

Modeling Action and Intention for the Production of Coordinating Communication in Human-Robot Task Collaborations

Aaron St. Clair and Maja J Matarić, *Fellow, IEEE*

Abstract—In this paper we consider a method of increasing team fluency in co-located human-robot task collaboration. In such scenarios, the robot can use explicit communication actions to convey internal state and intention, thus providing the collaborating partner with more information than s/he could glean from observation of a robot that does not communicate. We discuss the fundamental requirements of effective coordinating communication as derived from human-human collaboration tendencies. A preliminary approach, based on probabilistic modeling of intentional behavior for production of multimodal communication collaborative tasks, is presented with the aim of reducing idle time and increasing team situational awareness.

I. INTRODUCTION

AS robots become increasingly capable of performing a variety of tasks in real-world, dynamic environments, around people, the need arises for a principled means for robots to coordinate their behavior and signal their intentions to those around them. Effective coordination is a particularly acute requirement in human-robot task collaborations, in which mixed groups of human and robot agents have a common goal and work together as a team to achieve it. In this paper we focus on physically decoupled collaborations between a single robot and a single person interacting in a shared space. This scenario was selected since we are fundamentally concerned with the form and planning of inter-agent communication, which is further complicated by issues of scalability and observability when considering larger teams. Furthermore, these types of one-on-one interactions are frequently encountered and studied in areas such as socially assistive robotics and personal, in-home robots. We exclude the physically coupled case, in which both the robot and the person manipulate the same object at the same time, since, in such contexts, the dynamics of the system are likely a preferable means of achieving coordination than is social interaction.

Coordinating communication, while not strictly necessary in all tasks, serves to reinforce the perspective taking and intention attribution processes that people naturally employ when watching and interacting with others [1] as well as to resolve potential conflicts at the earliest opportunity. People are naturally very good at observing a scene and attributing motivating desires and intents, even in

the absence of any underlying purposeful action [2]. Based on this, we believe clearly conveying desires and intentional action to be a key component of coordination and consequently of the robot's communication capabilities.

II. ON COORDINATING COMMUNICATION

In this section we motivate the selection of robot communication modalities by drawing upon relevant literature from the social sciences and develop a definition of coordination that will inform the design of our communication model.

People exhibit a number of social behaviors during task collaboration, including a tendency to align their frames of reference [3] as well as linguistic representation at both syntactical and lexical levels [4]. Beyond verbal modalities, there is evidence to support that people account for the relative position of an observer when producing gestures, in particular during deictic (pointing) referencing [5, 6]. There is also evidence that speech and gesture production are parts of a single system and are tightly coupled [7], thus often conveying redundant information when used together. This is the case of emblematic gestures that can be directly translated to speech, including the subcategories of iconic gestures, which replicate the form of a particular object, or deictic gestures that are often associated with referential linguistic constructs such as 'that' or 'there'.

Despite a reasonable understanding of human behavior in similar circumstances, it is inconclusive whether or not robots should behave exactly like people do. We do however know that people are capable of employing multiple different communication modalities, often in combination, during social interaction and collaboration. Furthermore, since the robot is assumed to be an active participant in the task, it is possible that some portion of its sensing or actuation capabilities could be in use at any given time. For example, in a manipulation task, it might need to observe a specific portion of a table in detail, or, while carrying an object, it may not be able to use its arms. This argues for the robot to make use of multimodal communication capabilities comprising both speech and embodied gesture.

The actual form of communicative acts generated using these modalities, especially embodied gesture, will be intrinsically tied to the form and function of the robot's physical embodiment. Although we have made the assumption that the robot and the person have a shared, well-defined goal *a priori*, we do not assume that they have seen each other or otherwise colluded before the collaboration begins. We base on the desire for the robot to be able to effectively interact with diverse populations of non-expert users without requiring extensive training, as well as the

Manuscript received July 16, 2012. This work was supported in part by National Science Foundation (NSF) grants CNS- 0709296, IIS-0803565, and IIS-0713697 and ONR MURI grant N00014-09-1-1031.

A. St. Clair (corresponding author) and M. J Matarić are with the Computer Science Department at the University of Southern California, Los Angeles, CA 90089-0781 USA (e-mail: astclair@usc.edu; mataric@usc.edu).

observation that most day-to-day activities that people do together can be completed with *ad hoc* communication between parties. This excludes, for example, highly-trained construction or military teams that employ a pre-defined system of signaling that must be learned and used by all members of the team.

The assumption that the team has not previously interacted has two important consequences: 1) the team will not have a pre-determined communication strategy and 2) the team will not have agreed on a strategy for achieving the goal. In this context, we wish to plan robot task actions according to the following constraints: 1) the robot should select actions that minimize conflicts with the expected next actions of the person, and 2) the robot should select actions that advance the progress of the team toward the goal state. Conflict can occur at the tactical level, such as selecting the actions requiring the same resource at the overlapping points in time, and at the strategic level, where the agents may disagree on the sequence actions should be executed, regardless of role assignment, in order to achieve the goal (Figure 1). We will focus on minimizing the first type of conflict by enabling the robot to issue communication expressing its intended atomic actions to the collaborating partner, thus improving the robot’s readability and increasing situational awareness.

III. PROPOSED APPROACH

The lack of a pre-defined communication policy to provide structure to communication between collaborating teammates highlights the need for a framework that can intelligently compose primitive communication modalities, such as deixis, emblematic gesture, and speech, into timely, meaningful communicative acts in response to changes in

task-level contexts such as the robot’s internal representation. Such a framework should enable the robot to use combinations of communication modalities to express its state in a natural and timely manner to a person, with the goal of improving the collaboration and making the robot more readable. Such a framework will rely on both well-studied, validated communication production capabilities, and on methods for activity recognition and understanding to monitor the state of the world and the human partner. While there is a growing body of work on specific categories of embodied communication in HRI and embodied agent contexts [9], including our prior work on deictic gesture [10], it is non-trivial to incorporate disparate results into a validated communicative repertoire for any given robot platform. There is also a large body of existing work in human activity recognition and understanding [11, 12] to draw upon when designing perception systems and deciding on an underlying representation for task-related activity.

We propose a data-driven, probabilistic approach for constructing and generalizing models of communicative behavior based on goal-directed, intentional behavior. This basis on intentional action (sub-goal achievement) is motivated by evidence that people, even in infancy, tend to segment and process action sequences through attribution of goals and desires [2, 13]. Recognizing human intention has been shown to be feasible in a human-robot interaction scenario from an on-board camera and range sensing [14]; by definition the robot has direct access to its own intent. We assume *a priori* that the robot possesses the set of actions it can execute, A_R , and the set of observable atomic actions that the human can perform during the task, A_H . We also assume that the robot has a representation of the task in terms of expected pre- and post-conditions of each primitive task action and the final goal. From this we construct a Markov decision process, consisting of a set of S states, a set of A actions, a state-transition function $T(s, s')$, and a reward function $R(s, a)$ of the following form:

- S represents the state of the task and the set of possible intended role assignments for the robot and the person at a given point in the task, as defined by the cross product of the sets A_R and A_H .
- A , a set of communication actions the robot can execute.
- T is the transition function mapping states to action and accounts for uncertainty in human action selection given the current state and the communication action selected by the robot.
- R is the reward for selecting a given communication action in a given intentional state and is to be learned from human-human and human-robot training sessions.

This model was selected due to a number of desirable properties. First, it allows us to draw on the large body of existing work in adaptation and reinforcement learning, which is useful for adapting a trained model to account for inter-person personality, hierarchy, and stress differences, which influence preferred communication style [15]. There are also well-defined methods for scaling up to multiple agents as in multi-agent MDPs and in relaxing observability assumptions. Furthermore, the probabilistic nature of the model provides a robust means of handling uncertainty and

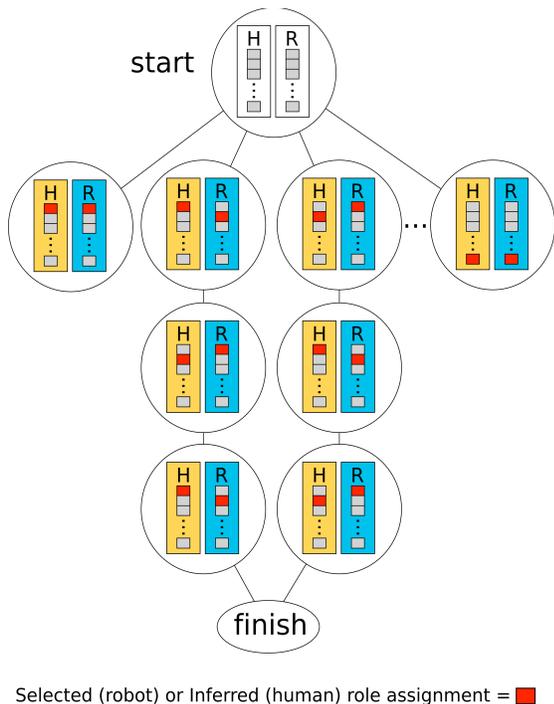


Figure 1. Action role assignment in a simple task requiring the first and second actions to be performed simultaneously. The left- and rightmost nodes are in tactical conflict, while the two longer sequences represent a potential for strategic conflict, if the agents select different strategies.

noise, which are unavoidable when interacting with people. Finally, similar systems based on Markov decision processes have shown promise in modeling high-level human social goal attribution in experiments where participants observed 2D animated stimuli [16].

Communication from the robot will consist of a repertoire of combined verbal and nonverbal feedback actions designed to convey intention-to or intention-that a specific action from the action be executed. That is, conveyed intentions are either in the first-person/self-referential (“I intend to do action a”) or second-person (“You should do action b”). The physical instantiation of these actions in terms of actual words spoken and motion trajectories will be motivated from literature and from annotated data collections of person-person interactions. The output from this system consists of planned communicative actions at timescales that correspond with changes in the task state and user’s estimated intentions. This will be integrated with the overall control system via an executive node that will integrate task and communication performance. As the robot performs the task, some communication actions, such as speech, can easily be executed while simultaneously performing an unrelated task action (e.g., navigation or manipulation), while other communication modalities may conflict with sensing and action used to perform the task. Blending task and communication actions and the micro-timing of their use during the interaction is an area of further research. Our initial implementation will favor communicating when possible by masking in-use modalities and using as many of the planned communication modalities as possible.

We have constructed an augmented reality task simulation system made up of a room with a set of overhead projectors and Microsoft Kinects for person-tracking, to simulate various task environments. A task simulation engine updates the state of the task based on locations of tracked people, our current implementation consists of a sheep herding game. By using a simulated environment, we can automatically tag role assignment and completed actions for each agent, to serve as our ground truth for intentional actions during training. In-room video and audio recording will be used to capture two people performing the task, and communication actions will be annotated using a mix of automatic (text-to-speech) and manual annotation. This data stream will be used to learn the expected utility of issuing a given communication action in a given state. We will subsequently map human communication modalities to robot modalities, enabling the system to be used on a physical robot. Performance will be validated via comparison to a non-communicating robot as well as the two-person team performing the task. We will use both objective metrics (idle time, task efficiency, time-on-task) and subjective measures (user preference, perceived intelligence of the robot) to evaluate performance.

IV. CONCLUSION & FUTURE WORK

We have considered the role of coordinating communication in increasing team fluency with application to human-robot task collaborations, and proposed a probabilistic framework for capturing the link between

multimodal communicative actions and contextually dependent intentional behavior. A data collection of two-person interaction is currently in progress and will be used to train and validate the system on a physical robot. In future work, we hope to address issues relating to generalizing the approach across different tasks with similar physical properties. Also, as different users or user populations will likely have different preferences in communication styles, there is also a need for adaptation methods for adjusting the communication style to the particular user or user population. Additionally, we are interested in emergent differences in communication tendencies as a result of changes in stress level or hierarchical relationships between teammates. The concept of high-level, strategic alignment in human-robot teams also invites further research and will likely depend on meta-level negotiation behaviors.

REFERENCES

- [1] S. Fussell and R. Krauss, “Coordination of knowledge in communication: Effects of speakers’ assumptions about what others know.,” *Journal of personality and Social Psychology*, vol. 62, no. 3, p. 378, 1992.
- [2] F. Heider and M. A. Simmel. An experimental study of apparent behavior. *American Journal of Psychology*, 57:243–249, 1944.
- [3] S. Whittaker, “Things to talk about when talking about things,” *Human–Computer Interaction*, vol. 18, no. 1, pp. 149–170, 2003.
- [4] C. Lockridge and S. Brennan, “Addressees’ needs influence speakers’ early syntactic choices,” *Psychonomic Bulletin & Review*, vol. 9, no. 3, p. 550, 2002.
- [5] A. Bangerter, “Accuracy in detecting referents of pointing gestures unaccompanied by language,” *Gesture*, vol. 6, no. 1, pp. 85–102, 2006.
- [6] M. Louwerse and A. Bangerter, “Focusing attention with deictic gestures and linguistic expressions,” in *Proceedings of the 27th Annual Meeting of the Cognitive Science Society*, 2005.
- [7] D. McNeill, *Hand and mind: What gestures reveal about thought*. University of Chicago Press, 1996.
- [8] P. Ekman and W. Friesen, “The repertoire of nonverbal behavior: Categories, origins, usage, and coding,” *Semiotica*, vol. 1, no. 1, pp. 49–98, 1969.
- [9] T. Fong, I. Nourbakhsh, and K. Dautenhahn, “A survey of socially interactive robots,” *Robotics and autonomous systems*, vol. 42, no. 3, pp. 143–166, 2003.
- [10] A. St Clair, R. Mead, and M. Mataric, “Investigating the effects of visual saliency on deictic gesture production by a humanoid robot,” in *Proceedings, Ro-Man, 2011*, pp. 210–216, IEEE, 2011.
- [11] P. Turaga, R. Chellappa, V. Subrahmanian, and O. Udrea, “Machine recognition of human activities: A survey,” *IEEE Trans. on Circuits and Systems for Video Tech.*, vol. 18, no. 11, pp. 1473–1488, 2008.
- [12] R. Poppe, “A survey on vision-based human action recognition,” *Image and Vision Computing*, vol. 28, no. 6, pp. 976–990, 2010.
- [13] G. Gergely, Z. Nádasy, G. Csibra, and S. Bíró, “Taking the intentional stance at 12 months of age,” *Cognition*, vol. 56, no. 2, pp. 165–193, 1995.
- [14] R. Kelley, C. King, A. Tavakkoli, M. Nicolescu, and G. Bebis, “An Architecture for Understanding Intent Using a Novel Hidden Markov Formulation,” *International Journal of Humanoid Robotics - Special Issue on Cognitive Humanoid Robots*, vol. 5, no. 2, pages 203–224, June 2008.
- [15] J. Shah and C. Breazeal, “An empirical analysis of team coordination behaviors and action planning with application to human–robot teaming,” *Human Factors*, vol. 52, no. 2, pp. 234–245, 2010.
- [16] T. D. Ullman, C. L. Baker, O. Macindoe, O. Evans, N. D. Goodman, and J. B. Tenenbaum, “Help or hinder: Bayesian models of social goal inference,” *Advances in Neural Information Processing Systems*, vol. 22, pages 1874–1882, 2009.