

Effect of a Humanoid’s Active Role during Learning with Embodied Dialogue System

Matthias Kerzel, Hwei Geok Ng, Sascha Griffiths and Stefan Wermter

Knowledge Technology Institute

Department of Informatics

Universität Hamburg

Vogt-Kölln-Str. 30

22527 Hamburg, Germany

Email: {Kerzel, 5ng, griffiths, wermter}@informatik.uni-hamburg.de

Abstract—When humanoid robots learn complex sensorimotor abilities from interaction with the environment, often a human experimenter is required. For a social companion robot, it is desirable that the learning can also be assisted by non-expert users. To achieve this aim, we present an embodied dialogue system which enables a humanoid to take on an active role during learning by guiding its user with verbal communication and through the display of emotions. We suggest an experimental setup for evaluating how the active role affects the learning result and the subjective evaluation of the humanoid by human participants.

I. INTRODUCTION

Humanoid robots are designed to operate alongside or together with humans in complex and unstructured everyday environments. In contrast to industrial robots, their tasks require adaptation to novel challenges that can be overcome through learning. State-of-the-art learning approaches, such as deep neural networks [1], [2], rely on a large quantity of training data. Apart from using costly human annotation, these data can be gathered through interaction with the environment [3]–[5]. Similar to a child, a humanoid can incrementally develop complex visuomotor skills through interacting in the physical world [6]. Due to their human-like appearance, humanoids can be assisted in this learning by non-expert users in a similar way these users would teach a child. This proposal of child-like learning for successful artificial intelligence goes back to Alan Turing [7]. Recent work, however, has shown that especially in teaching robots, people find it easier to teach a robot if it behaves similarly to a child [8]. Human assistance is often useful, as a humanoid occasionally needs help during the process, e.g. when objects roll off tables or are out of reach.

However, there are large differences between a child and a robot. While a child needs to develop all of its cognitive abilities, we can employ a mixed approach for a humanoid: some abilities develop through learning and interaction, other abilities are designed. Moreover, these capabilities can be used to enable an active role of the humanoid in learning.

In this paper, we present initial work on an embodied dialogue system that enables a humanoid to guide users through the training steps for developing visuomotor abilities, see Figure 1. We aim to evaluate the feasibility of the approach in guiding non-expert users to successfully collect training

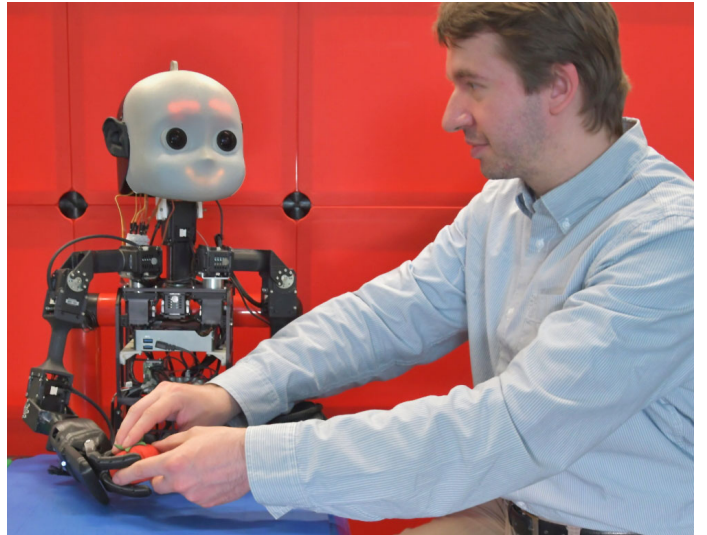


Fig. 1. A robot asks for help from a human assistant to learn grasping.

data and also to observe the effect of the active role of the humanoid with regard to the users’ subjective perception of the humanoid. By integrating a state-of-the-art approach for sensorimotor learning into our human-robot interaction (HRI) research, we increase the realism of the user interaction. Participants will take the role of a teacher of a physical robot in a real developmental task.

II. RELATED WORK

A. Learning Visuomotor Skills

Acquiring visuomotor skills with neural network approaches has gained interest in recent years. Mnih et al. [9] applied deep reinforcement learning for human-level control in computer games. Lillicrap et al. [10] extended the approach to continuous deep reinforcement learning for simplified robot arms in virtual environments. Neural reinforcement learning approaches rely on a large number of trial and error attempts to solve a task. They are successful in virtual environments, which can provide a large number of training samples in a short time at a low cost, both in terms of human supervision and damage to the robot (by wear or accidents).

When adapting neural learning approaches to physical robots, the time it takes the robot to perform actions becomes a critical factor. As shown by Gupta and Pinto [4], it takes a robot 700 hours of training to learn grasping positions and angles. For non-industrial robots, this long training time often exceeds the life-expectancy of hardware components. Levine et al. [3] and Kerzel and Wermter [5] suggest approaches for reducing the necessary training time by transforming the reinforcement learning into a supervised learning task. Both approaches rely on generating annotated training data, i.e. it is not learned through trial and error, but only from correct examples which leads to a shorter learning time. Levine et al. employ computation of forward kinematics to make the robot's state fully observable during training time. Kerzel and Wermter use the robot's ability to autonomously place objects to generate training samples. Though both of these approaches combine the advantages of mostly autonomous reinforcement learning with the short training time of supervised learning, both approaches ultimately rely on human experts for initiating the training and assisting the robot in case of errors.

Cruz et al. [11] have demonstrated that learning skills via accordances is improved by interactive learning versus regular reinforcement learning. Especially, speech and multimodal feedback are useful for adding interactivity to the learning process [12].

B. Spoken Dialogue Systems

Spoken dialogue systems (SDS) [13] are modules in HRI systems which receive speech as input and produce the corresponding replies [14]. An SDS solves five main tasks: Automatic Speech Recognition, Spoken Language Understanding, Dialogue Management (DM), Natural Language Generation and Text-to-Speech Synthesis [14]. In this paper, we focus on the integration of robot perception into the Dialogue Management to facilitate learning of visuomotor skills. DM is a decision-maker in SDS: It integrates information about the previous dialogue, internal states of the conversation agent, the robot's perception and agenda to decide on actions – which among others can be spoken utterances or motor actions. There are two main types of the dialogue systems: reactive and agenda-driven systems [15]. Reactive systems generate interactive responses based on what the user said with the purpose of producing a meaningful conversation. Agenda-driven dialogue systems do not produce output responses merely from the user's inputs but changes the context of the conversation to achieve its goals [15], which in our case is the realization of sensorimotor ability learning phases.

The variable and open-ended nature of language is now making data-driven methods more prominent in spoken dialogue systems [16] with deep learning approaches now also being explored [17]. However, state-of-the-art dialogue systems for HRI are often still built using a pipeline of tools and are mostly symbolic approaches to a large extent [18], [19]. Knowledge-based approaches are still efficient in restricted domains when data is not very variable [20] and the cost-benefit ratio of collecting and annotating data does not lend

itself to a data-dependent approach.

III. RESEARCH QUESTIONS AND EXPERIMENTAL SETUP

We aim to evaluate the effect of an active role during the learning of a humanoid robot that is assisted by non-expert users. In our experimental setup for grasp learning, a child-sized humanoid interacts with non-expert participants to learn sensorimotor skills for grasping. The grasp learning follows Kerzel and Wermter [5]: The robot interacts with an object on a table. It looks at, grasps and places an object repeatedly to gather samples for training its artificial neural architecture.

This learning process is only partially autonomous: Initially, the robot can move its hand to random positions on the table surface. Once the robot has an object in its hand, it places the object at a randomly chosen position and then associates the joint values during placing the object with how the scene looks after placement. Thus, the robot collects samples that link motor actions to visual inputs for developing hand-eye coordination. The robot, however, needs human assistance during the process: It only develops the ability to link its vision to its actions during training - it lacks this capability at first. Thus, a human assistant must initially place an object into the robot's hand. The robot will then begin its autonomous learning cycle of placing and picking up the object at random positions. However, it can happen that the robot accidentally moves the object during grasping or releasing. This leads to failed grasps where human assistance is needed.

We will realize the scenario in two conditions: In the *human-guided* condition, we will establish a baseline for non-active learning. A human experimenter will introduce the robot and the learning scenario, explain all the steps to the user, openly operate the robot by executing command-line programs and alert the user to situations that require assistance. To avoid influencing the participant, pre-recorded instructions in the same voice as used for the robot can be used. The robot will remain silent and perform the necessary actions for learning. It will however express emotions on its face to equalize the experimental conditions.

In the *robot-guided* condition, the robot will take on an active role: It will greet the participant and introduce the learning scenario; it will comment verbally on its actions and also use emotion expressions to indicate success or problems. All of the robot's actions will run autonomously, controlled by the embodied dialogue system. Both scenarios offer an authentic human-robot interaction based on a state-of-the-art neural deep learning approach. Our research is guided by three main questions:

- Does the embodied dialogue system enable non-expert users to undergo training with the humanoid? We will evaluate this objectively by comparing the number of successfully collected training samples in the *robot-guided* and *human-guided* experimental conditions.
- How comfortable and easy does the learning scenario feel to users? We will identify possible problems with a questionnaire and use this information to improve the

human-robot interaction and the learning scenario in future iterations.

- How does the humanoid’s active role in learning influence the user’s subjective perception of the robot? For this, we will employ the established GODSPEED test from Bartneck et al. [21] for both experimental conditions. We will compare these results to former studies that involved the same humanoid [22].

IV. NICO ROBOT AND GRASP LEARNING

We realize the experimental setup with NICO (Neuro-Inspired COmpanion), introduced by Kerzel et al. [23]. NICO’s child-like design is aimed to elicit a high user acceptance and make users intuitively adopt the role of a teacher.

NICO is primarily endowed with capabilities for human-like perception and interaction as well as object grasping and manipulation: Its two arms have six degrees of freedom and a human-like range of motion. NICO can grasp and manipulate small objects with its three-fingered hands. The fingers utilize a tendon mechanism which enables them to wrap themselves around objects of various sizes. Additionally, the state of tendons can be used as haptic feedback to evaluate if a grasping action has been successful.

NICO’s head can perform tilt and yaw movements; it features two cameras and two microphones. The child-like design of the head is adapted from the iCub [24]. The head features LED arrays around the mouth and eye regions that display stylized facial expressions [22].

A. Neural Grasp Learning

Visuomotor skills are acquired by associating a state of the environment with the desired action. In this paper, we follow the approach by Kerzel and Wermter [5], where the state of the environment is represented by images from the humanoid’s cameras and the action equals a joint configuration that moves the humanoid’s arm into a grasp position.

This association is facilitated by a deep neural network that can generalize from a limited number of training samples. The network architecture consists of two convolution layers that process input from the two cameras in the humanoid’s head and two dense layers that further transform this input into a joint configuration. The neural network is trained end-to-end, i.e. the training data consists of images and the output is the corresponding joint configuration for grasping the objects.

The training data is collected in a semi-autonomous training cycle, as shown in Figure 2. After the training object is placed in the humanoid’s hand, it moved to a random position on the table. The joint configuration that leads to this placement is memorized. The humanoid releases the object and moves the hand to the side to record images. These images are saved along with the memorized joint configuration to form one training data point. The humanoid then moves back to the memorized joint configuration to grasp the object again. If the grasp attempt was successful, the training cycle continues.

The approach takes advantage of the fact that placing an object is equal to the reversed act of grasping an object. A joint

configuration used to place an object can likely also facilitate grasping.

V. EMBODIED DIALOGUE SYSTEM

Moore [25] recently remarked that ”many roboticists regard a speech-enabled interface as a somewhat independent, bolt-on goody rather than a natural extension of a robot’s perceptuo-motorsystem”. This is to be viewed as problematic since it is much more advantageous to treat the language capabilities of a robot as part of the overall system and it is further much more in line with what is known about the role of language in humans. Feldman [26] indeed points out that language and cognition are best understood as a result of the brain being shaped for control of a physical body which navigates within a social world. However, this tighter coupling of the body, its control, and the dialogue processing system have not been fully explored in human-robot interaction research with respect to system design.

In this paper, an embodied dialogue system is implemented as a command center connecting all components that are involved in accomplishing visuomotor tasks. The Dialogue System is embodied as the decision-maker connecting each component together instead of being an independent module itself. The agenda-driven dialogue system guides the humanoid robot in achieving a goal, such as to test its grasping ability or to perform an object learning training. The goal is achieved by the joint-task agenda approach [15] in which tasks are accomplished by collaboration, combining effort from the humanoid robot and the user. The humanoid robot carries out its motor actions and reports its progress throughout the process while the user has hands the object to the robot upon request and provides assistance in case of errors.

The structured dialogue model is an effective model to be implemented in our goal-oriented dialogue system, as the states are atomic and finite, with its structure, position and information of each state fixed and domain-oriented [27]. The transition from one state to another is predefined, much like an if-else function: if object grasping is successful, perform action A; else, perform action B. This approach is useful in limiting the search space, thus increasing efficiency. Besides, the dialogue flow is controlled by restricting the flexibility. As the goal would require the humanoid robot to perform certain tasks in sequence, such as loading the neural network before getting joint values, implementing a finite state approach simplifies the interaction design.

As a decision-making module, the dialogue system embodies six components in performing tasks: Motion, Vision, Emotion, Computation, Knowledge and Natural Language Generation. The Motion component controls the sensorimotor ability of the humanoid robot such as moving the robot’s hand towards the object. The Vision component, the eyes of the humanoid robot, captures stereo images as inputs for the computation. The Emotion component shows facial expressions on the robot’s face using embedded LED lights, such as happy and sad expressions. The Computation component loads the trained model to the neural network and computes respective joint

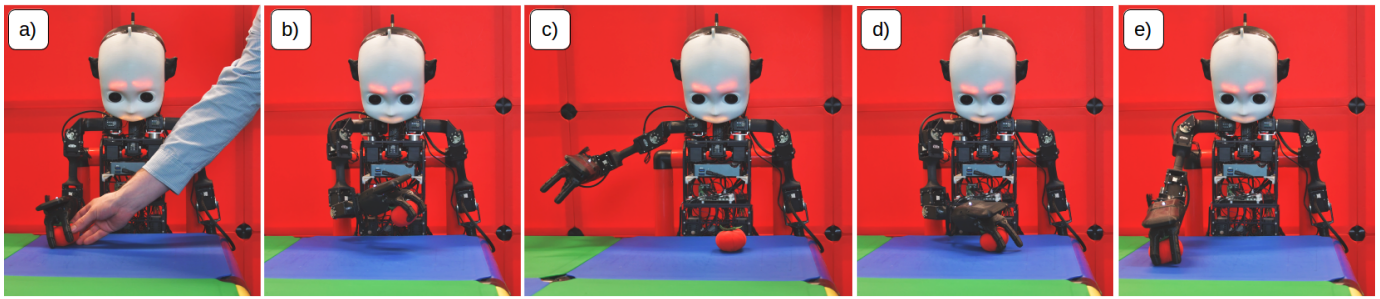


Fig. 2. Training cycle for grasping: a) A human experimenter places the training object in the humanoid’s hand. b) The humanoid moves the object to a random position on the table c) The humanoid places the object on the table and records an image. d) The humanoid’s arm moves back to the last joint configuration to grasp the object again. Steps b) to e) are repeated to gather more samples.

values for grasping. The Knowledge component stores and provides information for the tasks, and the Natural Language Generation component outputs speech response to the user via text-to-speech synthesis. The Dialogue Manager is implemented using SMACH¹, a state machine library developed by ROS. The Motion and Emotion components are implemented using NICOmotion, a library to execute the NICO robot, developed by the Knowledge Technology team [23]. The Vision component uses a common USB protocol. The Computation component is a Convolutional Neural Network developed using Theano and Lasagne². The Knowledge component is build using PyKE³, a Python-based knowledge engine. The Natural Language Generation component is implemented using the Python Google Text-to-Speech library⁴.

Combining the components’ functionalities in a specific order for each task, the embodied dialogue system decides which action to perform next, according to which task it is currently doing and which input it has received. There are ten dialogue states for the system: *Control*, *Perception*, *Grasp*, *Fail*, *Success*, *NLG*, *Train*, *Relax*, and *Test*, followed by a *termination state* at the end (Figure 3). The *Control* state receives a command from the user and decides which task is to be performed by the humanoid robot among four available ones: test object grasping, train object grasping, release motor torques or load information. For example, if the train object grasping task is requested, the *Control* state executes the *Perception* state. In that state, the Motion component is called to move the robot’s hand and lower the head, followed by the Vision component to capture and save stereo images to file. The next state *Grasp* loads the pre-defined model to the neural network, computes joint values based on the stereo images and performs the motion of grasping. Depending on the grasp outcome, if no grasp object is detected, the *Fail* state is executed which passes a dialogue ID to the *NLG* state. On the other hand, if a grasp object is detected, the *Success* state is executed which passes a different dialogue ID to the *NLG* state. The *NLG* state maps the dialogue ID to the

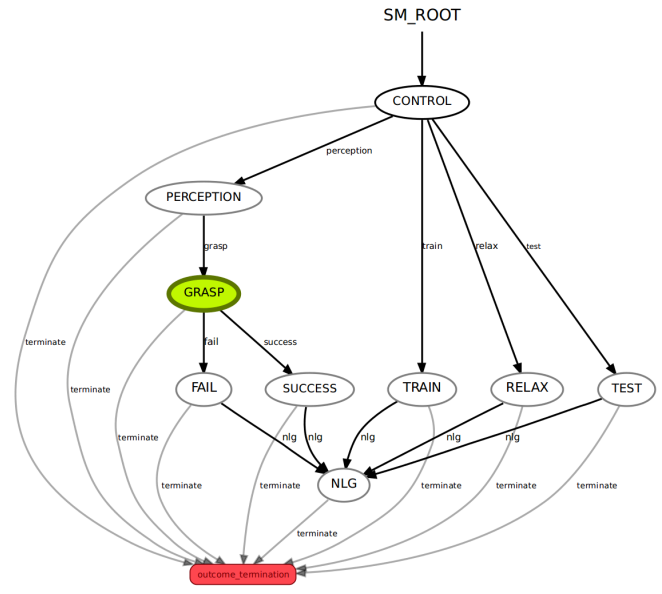


Fig. 3. Dialogue States for Object Grasping

respective sentence which is then conveyed to the user through speech using text-to-speech synthesis. The NLG function can be executed concurrently while another state is being executed, for the purpose of reporting progress without interfering with the executing action.

VI. EXPERIMENTAL PROTOCOL

A grasp-training task will be used to compare the effect of *robot-guided* and *human-guided* learning of a humanoid. After informing the participants about the experiment and gaining written consent, we will use a questionnaire to evaluate their previous experiences with robots. In the next step, the participants will be randomly assigned to one of the two conditions. The same protocol will be used for both the *robot-guided* and the *human-guided* learning conditions, with distinction in the way of interaction: in the *robot-guided* learning scenario, the robot will communicate with the user using Natural Language Generation, gaze, and display of emotions throughout the process whereas, in the *human-guided*

¹<http://wiki.ros.org/smach> [Accessed: 14.06.2017]

²<https://lasagne.readthedocs.io/en/latest/> [Accessed: 14.06.2017]

³<http://pyke.sourceforge.net/index.html> [Accessed: 14.06.2017]

⁴<https://pypi.python.org/pypi/gTTS> [Accessed: 14.06.2017]

Active robot-guided Scenario	Human-guided Scenario
I am ready to look.	The robot is ready to look.
Please put the learning object onto the table for me.	Please put the learning object onto the table for the robot.
I am looking at the object.	The robot is looking at the object.
I have loaded the neural network.	The robot has loaded the neural network.
I computed the joint values.	The robot computed the joint values.
I am ready to grasp.	It is ready to grasp.
I grasp the object.	The robot grasps the object.
(success) Here you go, this is for you.	(success) The robot lifts the object.
(failure) Oh no, I failed to grasp the object.	(failure) Oh no, the robot failed to grasp the object.
I will try again.	It will try again.

TABLE I

Dialogues for Grasp Evaluation of Active Robot-Guided and Human-Guided Learning Scenario

learning scenario, the same dialogue will be given by the experimenter to the user. Table I shows the dialogue for both scenarios. In the *human-guided* condition, the robot will not engage in dialogue interaction with the participant. Other than that, both conditions will have the same steps:

1) *Step 1*: The learning phase begins by handing the training object to the humanoid: The humanoid’s hand moves to the starting position and opens. The participant is asked to

place the training object into the hand. The humanoid then closes its hand and places the object in a random position on the table. Upon placing the object, the robot moves its hand away from the table to capture pictures. The hand is then moved back to re-grasp the object and continue the training cycle. As described in section IV-A, the learning phase is mostly autonomous after the participant has initially handed the training object to the robot. Should an error occur, like the training object falling off the table or being shoved away during grasping, the participant is alerted, and instructions are given to hand the object back to the humanoid before the training is resumed. In the *robot-guided* condition, the humanoid uses haptic perception to detect failed grasps, moves its head to an upright position to face the participant, displays a sad face, and requests help. In the *human-guided* condition, the experimenter stops the training cycle and instructs the participant to hand back the object. The learning phase lasts for 10 minutes; we use a fixed time to evaluate how many samples are collected in this time frame.

2) *Step 2*: Next, the participant is asked to evaluate the learned visuomotor skills by placing the training object repeatedly in front of the humanoid. The participant is truthfully informed that the evaluated sensorimotor skills were trained in the same way as in the learning phase, but more samples were necessary and thus an already trained neural model is used. After the participant places the object on the table, the humanoid looks down and records images. Using its neural network, the robot processes these images to compute joint values for grasping. In the *robot-guided* condition, the humanoid differentiates successful and failed grasps. It reacts accordingly by looking up, displaying a smiling face and offering the object to the participant, or by looking up and making a sad face, see Figure 4. In the *human-guided* condition, the robot lifts its hand regardless of the success of the grasp. Table I shows the dialogue for the *robot-guided* and *human-guided* conditions. This phase lasts for 5 minutes.

3) *Step 3*: Finally, the participant is asked to evaluate both the interaction and the humanoid. To evaluate the interaction, a specialized questionnaire is used. The humanoid is evaluated with the GODSPEED questionnaire [21]. Participants are asked to rate the humanoid on 24 five-point scales between pairs of opposed adjectives, e.g. artificial vs lifelike. The items cover the five categories anthropomorphism, animacy,

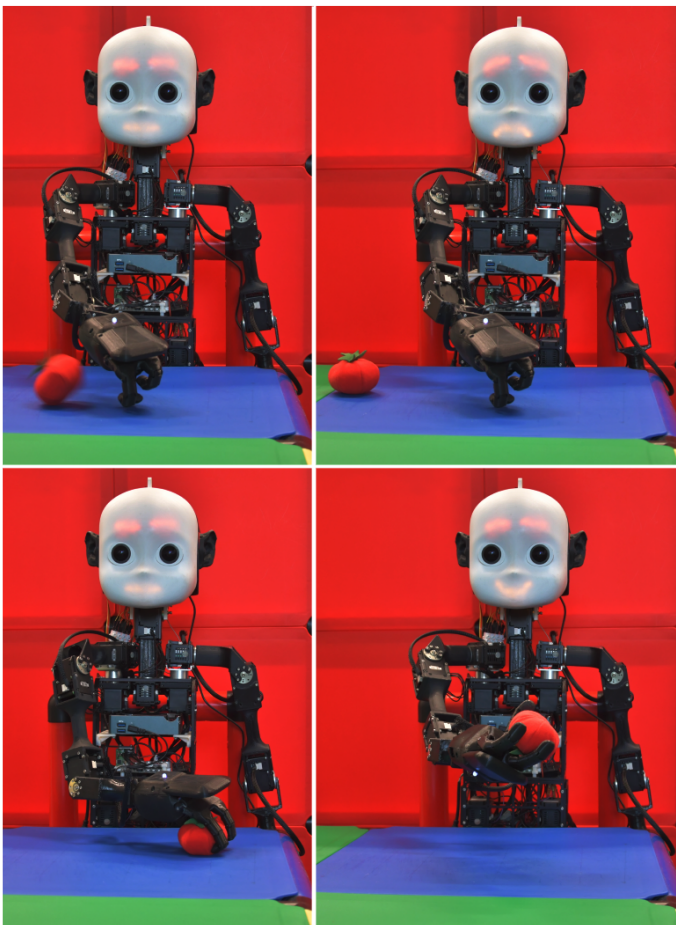


Fig. 4. Top: Failed grasp attempt; the humanoid displays a negative emotion to alert the human interaction partner. Bottom: Successful grasp attempt; the humanoid displays a positive emotion and offers the object to its interaction partner.

likeability, perceived intelligence, and perceived safety. After the experiment, participants are debriefed.

VII. CONCLUSION AND FUTURE WORK

We aim at researching the effect of an active role in learning of a humanoid robot. We want to evaluate how well a learning scenario that is solely mediated by the humanoid works, how well users accept such a scenario and how the robot's active role in learning influences the participant's perception of the robot. To answer these research questions in a principled way we designed and realized an experimental setup. We chose the child-sized humanoid NICO [23] which, due to its appearance, should enable participants to easily adopt the role of a teacher. NICO's arms have a human-like range of motion, which enables it to manipulate an object in front of the body. A haptic sensing mechanism in the hands informs the robot of successful grasp attempts. NICO can display emotion on its face to further enhance the human-robot interaction.

We employ a state-of-the-art approach for visuomotor skill acquisition based on deep neural learning to increase the authenticity of the scenario. The participants will train the robot using the same way it has been trained by researchers. All system components are integrated into an embodied dialogue system that not only handles the verbal interaction with the user but also uses knowledge about the learning progress and haptic sensing to control a multimodal interaction that includes physical actions and display of emotions.

In future work, we will evaluate the active learning humanoid against a baseline scenario where a human experimenter instructs the participants to assist a robot. We will use the insights gained from this study to improve the entire experimental setup, ranging from the robotic hardware to the embodied dialogue system.

ACKNOWLEDGMENT

This work was partially funded by the German Research Foundation (DFG) in project Crossmodal Learning (TRR-169) and the Hamburg Landesforschungsförderungsprojekt CROSS.

REFERENCES

- [1] J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural networks*, vol. 61, pp. 85–117, 2015.
- [2] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [3] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *arXiv preprint arXiv:1504.00702*, 2015.
- [4] L. Pinto and A. Gupta, "Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours," in *Robotics and Automation (ICRA), 2016 IEEE International Conference on*. IEEE, 2016, pp. 3406–3413.
- [5] M. Kerzel and S. Wermter, "Neural end-to-end self-learning of visuomotor skills by environment interaction," in *International Conference on Artificial Neural Networks (ICANN)*, 2017 accepted.
- [6] A. Cangelosi and M. Schlesinger, *Developmental Robotics: From Babies to Robots*. MIT Press, 2015.
- [7] A. M. Turing, "Computing machinery and intelligence," *Mind*, vol. 59, no. 236, pp. 433–460, 1950.
- [8] A.-L. Vollmer, M. Mühlhig, J. J. Steil, K. Pitsch, J. Fritsch, K. J. Rohlfing, and B. Wrede, "Robots show us how to teach them: Feedback from robots shapes tutoring behavior during action learning," *PLoS one*, vol. 9, no. 3, p. e91349, 2014.
- [9] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [10] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [11] F. Cruz, S. Magg, C. Weber, and S. Wermter, "Training agents with interactive reinforcement learning and contextual affordances," pp. 271–284, Dec 2016. [Online]. Available: <https://www2.informatik.uni-hamburg.de/wtm/publications/2016/CMWW16/>
- [12] F. Cruz, G. I. Parisi, J. Twiefel, and S. Wermter, "Multi-modal integration of dynamic audiovisual patterns for an interactive reinforcement learning scenario," *IEEE*, pp. 759–766, Oct 2016. [Online]. Available: <https://www2.informatik.uni-hamburg.de/wtm/publications/2016/CPTW16/>
- [13] D. Gibbon, R. Moore, and R. Winski, Eds., *Handbook of standards and resources for spoken language systems*. Berlin: Walter de Gruyter, 1997.
- [14] D. Griol, Z. Callejas, R. López-Cózar, and G. Riccardi, "A domain-independent statistical methodology for dialog management in spoken dialog systems," *Computer Speech & Language*, vol. 28, no. 3, pp. 743–768, 2014.
- [15] P. Piwek, "Dialogue with computers: dialogue games in action," in *Dialogue across media book*, ser. Dialogue Studies, J. Mildorf and B. Thomas, Eds. Amsterdam: John Benjamins, 2017, vol. 28.
- [16] H. Cuayhuil, "Simpleds: A simple deep reinforcement learning dialogue system," in *Dialogues with Social Robots*, K. Jokinen and G. Wilcock, Eds. Singapore: Springer, 2017, pp. 109–118. [Online]. Available: http://link.springer.com/chapter/10.1007/978-981-10-2585-3_8
- [17] R. T. Lowe, N. Pow, I. V. Serban, L. Charlin, C.-W. Liu, and J. Pineau, "Training end-to-end dialogue systems with the ubuntu dialogue corpus," *Dialogue & Discourse*, vol. 8, no. 1, pp. 31–65, 2017.
- [18] A. Perzylo, S. Griffiths, R. Lafrenz, and A. Knoll, "Generating grammars for natural language understanding from knowledge about actions and objects," in *Proceedings of the IEEE International Conference on Robotics and Biomimetics (ROBIO)*, Zhuhai, China, December 2015, <http://youtu.be/mgPQevfTWP8>.
- [19] M. Eppe, S. Trott, and J. Feldman, "Exploiting deep semantics and compositionality of natural language for human-robot-interaction," in *Intelligent Robots and Systems (IROS), 2016 IEEE/RSJ International Conference on*. IEEE, 2016, pp. 731–738.
- [20] S. Griffiths, F. A. Eyssel, A. Philippesen, C. Pietsch, and S. Wachsmuth, "Perception of artificial agents and utterance friendliness in dialogue," in *Proceedings of the 4th International Symposium on New Frontiers in Human-Robot Interaction at the AISB Convention 2015*. The Society for the Study of Artificial Intelligence and Simulation of Behaviour, 2015.
- [21] C. Bartneck, D. Kulić, E. Croft, and S. Zoghbi, "Measurement instruments for the anthropomorphism, animacy, likeability, perceived intelligence, and perceived safety of robots," *International journal of social robotics*, vol. 1, no. 1, pp. 71–81, 2009.
- [22] N. Churamani, M. Kerzel, E. Strahl, P. Barros, and S. Wermter, "Teaching emotion expressions to a human companion robot using deep neural architectures," in *2017 International Joint Conference on Neural Networks (IJCNN)*, May 2017, pp. 627–634.
- [23] M. Kerzel, E. Strahl, S. Magg, N. Navarro-Guerrero, S. Heinrich, and S. Wermter, "NICO – Neuro-Inspired COmpanion: A developmental humanoid robot platform for multimodal interaction," in *Proceedings of the IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*. IEEE, 2017 accepted.
- [24] G. Metta, L. Natale, F. Nori, G. Sandini, D. Vernon, L. Fadiga, C. von Hofsten, K. Rosander, M. Lopes, J. Santos-Victor, A. Bernardino, and L. Montesano, "The iCub humanoid robot: An open-systems platform for research in cognitive development," *Neural Networks*, vol. 23, no. 8–9, pp. 1125–1134, 2010.
- [25] R. K. Moore, "From talking and listening robots to intelligent communicative machines," in *Robots that talk and listen*, J. Markowitz, Ed. Berlin: de Gruyter, 2015, pp. 317–335.
- [26] J. Feldman, "Embodied language, best-fit analysis, and formal compositionality," *Physics of Life Reviews*, vol. 7, no. 4, pp. 385–410, 2010.
- [27] D. Schlangen, "Modeling dialogue: Challenges and approaches," *Künstliche Intelligenz*, vol. 3, no. 5, pp. 23–28, 2005.