inferred, hence a probabilistic model enables richer predictions.

This Bayesian network would contain the probabilities that *abstract variables* take values given the values of *perceived variables*.

For example, if the robot points at an object and detects that the child saw the movement, then it expects the child to look at the targeted object immediately after. In other terms, if the child looks at the hand of the robot but does not look at the target object, the probability that the child understood the pointing movement is expected to be small. Knowing that, the robot can make the decision to exaggerate its pointing gesture.

### IV. DESCRIPTION OF THE ARCHITECTURE

The picture 1 visually summarizes the global design of our architecture.

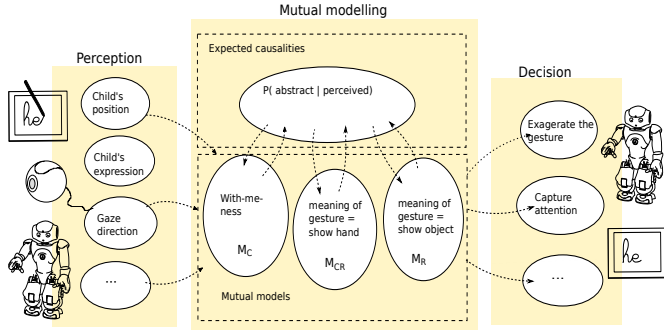


Fig. 1. **Overview of the cognitive architecture.** Yellow squares represent the main parts. White ellipses represent modules. It shows possible devices used for perception and decisions with the context of the CoWriter Activity. We illustrate the architecture with a situation of misunderstanding : the child has a bad interpretation of the gesture of the robot. In order to resolve this misunderstanding, a possible decision could be to exaggerate the movement.

Our cognitive architecture for mutual modelling contains three main parts. The **perception part** (see IV-A) regroups all the modules that measure the values of the perceived variables using sensors. These values are sent to the **mutual modelling part** (see IV-B) that updates models with measured values of perceived variables and infers the value of abstract variables in real-time. Finally, the **decision part** (see IV-C) contains all the modules associated to the control of the robot (and other active devices like tablets in CoWriter). These modules can read values given by mutual models in order to compute decisions. In the example of the CoWriter activity, these modules are given by the system that learns and generate letters, but we can add a module that generates micro-behaviours, another that decides to switch to a new activity (e.g. drawing with the robot),... The following subsections explain in detail the content and operation of each part of the architecture.

#### A. Perception modules

The sensitive modules measure values of relevant perceived variables. While the agent's gaze direction and facial expression can be used in any interaction, some additive variables can be specific to the activity: in CoWriter, a module takes as input perceived variables from a tablet to compute the new state of robot's writing. It defines a sensitive module, and the value of the new state of the robot to write a letter defines a perceived variable. The evaluation of the robot by the child via the feedback buttons on the tablet defines another perceived variable provided by the modules of the activity. Other modules are independent of the activity: the system that estimates the target objects looked at by the child provides additive information not directly used by the modules of the activity.

#### B. Mutual modelling modules

Each perceived value measured by the sensitive modules are associated with the model of an agent (or a  $n^{th}$ -order agent). Each mutual model can be designed as a module that deals with a list of associated perceived variables and watch if the value of one of these variables has been changed. An additive module knows all the expected causalities and computes the values of abstract variables. Some expected causalities can be empirically learned and others pre-programmed.

#### C. Decision making

The values of mutual modelling variables will provide rich and useful information for decision making. Taking in account these values to elaborate decision should improve the realism and the efficiency of the robot in the interaction. Similarly to the sensitive ones, the modules that make decisions can be specific to the activity (in CoWriter, the choice of a new learning curve or the decision to suddenly make a big mistake), or can govern a general behaviour (for example the exaggeration of a misunderstood gesture). Some decisions can have a high impact on the interaction: to stop an activity and to switch to a new one can frustrate a child that was committed. The conditions to make such a decision are not directly assessable, but must be learned by the robot. In order to make these decision cautiously, we propose to start by a Wizard-of-Oz approach and to move towards an autonomous approach following these steps:

- 1) **Wizard-of-Oz:** A human takes decisions; the robot learns
- 2) **Mixed-initiative:** The robot makes suggestions; a human agrees or disagrees
- 3) **Autonomous:** The robot makes decisions

#### V. CONCLUSION

Educational HRI based on learning by teaching approach needs robots to be able to perform second-level mutual modelling. We introduced a new approach to implement mutual modelling into a cognitive architecture. We used the CoWriter activity as an example of application, but our architecture could be easily generalised for any kind of interaction. We believe that this step must be reached in other contexts of HRI, in order to develop higher realism of behaviours and to improve the quality of interactions.

#### ACKNOWLEDGMENT

This research was partially supported by the Fundação para a Ciência e a Tecnologia (FCT) with reference UID/CEC/50021/2013, and by the Swiss National Science Foundation through the National Centre of Competence in Research Robotics.

#### REFERENCES

- [1] S. Baron-Cohen, A. Leslie, and U. Frith, "Does the autistic child have a "theory of mind" ?" *Cognition*, 1985.
- [2] P. Dillenbourg, "What do you mean by collaborative learning?" *Collaborative-learning: Cognitive and Computational Approaches.*, pp. 1-19, 1999.

- [3] D. Hood, S. Lemaignan, and P. Dillenbourg, "When children teach a robot to write: An autonomous teachable humanoid which uses simulated handwriting," in *Proceedings of the Tenth Annual ACM/IEEE International Conference on Human-Robot Interaction*, ser. HRI '15. New York, NY, USA: ACM, 2015, pp. 83–90.
- [4] A. Jacq, S. Lemaignan, F. Garcia, P. Dillenbourg, and A. Paiva, "Building successful long child-robot interactions in a learning context," in *Proceedings of the 2016 ACM/IEEE Human-Robot Interaction Conference*, 2016.
- [5] C. A. Rohrbeck, M. D. Ginsburg-Block, J. W. Fantuzzo, and T. R. Miller, "Peer-assisted learning interventions with elementary school students: A meta-analytic review," *Journal of Educational Psychology*, vol. 95, no. 2, pp. 240–257, 2003. [Online]. Available: <http://doi.apa.org/getdoi.cfm?doi=10.1037/0022-0663.95.2.240>
- [6] S. Lemaignan and P. Dillenbourg, "Mutual modelling in robotics: Inspirations for the next steps," in *Proceedings of the 2015 ACM/IEEE Human-Robot Interaction Conference*, 2015.
- [7] J. H. Flavell, F. L. Green, and E. R. Flavell, "Developmental changes in young children's knowledge about the mind," *Cognitive Development*, vol. 5, no. 1, pp. 1–27, 1990.
- [8] V. Hendricks and J. Symons, "Epistemic logic," in *Stanford Encyclopedia of Philosophy*, 2008.
- [9] J. Roschelle and S. D. Teasley, "The construction of shared knowledge in collaborative problem solving," in *Computer supported collaborative learning*. Springer, 1995, pp. 69–97.
- [10] M. Sangin, N. Nova, G. Molinari, and P. Dillenbourg, "Partner modeling is mutual," in *Proceedings of the 8th international conference on Computer Supported Collaborative Learning*. International Society of the Learning Sciences, 2007, pp. 625–632.
- [11] B. Scassellati, "Theory of mind for a humanoid robot," *Autonomous Robots*, vol. 12, no. 1, pp. 13–24, 2002.
- [12] C. Breazeal, M. Berlin, A. Brooks, J. Gray, and A. Thomaz, "Using perspective taking to learn from ambiguous demonstrations," *Robotics and Autonomous Systems*, pp. 385–393, 2006.
- [13] J. Trafton, N. Cassimatis, M. Bugajska, D. Brock, F. Mintz, and A. Schultz, "Enabling effective human-robot interaction using perspective-taking in robots," *IEEE Transactions on Systems, Man and Cybernetics, Part A: Systems and Humans*, vol. 35, no. 4, pp. 460–470, 2005.
- [14] R. Ros, E. A. Sisbot, R. Alami, J. Steinwender, K. Hamann, and F. Warneken, "Solving ambiguities with perspective taking," in *5th ACM/IEEE International Conference on Human-Robot Interaction*, 2010.
- [15] S. Lemaignan, "Grounding the interaction: Knowledge management for interactive robots," Ph.D. dissertation, CNRS - Laboratoire d'Analyse et d'Architecture des Systèmes, Technische Universität München - Intelligent Autonomous Systems lab, 2012.
- [16] J. H. Flavell, "The development of knowledge about visual perception," *Nebraska Symposium on Motivation*, vol. 25, pp. 43–76, 1977.
- [17] C. Breazeal, J. Gray, and M. Berlin, "An embodied cognition approach to mindreading skills for socially intelligent robots," *The International Journal of Robotics Research*, vol. 28, no. 5, pp. 656–680, 2009.
- [18] M. Warnier, J. Guitton, S. Lemaignan, and R. Alami, "When the robot puts itself in your shoes. managing and exploiting human and robot beliefs," in *Proceedings of the 21st IEEE International Symposium on Robot and Human Interactive Communication*, 2012, pp. 948–954.
- [19] H. Wimmer and J. Perner, "Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children's understanding of deception," *Cognition*, vol. 13, no. 1, pp. 103–128, 1983.
- [20] S. Lemaignan, F. Garcia, A. Jacq, and P. Dillenbourg, "From real-time attention assessment to with-me-ness in human-robot interaction," in *Proceedings of the 2016 ACM/IEEE Human-Robot Interaction Conference*, 2016.
- [21] K. Sharma, P. Jermann, and P. Dillenbourg, "with-me-ness: A gaze-measure for students attention in moocs," in *International conference of the learning sciences*, no. EPFL-CONF-201918, 2014.
- [22] H. H. Clark and S. E. Brennan, "Grounding in communication," *Perspectives on socially shared cognition*, vol. 13, no. 1991, pp. 127–149, 1991.