COORDINATING SOCIAL COMMUNICATION IN HUMAN-ROBOT TASK COLLABORATIONS

by

Aaron B. St. Clair

A Dissertation Presented to the FACULTY OF THE GRADUATE SCHOOL UNIVERSITY OF SOUTHERN CALIFORNIA In Partial Fulfillment of the Requirements for the Degree DOCTOR OF PHILOSOPHY (COMPUTER SCIENCE)

August 2015

Copyright 2015

Aaron B. St. Clair

Acknowledgements

I would like to first thank my advisor, Prof. Maja Matarić, for her support and guidance throughout my time at the University of Southern California and for fostering and encouraging my growth as a researcher.

I would also like to acknowledge the members of my committee: Prof. Gaurav Sukhatme, Prof. Nora Ayanian, and Prof. Aaron Hagedorn) as well as my qualiying examination committee, Prof. Milind Tambe. Their insightful advice and recommendations have elevated this work and guided my research directions.

I would also like to thank all of my fellow members of the Interaction Lab. Their support, assistance, and advice have been a constant source of inspiration.

Finally, I would like to express my gratitude to my family for their constant support, encouragement, and understanding throughout my life. I would not have made it this far without them.

Table of Contents

Acknow	ledgem	ents	ii
List of]	Figures		vi
List of '	Tables		ix
List of .	Algorit	hms	х
Abstrac	t		xi
Chapter	r 1: Ir	troduction	1
1.1	Challe	enges in Human-Robot Collaboration	1
1.2	Motiv	ation and Problem Statement	3
	1.2.1	Robot Control and Communication in Structured Tasks	4
	1.2.2	Human Behavior in Structured Tasks	5
	1.2.3	Providing Communicative Feedback for Improved Interactions	7
1.0	1.2.4	Evaluation of Human-Robot Collaboration	9
1.3	Metho	odology for Producing Effective Coordinating Communication	10
1.4	Disser	tation Contributions	11
1.5	Disser	tation Outline	13
Chapter	r 2: B	ackground and Related Work	14
2.1	Huma	n-Human Collaboration	14
	2.1.1	Theory-of-Mind and Perspective-Taking	15
2.2	Multi-	Robot Coordination	16
	2.2.1	Applicability to the Human-Robot Case	18
	2.2.2	Task Classifications	19
2.3	Huma	n-Robot Collaboration	22
	2.3.1	Collaborative Human-Robot Interaction	23
	2.3.2	Task Adaptive Collaborative Robots	24
2.4	Summ	ary	26

Chapter	r 3: Production of Embodied Coordinating Social Communication	27
3.1	Speech	28
3.2	Deictic Gesture	32
3.3	Other Gesture Modalities	35
3.4	Evaluating Human Interpretation of Robot Deictic Gestures	36
	3.4.1 Study Design	37
	3.4.2 Participants	38
	3.4.3 Hypotheses and Outcome Measures	39
	3.4.4 Results	41
	3.4.5 Discussion \ldots	46
3.5	Summary	48
Chapter	r 4: Approach to Coordinating Social Communication	49
4.1	Human-Robot Task Representation	52
4.2	Representing and Recognizing Human Activity	57
	4.2.1 User Policy Recognition	58
4.3	Planning Robot Feedback	61
	4.3.1 Self-Narrative Feedback	61
	4.3.2 Role-allocative Feedback	64
	4.3.3 Empathetic Feedback	67
4.4	Communication Executive	68
4.5	Extensions of the Approach	69
4.6	Summary	71
Chapte	r 5: Evaluation of Coordinating Social Communication	72
5.1	Augmented Reality Task Environment	74
5.2	Pseudo-herding Task	76
-	5.2.1 Application of the Approach in the Pseudo-Herding Task	78
5.3	Study: Communication and Role-Usage in the Pseudo-Herding Task .	79
	5.3.1 Study Design	82
	5.3.2 Participants	83
	5.3.3 Hypotheses and Outcome Measures	83
	5.3.4 Results	83
	5.3.5 Discussion	85
5.4	Study 2: Coordinating Social Communication in Pseudo-Herding Task .	86
	5.4.1 Study Design	86
	5.4.2 Task Description and Setup	87
	5.4.3 Participants	88
	5.4.4 Hypotheses and Outcome Measures	88
	5.4.5 Results	89
	5.4.6 Discussion	91
5.5	Study 3: Coordinating Social Communication with Older Adults	93

5.6	Simula	ated Cooking Task	93
	5.6.1	Application of the Approach in the Cooking Task	95
	5.6.2	Task Description and Setup	97
	5.6.3	Study Design	98
	5.6.4	Participants	98
	5.6.5	Hypotheses and Outcome Measures	99
	5.6.6	Discussion	99
5.7	Summ	ary	100
Chapte	r 6: Si	immary and Conclusions	102
6.1	Major	Contributions, Findings, and Insights	102
6.2	Open	Problems and Future Work	104
	6.2.1	Extension to multi-robot and multi-person teams	105
	6.2.2	User role specification and generalization across tasks	105
	6.2.3	Communication personalization and adaptation	105
Bibliog	raphy		107

List of Figures

2.1	A selection of most relevant prior work in human-robot task collaboration: (a) AUR robotic desk lamp (Hoffman and Breazeal, 2010), (b) Chaski (Shah et al., 2011), and (c) Wakamaru humanlike robot (Mutlu et al., 2013).	23
3.1	The PR2 mobile manipulator attempting to produce a "step back" ges- ture. Producing gestures on a robot platform that was not designed with this purpose in mind is a difficult problem and currently requires skilled hand animation.	34
3.2	Bandit pointing with its (a) head; (b) straight-arm; (c), bent-arm; and (d) and head+arm.	36
3.3	The experimental setup showing a participant indicating the perceived location of the Bandit robot's referent target with a laser pointer among a set of salient targets.	38
3.4	(a) Mean angular pointing error and (b) Mean angular error between perceived and actual targets.	42
3.5	(a) Mean angular error between perceived and desired targets and (b) Mean angular error between perceived and desired targets for cross-body and straight-arm points.	42
3.6	(a) Mean time from start of pose to marking of estimated point and (b) Mean error by saliency condition, human pointer.	43
3.7	(a) Mean angular error with respect to horizontal target position as seen from the participant's perspective and (b) Mean angular error with re- spect to vertical target position as seen from the participant's perspective.	44

4.1	A brief over of the major components of the approach including the user policy recognition system, the communication planner with three types of communicative intents and the communication realizer.	52
4.2	A system diagram showing the interaction of the components of the approach with human activity tracking in green, robot task planning in pink, robot communication planning in blue, and external input and output modules in orange and purple, respectively.	56
4.3	An example of the effect of implicit feedback on model-based action se- lection. For simplicity, the user and robots actions are combined in the transitions. If the robot clearly conveys that it will take action a_{r2} then the user can safely select action a_{h2} without fear of reaching S_3	63
5.1	Diagram of the augmented reality task simulation environment with a person and Pioneer 2 robot	74
5.2	A person and a Pioneer 2 AT robot collaborating on the pseudo-herding augmented reality task with the virtual elements shown projected on the floor of the room. In a) the lock is unlocked and only one sheep is in the pen (light blue) while in b) the robot has just finished the game by collecting the last sheep	76
5.3	Example of the Wizard of Oz interface used by the experimenter to con- trol the Pioneer 2 robot during the pilot study. The control interface in a) contains buttons to trigger each of the autonomous actions from the robot's action set and a cancel button to stop all motion, as well as text displays for completion status and error output. In b) the visualization of the robot, person, and simulated objects potions is depicted	80
5.4	An overhead diagram of the experimental setup used when constructing the action set of the robot. Users were presented various diagrams and asked to instruct one of the participants what to do next	81
5.5	Bar chart showing sheep herding allocation across agents for person- person and person-robot teams. Note the less equitable allocation of sheep herding in the human-robot teams compared to the human-human teams	85
5.6	Time to completion and allocation counts of agent actions in communi- cating and silent conditions	90

- 5.7 Screen capture of the interface used for the virtual cooking task. Users are asked to fill all the orders in the top row by moving food items from the green ingredients box through two stages (preparation and cooking) until all items are fulfilled as indicated by green check marks.

94

List of Tables

3.1	Samples collected in the deixis experiment for each condition; arm is over- represented in the non-salient condition to compare away-from-body and	
	cross-body	39
3.2	Results of deixis experiment ANOVA by condition	41
5.1	Post-experiment survey results	90
5.2	Survey results on the pseudo-herding evaluation comparing communicat- ing and non-communicating robots. Scale is 0 for strongly disagree, 6 for	
	strongly agree.	90

List of Algorithms

4.1	User role recognition algorithm	60
4.2	Alternative role-allocation suggestions	66
4.3	Role-allocative feedback - single-agent MDP - non-conflicting action $\ .$.	69

Abstract

Robots have become increasingly capable of performing a variety of tasks in real-world dynamic environments, including those involving people. Beyond competently performing the tasks required of them, service robots should also be able to coordinate their actions with those of the people around them in order to minimize conflicts, provide feedback, and build rapport with human teammates in both work environments (e.g., manufacturing) and home settings. Humans coordinate their actions in various task settings through structured social interaction aimed at representational alignment and intentional feedback. In order for robots to coordinate their actions using similar modalities, they must be capable of contextualizing the actions of human partners and producing relevant natural communicative behaviors as the task progresses. This dissertation is motivated by the high-level goal of *producing effective social feedback during task performance*, and alleviating the burden of coordinating the team's joint activity by allowing human users to interact with robots through natural social modalities as partners rather than as operators.

This dissertation develops an approach for constructing and generalizing models of role-based coordinating communication during physically-decoupled human-robot task scenarios, specifically pairwise collaborations in which a person and a robot work together to achieve a shared goal. The approach is validated in different task contexts with different user populations using objective and subjective measures of task performance and user preferences. To support role-allocative communication observed in our pilot experiments with two-person teams, the human-robot collaboration problem is formulated as a Markov decision process in which roles are represented by a set of policies capturing different action selection preferences and accounting for unequal capabilities between human and robot collaborators. A probabilistic method is used to track the user's activity over time and to recognize the role assumed by the user, communication is then planned given the expected policy of the user, the policy of the robot, and the current task state. The communication generated by the robot consists of three types of speech actions and associated co-verbal behavior: 1) self-narration of the robot's activities, 2) role allocation suggestions for the user, and 3) empathetic displays when positive and negative state changes occur.

The approach was validated initially on a dynamic augmented reality herding task with a population of convenience users using objective metrics (idle time, distance traveled) as well as subjective evaluations (user preference, perceived intelligence of the robot), where a higher utilization of the robot and more equitable path distance was observed in comparison to a non-communicating robot. The generalizability of the approach to a different task setting and user population was also evaluated on a cooking task with an elderly user population. The contributions of this dissertation lie in the development of an approach for modeling human-robot task performance for the planning and production of effective robot verbal feedback.

Chapter 1

Introduction

This chapter provides an overview of human-robot collaboration and motivates the use of social communication by a robot as a means of coordinating human-robot teamwork. The task representation and its applicability to different human-robot teamwork scenarios is introduced, as are the types of embodied robot communication modalities used by the robot. The communication planning approach is presented. The chapter concludes with an outline of the rest of the dissertation and a list of primary and secondary contributions of this work.

1.1 Challenges in Human-Robot Collaboration

Robots are becoming increasingly competent at performing tasks such as navigation and manipulation in real-world human environments including homes, the workplace, and other settings where they are around people. This expanded range of capabilities makes possible many scenarios in which robots not only take action around people but also collaborate with them to accomplish tasks together. In a home environment a robot might help a person to clean up a messy room, in the workplace a robot might fetch and deliver objects and tools, in a manufacturing setting a robot might help assemble a product, and in an assistive scenario the robot might provide instructions guiding a user through an unfamiliar activity. These collaborative scenarios present a number of challenges; not only must the robot avoid physical collision with the person, it must also be a productive part of the collaboration accounting for the person's activity when deciding what to do next, and provide a means for users to understand and control its behavior.

This dissertation addresses a number of problems related to generating appropriate coordinating social communication from a robot to a person during the course of a pairwise task collaboration with the aim of allowing the robot to provide helpful and guiding feedback to a human teammate. Effective feedback is important in these scenarios since it has a direct impact on the situational awareness of the robot's collaborator. A general methodology was developed that allows the robot to produce coherent, meaningful social communication during task performance, in order to efficiently and intelligently provide feedback to users during a task collaboration. The problem of producing coordinating social communication actions with a robot in order to provide feedback and support the in situ decision-making of a human teammate requires careful consideration of the kind of feedback the robot should provide as well as the task and user model employed by the robot. Social science literature and existing work in human-robot interaction indicates many communication actions people use that could be beneficial for a robot to produce during a collaborative task, including attentional cues to indicate an area of focus Mutlu et al. (2012), staging actions to maximize shared visual information (Fussell et al., 2000), gestural and speech cues indicating intentional goals or instructions (Grosz and Sidner, 1986, Huang and Mutlu, 2013a), and coaching actions such as instructional feedback Fasola and Matarić (2013), Huang and Mutlu (2013b), and empathetic displays to build team rapport (Lamm et al., 2007).

1.2 Motivation and Problem Statement

The aim of the work is to enable robots to provide effective task-relevant feedback to colocated teammates using social communication modalities such as speech and gesture. Rather than requiring users to gain competence in operating an autonomous robot via a screen-based interface, the goal is to enable the robot to use human-like coordination mechanisms, specifically embodied social communication such as speech and embodied gesture. The key supposition, on which this work is based, is that a social robot teammate will make a more effective work partner and improve quantitative measures of team performance by supporting the natural coordination mechanisms people use when working together in a co-located setting. In home environments, for instance, the robot should be able to clearly communicate its intentions without requiring a nonexpert user to consult a tablet to make sure it is doing what they asked. Similarly, in manufacturing settings, where safety is the chief concern, allowing workers to monitor a robot's task performance directly via visual observation and verbal feedback instead of visually monitoring the robot and a graphical user interface may result in reducing the operator's cognitive load. The work in this dissertation is aimed at supporting the user's natural collaboration modalities and preferences in a broad range of co-located, pairwise human-robot task collaboration scenarios.

Allowing users to interact with the robot as a teammate rather than an operator has the potential to offload the burden of coordination from being solely the person's responsibility, resulting in quantitative improvements in team performance. Beyond improved coordination, a robot collaborator that is treated as a teammate has the potential to impact the joint decision-making process in ways that it would not if it were merely controlled by the person, potentially allowing a robot to guide the team to a desired outcome through a specific trajectory. Robot-guided, joint activity is applicable in many scenarios including training inexperienced users, addressing assistive opportunities, and adding variability to repetitive tasks.

The approach employed in this work generated three types of situated verbalizations aimed at providing useful information to a teammate in co-located joint activities. To make the approach generalize across many task settings, it is based on existing task formulations from multi-robot systems. Unlike existing work in human-robot collaboration, it does not rely on a specific conversational structure such as turn-taking, allowing for it to be used in dynamic tasks where users may not adhere to strict conversational structures. By providing these types of feedback, in the form of coordinating social communication, the approach is able to improve quantitative metrics of team performance and subjective evaluation of the task itself by users, as well as enabling the robot to provide guidance for the user. Next, an overview of how people collaborate in these types of settings is presented along with a discussion of how these behaviors affect the human-robot case.

1.2.1 Robot Control and Communication in Structured Tasks

As this work assumes that the robot is autonomous and an active participant in the task, the robot must have a controller that allows it to perform the task competently in a real-world scenario. Existing work has explored the interplay of robot task control with human collaborators and has demonstrated that robots that account for the actions of their collaborators yield better team performance (Gombolay et al., 2013, Nikolaidis and Shah, 2012, Shah et al., 2011). Providing coordinating communication is supplemental to these approaches as the communication the robot provides must always incorporate the robot's current and planned actions, even if these actions are not optimal or are changing over time. Also, existing task models used in multi-robot coordination the requirements do not extend to the human-robot case since communicating and sharing

task-relevant observations and actions with other robots is much different compared to communication with a person. In order for any communication approach to generalize across various task environments, it must support a task model that also allows the robot to act autonomously in different environments. The work in this dissertation employs task models from robot control and multi-robot systems in planning coordinating communication.

1.2.2 Human Behavior in Structured Tasks

People have a natural ability to collaborate comprised of a complex series of behaviors and attribute agency to the people and things around them (Heider and Simmel, 1944). People also recognize teammates' actions via visual observation the mirror neuron system (Buccino et al., 2004) and theorize about the mental models of our teammates including beliefs, desires, and intentions (Baron-Cohen, 1999). Finally, people issue a variety of implicit and explicit social cues aimed at helping others understand each other's actions (Hanna and Tanenhaus, 2004, Keysar et al., 2000, Whittaker, 2003). Many of these processes, with the exception perhaps of strategic planning, occur automatically without conscious thought on our part and are active when the person is collaborating with a robot (Rizzolatti and Sinigaglia, 2008).

To gain more insight into how human-compatible patterns of robot-collaborative activity might be defined, consider the types of activities the robot is likely to be doing alongside a person. For home scenarios, surveys of older adults (Beer et al., 2012) indicate that object manipulation including fetching, moving, and organizing household objects is a priority. In manufacturing settings, platforms such as Rethink Robotics Baxter (Rethink Robotics, 2012) have pick-and-place capability allowing them to manipulate items and pack boxes or do assembly tasks.

In human-human collaborations, the primary use of social mediation is to determine who is responsible for doing what parts of the task. To accomplish this people often make use of the notion of *roles* to denote expected patterns of behavior (Smith et al., 2001). In work organizations, emergency preparedness, team sports, and many other contexts featuring organized group activity people are familiar with assigning and assuming roles that constrain each person's area of responsibility. Unfortunately, there is not a formal definition of roles or method for decomposing a given task definition into constituent behaviors that can be reliably identified as a role across individuals (Smith et al., 2001). In a human-robot task collaboration the robot will never have access to the mental model of the person and the corresponding understanding of the task structure. This is distinct from the multi-robot coordination case, which is well-studied (Gerkey and Matarić, 2004, Roth et al., 2006), where a common representation of the world makes sharing observations and actions easier. In these scenarios one robot can use network communication to transmit a series of symbols that uniquely denotes to another robot that at some point in time it received a specific observation and took a specific action. To successfully coordinate a robot's action with a person requires a representation for common patterns of human behavior (roles) that will allow the robot to 1) identify which role a person has assumed at a given point in time and 2) to produce speech that reliably guides a person to a desired role that coordinates well with the robot's course of action in the future.

Considering these scenarios, roles in most tasks can be made up of spatial constraints such as the positions of work objects and other agents and temporal constraints, i.e., that a person performing a certain role makes use of certain objects at a given place and time. Another constraint for defining roles is that they be easily integrated with the robot task planning process. This is necessary because the robot must have some means of evaluating a given set of role allocations and proposing better alternatives if there are any available. Markov decision processes and related models have been extensively studied and applied to many types of task planning (Boutilier, 1996). Roles, as defined in this work, have many commonalities with policies in decision-making systems in that they indicate a preference for taking certain actions in certain states. The set of roles available for agents to assume will thus be represented in this work as a set of policies. These are not assumed to be either optimal or even capable of completing the task when executed individually. Rather they are used to capture the action-selection preferences of a person undertaking a given role.

1.2.3 Providing Communicative Feedback for Improved Interactions

The primary focus of this work is to improve coordination and assist the user by enabling the robot to issue role-based communicative feedback to guide user behavior and improve the team's situational awareness. A review of the relevant social science literature covering human-human collaborations, as well as observations of person-person task collaboration in pilot experiments, yields a large number of social cues that could be employed by a robot collaborator. These include appropriate use of attentional cues, staging actions to maximize shared visual information, providing well-formed gestural and speech cues, and coaching actions such as feedback, encouragement, and empathetic displays to build team rapport. Verbal feedback is also primary modality for people to convey task allocation information and could also be used by a robot as it is easy to produce across many different robot embodiments although speech recognition and natural language processing remain difficult. Speech also does not require the person to be looking directly at the robot, which is useful in scenarios where the user's visual attention is concentrated on another activity, as is often the case in a distributed task collaboration.

Enabling a robot to effectively produce all these forms of communicative feedback at a human level of competency, while the robot simultaneously performs parts of a task, is not currently feasible due to limitations in robot embodiments and sensing technologies. In some cases, for example, gesture production might require the robot to use part of its embodiment that is otherwise employed for object manipulation. It is also difficult to generalize communication production, in particular embodied gesture, across different robot embodiments; for example, indicating attentional focus is different with humanoid and non-humanoid robots. As another example, iconic gesture consists of making a visual representation of the intended concept, such as making a hammering motion to indicate a hammer (Mehrabian, 1977). Making reliable iconic gestures is highly domain- and platform-dependent and remapping motion trajectories from one platform to another is a difficult open problem. More straightforward gestures such as deixis (pointing) can be generated reasonably reliably on any robot where either arms or a pan-tilt unit can be used as an end-effector but error in referent accuracy may vary depending on the robot's appearance. For this reason, speech and a simplified set of associated coverbal deictic gestures were selected as the primary modalities by which the robot issues feedback in this work.

Speech works well across different embodiments and is effective in communicating intent. On the other hand, speech has obvious limitations in noisy environments and with users with hearing or understanding limitations. Nonetheless, speech is a natural human communication modality that addresses a range of use cases in home and work environments and has been used previously as a method for controlling robots on a variety of tasks including navigating a wheel chair (Levine et al., 1999) to simple object manipulation tasks (Breazeal et al., 2005). Other work. (Scopelliti et al., 2005) has also shown speech as a preferred method of human-robot communication when compared to keyboard input, mouse-based interfaces, and touch screens. Deictic gesture, can be produced coverbally with linguistic constructs such as "this", "that", and "there" or on its own and is used to indicate a referent target in the environment. Existing work has shown how, with special attention paid to timing, speech and deixis can be combined to produce human-like results that aid interpretation (Mutlu et al., 2012).

1.2.4 Evaluation of Human-Robot Collaboration

Evaluating human-robot collaborations typically consists of measuring individual task performance and assessment of contribution to the team performance as a whole. Typically these team scenarios are evaluated objectively based on task performance including measures such as idle time, time to task completion, number of accomplished sub-tasks per unit time, or relative workload. Assessing the efficacy of robot-issued coordinating communication presents unique challenges as the effect of the robot's communicative feedback on the team's task performance is evaluated. As a baseline, the robot task controller is executed without incorporating any communicative feedback, to simulate the scenario where the robot does task control only. Another evaluation approach employed in this work is to selectively enable different types of communicative feedback and attempt to assess the differences in the quality of the collaboration. Also, since one of the primary goals of this work is to make the robot as at teammate were measured by administering surveys and assessing the sentiment of users in textual descriptions of the experiment.

Approaches aimed at generalizing across different task scenarios must also be evaluated in more than one task context and ideally over an extended period of performing the task with the same user to overcome the novelty effect of working with the robot. For instance, more communication may be tolerated in a first interaction with the robot or when the user is unfamiliar with the task compared to subsequent performances when they have become acclimated to the robot's behavior and have mastered the task. Finally, a number of outside factors may have an effect on a particular user's communication usage including personality, whether they are under stress, and the perceived level of authority with respect to the robot.

1.3 Methodology for Producing Effective Coordinating Communication

The approach in this dissertation produces the combination of three types of coordinating social communication: 1) self-narration of the robot's activities, 2) suggestions for what the user should do next, and 3) empathetic displays when positive and negative events occur. The combination provides a balance of information aimed at improving the user's situational awareness and guiding user behavior through both implicit and explicit communicative feedback. The self-narrative feedback provides user's with information about what the robot is going to do next. The role-allocative feedback, offering suggestions to the user, informs the user that the robot is monitoring the user's progress and evaluating things from the user's perspective and also allows the robot to potentially reinforce or change the person's decision-making process. Empathetic feedback, is useful for indicating that the robot has evaluated the current state of the activity in a hopefully similar way to the person and that it is invested in achieving a good result. Producing each of these types of feedback requires specific robot capabilities primarily concerned with monitoring the task state and monitoring the user's behavior over time. These three types of communicative feedback will be generated using a model of rolebased collaborative behavior and then conveyed to the user by the robot in the form of speech and a limited set of coverbal embodied deictic gesture.

1.4 Dissertation Contributions

This dissertation addresses a number of problems related to generating appropriate coordinating social communication from a robot to a person during the course of a pairwise task collaboration with the aim of allowing the robot to provide helpful and guiding feedback to a human teammate. An approach to planning social communication based on a formalism of the human notion of role allocation is presented as well as an analysis of the limitations and applicability of the system to a variety of task scenarios. The implementation details for the system as applied to a sample task environment are presented along with the results of a user study evaluating the efficacy of the approach.

The following are the main contributions of this dissertation:

- 1. A novel representation for modeling user roles in coordinated activities capable of capturing the action selection preferences that people typically exhibit when allocating responsibilities during joint activity and compatible with traditional planning methods.
- 2. An approach for planning coordinating social communication actions in a situated pairwise task collaboration between a robot and a person aimed at providing useful feedback to the user and enabling the robot to actively participate in the joint planning. The planning approach incorporates activity recognition, planning under uncertainty, and principles from social science of human joint activity to produce narrative, role-allocative, and empathetic feedback in real-time during dynamic task activity.
- 3. A framework for executing social coordinating communication during robot task performance, making use of the given social communicative capabilities of the robot in use and potentially including speech and embodied gesture. The robot's communication actions are derived from human-human communication in similar

task situations and applies relevant design principles from the social sciences on joint attention and politeness theory to ensure that the robot provides coherent feedback.

The following are secondary contributions:

- 1. A robot communication executive that produces coordinated speech and gesture during a human-robot task collaboration and also supports speech interrupts with stop words for state changes.
- 2. User studies demonstrating the ability of the approach to improve both team performance and subjective measures of the robot as a teammate, as compared to a silent robot with an identical task control system.
- 3. A user study of the perceptual accuracy of robot head and arm deictic gesture production and control system for reducing grounding errors accounting for the noise introduced due to interpersonal differences in perception.
- 4. An augmented reality task simulator created using overhead projectors and depth sensors (Microsoft Kinect (Microsoft, 2010)) that allows for rapid-prototyping of different dynamic task scenarios supporting any combination of people, robots, virtual agents, and virtual objects.

1.5 Dissertation Outline

The remainder of this document is organized as follows:

- Chapter 2 provides background on existing work in human-robot task collaboration and a review of relevant work in multi-robot task allocation and humanhuman collaboration in relation to the human-robot case.
- Chapter 3 provides an overview of embodied social communication production for human-robot collaborations, describes specific communication modalities that can be used *in situ*, and details a user study evaluating the accuracy of robot deictic gesture production is presented.
- Chapter 4 describes the approach employed for planning coordinating social communication during human-robot task collaborations including the task representation, user activity modeling, role formalism, and communication planning approach for different task models.
- Chapter 5 discusses the application and evaluation of the system in an augmented reality task environment where a person and robot collaborate on a dynamic, pseudo-herding task as well as application to a real-world, physical task involving a human-robot team performing an object manipulation task.
- Chapter 6 provides a summary and concluding statements as well as potential open problems and extensions to the dissertation.

Chapter 2

Background and Related Work

In this chapter prior work in human-robot collaboration, multi-robot coordination, and relevant work from the social sciences on human-human collaboration are surveyed to provide a foundation for the dissertation approach and its contributions to robot production of coordinating social behavior.

2.1 Human-Human Collaboration

It has been demonstrated that people have a tendency to adapt both their speech and actions in response to the person they are interacting with, to be salient and sensible to a collaborating partner, especially in circumstances involving work objects in the environment and frames of reference (Tomasello et al., 2009, Whittaker, 2003). Collaborators align their linguistic representations of the environment, allowing for more effective communicative behavior with their partner(s). This alignment is achieved via a process in which local environmental representations, i.e., specific speech and gesture, are implicitly adopted and propagated to global representations via a priming mechanism (Pickering and Garrod, 2004). The same priming mechanism is similarly used to achieve lexical and syntactic alignment resulting in a consistent vocabulary and shared environmental representation at the task level for corresponding communication. Researchers disagree over how deeply people model their interaction partners and the role these models play in language and gesture production. Some argue that a speaker's model of their addressees plays a central role in production (Lockridge and Brennan, 2002), others suggest it functions as a late stage corrective mechanism for tailoring speech (Keysar et al., 2000), and still others suggest a hybrid approach is used depending on context (Hanna and Tanenhaus, 2004). Additionally, people rely heavily on speech for coordination. Nevertheless, it has been shown that shared visual information can result in more efficient, less verbose utterances (Fussell et al., 2000). This suggests that collaboration can be accomplished without relying exclusively on natural language processing, particularly in co-located scenarios. Robots that collaborate with people should take advantage of the capabilities of human teammates by providing appropriate signaling to improve the collaborating partner's situational awareness. By enabling a robot to produce humanlike communication during a task, this work aims to support effective coordination between people and robots without requiring users to be trained how to use the robot beforehand.

2.1.1 Theory-of-Mind and Perspective-Taking

In order to effectively coordinate a robot's actions with those of its human collaborator, the robot must be able to accurately estimate the human's planned actions from context or from explicit communication. Analogously, the robot must be able to effectively convey its planned actions clearly to a human. This ability to attribute mental state to others and use it to plan and predict behavior is called Theory of Mind (ToM) and has been extensively used for various capabilities with autonomous robots (Scassellati, 2002, Trafton et al., 2005). This work consists of a Theory of Mind-inspired model in which the robot contains estimates of its own state, the state of third parties', and those third parties' estimates of the robot's state. These states contain information relevant to the task including a world model and a partial task allocation i.e., assignments of various agents to sub-tasks. Previous work has demonstrated the viability of similar frameworks to model and learn from human activity at various levels of perceptual abstraction. At the task-level it is generally assumed that the environmental dynamics are accurately detectable by the robot and the modeling relies on some notion of symbolic state (Ullman et al., 2010) while other work focuses on bottom-up learning from raw or annotated sensor input (Hoffman and Breazeal, 2010, Kelley et al., 2008, Scassellati, 2003). This dissertation is focused on enabling natural collaborative communication in realistic collaborative task settings devoid of a strict dialog structure, such as turn taking. Thus, the task-level symbolic states are augmented to allow the robot to formulate collaborative communication and infer intention from human social signals from perceptual features, such as head direction and deictic gesturing, extracted from on-board sensing.

In order for the robot to make use of this state and task representation, it must be able to evaluate the environment from multiple points of view. Existing approaches to *perspective-taking* (Kelley et al., 2008, Trafton et al., 2005) generally rely on transforming sensor data to a local reference frame and planning under the estimated visibility constraints for another agent (Breazeal et al., 2009, Trafton et al., 2008).

2.2 Multi-Robot Coordination

Methods for planning coordinated behavior have been studied in both the planning and multi-agent systems research communities. The majority of the developed methods are not, however, readily and directly applicable in the human-robot interaction context for various reasons. The work in classical planning, for instance, typically relies on a symbolic representation in which the exact action that the person is taking is known at the time of planning. These techniques have been applied successfully in human-computer (HCI) and human-robot interaction (HRI) scenarios, in which a person provides complete verbal narration of their activities with a specialized vocabulary during execution, thereby eliminating perceptual and uncertainty issues. The work in multi-agent systems has proven a number of interesting theoretical bounds for multiagent planning with limited communication/observability (e.g., decentralized partially observable MDPs (Bernstein et al., 2002, Nair et al., 2003)) but often the notion of communication in those contexts consists of sharing a set of observations, and the decision of whether to communicate or not relies on the assumption that all agents' observations, state, and action selection is the same. These assumptions break down when working with people. Human and robot action are both noisy and unreliable, and benefit from methods that employ probabilistic representations and communication. Furthermore, such communication may take advantage of common representations that the robot might not share with the person.

In multi-robot scenarios a number of problem formulations exist that account for different communication types and environmental observability. These formulations typically involve selecting a joint action from the cross-product of possible actions $A_A \times A_B$ for robots, A and B. Communication typically involves communicating local state or observation sequences to other robots in cases which would change the policy of the receiving agent. In multi-robot collaboration scenarios, a number of techniques exist that account for different levels of communication and observability. The simplest assumes free communication problem is typically defined as the joint allocation of actions to each robot to maximize team reward. If robot A has a set of actions A_A and robot Bhas a set of actions A_B , then planning actions consists of selecting a pairwise allocation from $A_A \times A_B$ to maximize team reward. When considering non-free communication the cost of the communication must also be considered. Exact solutions to these types of problems are computationally intractable although many approximate and heuristic solutions exist as well as analyses of special cases that aid tractability.

Work in multi-agent systems has extensively explored representations and solution methods for modeling teams of agents cooperating on tasks with partial observability including interactive partially-observable Markov decision processes (iPOMDP) (Doshi et al., 2009), decentralized models (DEC-POMDP) (Bernstein et al., 2000), POIPSG (Peshkin et al., 2014), and COM-MTDP (Pynadath and Tambe, 2002). These systems typically maintain beliefs about teammates and may incorporate methods for planning when and what to communicate to teammates (Roth et al., 2006). The communication used in these scenarios often relies on a common underlying symbolic representation of the task state, actions, and/or observations, and is usually assumed to occur without delay or misinterpretation. In the human-robot scenario, the robot does not have access to the mental model of the person and must rely on humanlike communication modalities such as speech and coverbal gesture. Additionally, in many assistive HRI scenarios, optimization of task performance is often a secondary consideration to ensuring that the user is engaged and motivated, since the user may have deficits that impair their ability to perform parts of the task. The aim of this work is to contribute toward research in human-robot interaction by developing an effective methodology that allows the robot to issue task-relevant, speech-based communication to guide a human teammate during a dynamic, assistive collaboration, while supporting existing models for robust robot task planning.

2.2.1 Applicability to the Human-Robot Case

The human-robot case is distinct from multi-robot scenarios. First, people are independent of the robot and ultimately make their own decisions about what to do. Thus

to achieve coordination the robot should select the best action in expectation of the human's action. The robot also has no access to the mental model of the task used by its human teammates, making monitoring and communicating about the progress of the task more difficult. Unlike the multi-robot case, the robot cannot rely on a shared representation of the task and networked symbolic communication to transfer new observations or planned actions during the course of the task. Despite these disadvantages, there are a number of aspects of human-robot collaboration that can be used to the robot's benefit. The aspects of human behavior in joint activity that were outlined previously give us some expectation of the types of behavior people will exhibit during joint activity e.g., that they will position themselves to maximize shared visual information and otherwise be helpful teammates. Also, as many approaches for robot learning from demonstration and human instruction have shown, humans can be successfully queried by people and provide useful input to the robot (Rosenthal et al., 2010, Thomaz and Breazeal, 2006). This affords the robot the opportunity to ask for help in cases when things go wrong as well as relying on the person's creativity and expertise to recover from errors in ways that would not be possible in the human-robot case.

2.2.2 Task Classifications

The work in this dissertation is concerned with physically-decoupled tasks where the human and robot are in a shared environment. The focus on physically-decoupled tasks enforces the constraint that any information that the robot has about the person come through hands-off observation or social communication rather than force dynamics. Since this work is also focused on the robot's use of communication via natural human modalities it assumes that the collaboration takes place in a setting where modalities such as speech and gesture can be used effectively. In order to further scope the types of task settings under consideration, the scope of the work is clarified with respect to two different task classification taxonomies: one used in multi-robot task allocation (Gerkey and Matarić, 2004) and the other used in the social sciences to explain group processes (Steiner, 2007).

First, the approach employed in this work is classified in terms of multi-robot task allocation. The taxonomy proposed by Gerkey and Matarić classifies tasks along three dimensions. First, single-task (ST) versus multi-task (MT) robots distinguishes robots that only do a single task at a time versus those that can do multiple tasks simultaneously. The second dimension, single-robot (SR) versus multi-robot (MR) describes problems that require exactly one robot to achieve each task as compared to problems where each task may require more than one robot. Finally, instantaneous assignment (IA) versus time-extended assignment (TA), describes whether the problem allocates robots to tasks instantaneously or contains information for assigning tasks to robots over time. The scope of tasks addressed in this work can be described as ST-SR-TA. The single-task constraint comes from the assumption that many tasks in a shared environment will involve either object manipulation or physically navigating to the right place at the right time. This, combined with the fact that most robots are currently limited in their physical capabilities, i.e., actuation and navigation, leads us to the assumption that the robot will undertake a single task at a time. While the person may decide to multi-task, it is assumed that their actions are decomposable to a series of atomic individual actions that satisfy the single-task assumption. Since multi-tasking can degrade performance, and lead to increased stress on the person, will assume that it is not a primary mode of operation for the user. The single-robot constraint comes from the assumption that joint manipulation or otherwise tightly-coupled interaction between the teammates will not happen. Finally, the time-extended allocation assumption follows from the assumption that the robot and person have a shared goal and know the structure of the task and are working together, undertaking a series of sub-actions to achieve the goal.

Next, the scope of the work is defined according to Steiner's taxonomy of group tasks (Steiner, 2007) to further clarify the types of tasks amenable to this modeling approach. Steiner defines three categorical dimensions across which tasks are grouped that describe: divisibility of the task, type of goal, and interdependence of each individual's inputs in final output.

- Divisibility: By relying on role assignment as the means of coordination, the approach in this work implicitly assumes that the task is divisible i.e., can be effectively separated into subtasks that each agent can perform separately, in contribution to the greater task goals. Unitary tasks, as defined by Steiner, are an interesting special case in which the task is not divisible resulting in one teammate becoming a bystander. In this case, either the robot or the person could not directly perform the task and would become idle. While the approach could support an idle robot by setting the robot's task action set to be empty, the resulting communication generated would only be comprised of recommended actions for the person. This scenario, in which the robot monitors an activity and offers feedback, is an interesting future direction with applicability in rehabilitation and other domains, but isl not considered specifically here.
- Goal types: Two goal types are considered in Steiner's framework optimizing goals, with subjective ratings of the achieved result and maximizing goals, which typically have an objective evaluation metric built-in. From the robot's point of view this distinction is not very useful since ultimately the robot control system requires some quantitative metric with which it can select optimal behaviors in expectation. The tasks on which the approach has been validated belong to the maximizing goal category.

• Interdependence: This area of Steiner's categorization is concerned with how the input from various members is combined to form a single output and how group performance can be related to the best and worst individual performances. The addition of the robot in the evaluation of the approach has been shown to improve performance as compared to either the robot or person alone, largely due to the limitations on the types of tasks the robot can perform. These are typically divisible into separate subtasks with output combined in an additive manner. Adding more agents to these types of scenarios typically will increase the performance of the group.

2.3 Human-Robot Collaboration

Relevant prior work on human-machine collaboration includes approaches from HCI and HRI, as well as from cognitive science and linguistics. Approaches aimed at improving human-robot collaboration typically address problems in two broad categories. One category consists of work aimed at using social communication modalities such as speech and gesture to enable humanlike interaction with a robot. This includes a subset of the broader field of socially assistive robotics (Feil-Seifer and Mataric, 2005, Tapus et al., 2007), in which the main use of the robot is to provide assistance to a user through social interaction. This type of interaction occurs in collaborative settings where the robot is an active participant in team task performance. A second line of work consists modifying characteristics of the robot's task performance, such as path and motion planning and sequential action selection in order to better coordinate activity in the presence of human teammates.



Figure 2.1: A selection of most relevant prior work in human-robot task collaboration: (a) AUR robotic desk lamp (Hoffman and Breazeal, 2010), (b) Chaski (Shah et al., 2011), and (c) Wakamaru humanlike robot (Mutlu et al., 2013).

2.3.1 Collaborative Human-Robot Interaction

An extensive body of work from artificial intelligence exists, developing top-down deliberative approaches for modeling collaborations aimed at establishing and maintaining alignment, assuring coherent discourse, and constructing shared plans for collaboration between a person and an intelligent agent (Grosz and Kraus, 1999, Grosz and Sidner, 1986). Aspects of this model of human discourse and shared planning have been applied in human-robot interaction to produce more humanlike robot behavior by generating appropriate gaze cues during face-to-face interactions (Sidner et al., 2004). These approaches typically depend on certain patterns of behavior, such as turn-taking that may or may not be present during distributed task performance. There has also been extensive work on intent recognition relying on perspective taking or Theory of Mindinspired models to allow a robot to recognize a person's intentional behavior through observation. These types of approaches have been employed to learn the rules of a series of games by clustering estimated intents into roles that a robot can assume (Crick and Scassellati, 2008), as well as to recognize and generate intentional face-to-face meeting initiation behavior (Kelley et al., 2008), and to recognize helping and hindering social behavior via an MDP formalization utilizing inverse planning (Ullman et al., 2010). Other work has treated the collaborative process as a dialog, supporting verbal turntaking and sub-task assignment (Breazeal et al., 2005, Trafton et al., 2005). Speech has been demonstrated to be an effective input modality for commanding a robot (Kollar et al., 2010, Shah et al., 2011) as well as a means by which a robot can provide instruction to a person who is unfamiliar with a task (Sauppé and Mutlu, 2014a). Other work on human-robot collaboration has underscored the need for integration of social cues in understanding situated user behavior and in enriching the social environment in the workplace and achieving coordination with co-workers (Sauppé and Mutlu, 2014b).

Other relevant work demonstrated the ability of a robot to learn simple tasks through human tutelage and collaborate effectively via turn-taking (Breazeal et al., 2009). Still other work has focused on identifying coordination behaviors, particularly eye gaze cues (Mutlu et al., 2013) and implicit and explicit verbal communication (Shah and Breazeal, 2010). In other related work, the focus on enabling effective human-robot interaction by creating human-compatible input methods for teaching and instructing robots including collaborative settings (Chernova and Veloso, 2008, Kollar et al., 2010).

2.3.2 Task Adaptive Collaborative Robots

Another line of related work has focused on enabling better robot co-workers by making robots that work in assembly and manufacturing scenarios more intuitive to users in their task planning, motion planning, and use of humanlike social cues. While not strictly focused on collaborative scenarios, this work is an instructive area of HRI in enabling methods for robots to account for the needs of users while performing realworld tasks. Some existing work has focused on reducing the need for caging existing robotic systems in manufacturing and assembly domains, enabling potentially dangerous robots to be effective co-workers, and developing effective coordination mechanisms for robots by adapting characteristics of the robot's task performance. This involves
creating robots that move safely and predictably around people (Lasota et al., 2014). Task scheduling systems have been developed to change the task planning of a robot teammate in response to actions taken by a human teammate (Shah et al., 2011) as well in settings with complex temporal constraints (Gombolay et al., 2013). Other approaches use Markov decision processes to plan coordinating actions (Nikolaidis and Shah, 2012, Wilcox et al., 2013) in assembly tasks. This prior work has demonstrated that robots that account for the actions of its collaborators when deciding what to do are preferred and perceived as more intelligent (Shah et al., 2011) and that anticipatory action can play an important role in increasing team fluency (Hoffman and Breazeal, 2007).

A growing body of work also exists in adjusting the behavior of robots, primarily in manufacturing settings, to make the robot better anticipate and align its actions with people around them (Hoffman and Breazeal, 2010, Nikolaidis et al., 2015) as well as adjusting robot motion planning to improve its legibility to human observers (Dragan et al., 2013).

This dissertation work develops a methodology for the robot to produce social communication (speech and embodied gesture), to efficiently and intelligently allocate *roles* during a collaboration, and support the user's natural collaboration modalities and preferences. The approach relies on situated robot communication production via humanlike output modalities aimed at coordinating co-located joint activities. Unlike existing work in human-robot collaboration it does not rely on a specific conversational structure such as turn-taking allowing for it to be used in dynamic tasks. The evaluation of the approach demonstrates that providing feedback in the form of coordinating social communication leads to improved team performance and improved subjective evaluation of the robot and the task itself by users.

2.4 Summary

This chapter has given an overview of the relevant research on human behaviors for coordinating joint activity from the social sciences and its implications for work in human-robot contexts. Specifically, the tendency for people to take the perspective of teammates in order to ascribe intent and use a variety of social behaviors including speech and gesture motivates allowing a robot to use these humanlike modalities to coordinate a human-robot joint activity. Next, the chapter reviewed work on multi-robot coordination including a variety of representations based on Markov decision processes, including multi-agent MDPs, hierarchical task networks, and partially observable models that have useful properties for enabling robust decision-making, online learning, and learning from demonstration in complex task environments. The applicability of these methods to the problem of coordinating human-robot activities was discussed, highlighting the problems of differing task models and effective communication between agents in the human-robot case. The types of tasks the work in this dissertation applies to were clarified with respect to two different task classification taxonomies. Finally, a review of existing work on human-robot task collaborations demonstrated many systems that adjust the robot's task behavior in response to the user's behavior and case studies motivating the importance of social communication usage in human-robot tasks. Comparatively little work exists on enabling the robot to use humanlike social modalities. particularly speech, to coordinate joint activity with a person. The contribution of this dissertation is on producing robot feedback to coordinate and guide user behavior and improve situational awareness during collaboration with an autonomous robot.

Chapter 3

Production of Embodied Coordinating Social Communication

This chapter discusses social communication production from a situated robot and how planned coordinating communication, such as allocating a user to a given role, is ultimately translated into physical actions executed by the robot. It details how speech and deictic gesture were chosen as the output modalities of choice and describes an experiment aimed at verifying robot deictic gestural accuracy across many people.

The previous chapter presented three types of communication actions that could be employed by a robot to increase the situational awareness of a human-robot team. These actions are planned at timescales that correspond with changes in the task state and estimates of the user's intended action. To integrate this with the overall robot control system, including control for task performance, requires a communication executive node that integrates with a task control node that is assumed to exist. Some communication actions, such as speech, can easily be executed while the robot is simultaneously

performing a task action, such as navigation or manipulation. If the robot is carrying a box from one place to another, it can easily produce speech while navigating, by playing sound from an on-board speaker. Other communication modalities may conflict with sensing and actuation used to perform the task. If both arms and the head of a humanoid robot are in use for a manipulation task, perhaps requiring sensors in the head to be aimed at a specific target while the arms are moved to manipulate a target object, then neither the head nor the arms could be used to issue a deictic (pointing) gesture without compromising task performance. In other cases it may be possible to blend task and communication actions by combining degrees of freedom that are not in use to produce a gesture that can be correctly interpreted by a person. If the robot is carrying something in one of its manipulators, it could still perhaps use the arm carrying the object to make a gesture as long as the object was firmly grasped. Finally, timing the use of these modalities is important for two reasons: 1) inappropriate timing may make the gestures less human-like and result in interpretation difficulty or dissonance and 2) should the planned communication action change before the previous is finished executing the robot should insert stop words or behaviors to make it clear to the user that a transition has occurred.

The next sections discuss several types of social communication that are applicable in co-located human-robot activity and capable of being produced by a wide variety of robot platforms, with emphasis on speech and deictic gesture production and an overview of other gesture types and the difficulty in designing these gesture for robots.

3.1 Speech

Speech, and the coverbal gestures associated with it, is the primary means by which people coordinate complex shared activity (Grosz and Sidner, 1986) with each other and can be used in tasks where users are co-located or in remote presence applications. There is a large body of work on control of different robots with speech (Kollar et al., 2010) as well as work on how people use speech to coordinate joint activity with each other (Gergle, 2006, Grosz and Kraus, 1999, Ozyurek, 2002). Speech can be used during a collaboration to construct a plan, repair failures, issue feedback, and provide encouragement, among other things. The primary challenge in robot speech production is making sure the robot says things that are helpful to the user, that they are timed correctly, and that they are combined with appropriate embodied gestures when possible. The approach outlined in Section 4.2 generates three streams of communicative intents; narrative feedback, role-allocative feedback, and empathetic feedback. These streams are generated in parallel and must be processed into a final speech action i.e., the robot says a string of text using a text-to-speech system and this output is then multiplexed via a queue system so each verbalization is played without interruption.

In order to convert each communicative intent into a final speech action, the underlying meaning of the intent must be considered and several phrases specified that the robot can use to convey the selected intent. These intents are typically domaindependent with the exception of empathetic output, which can be produced using a set of positive or negative feedback phrases, like "Oh no!" or "Great!" The specific phrases can be collected either through a data collection with people performing the task or specified by the developer considering the roles and actions used. Appropriate phrases for each type of verbal feedback were collected by transcribing the phrases used by people in a pilot experiment of the same task (see Section 5.3). This approach was used on sample tasks where variations in the speech usage of a small number of teams yielded a suitably large number of phrases for each type of feedback to prevent overly repetitive use of phrases during a five minute interaction with the robot. A large-scale data collection of this type of role assignment and narration data was also conducted by administering surveys in which users are shown a simplified diagram of the task and asked to either describe a specific action or to command one of the collaborators to perform a specific action (Section 5.3). This is a useful way to collect a large amount of data on human-like verbal feedback in a given task environment as long as it is amenable to depiction as a diagram, and yielded similar phrases with less variation as compared to the human-human speech transcription approach.

Phrases can also be generated automatically in scenarios where the robot's actions consist of manipulating a specific set of work objects or performing a well-defined set of tasks in an environment. In practice, according to pilot experiments conducted on a sample task (Section 5.3) it was found that people have some common forms of describing responsibilities as they pertain to physical objects or locations. Phrases such as "I'll take care of...", "I've got..." and "Can you get..." followed by the name of an object or location are commonly used. Variations on these types of phrases can be generated by combining a list of synonyms for the relevant objects or behaviors with these generic phrases assigning or accepting of responsibility. Empathetic feedback can also be produced in a task-agnostic manner as it only requires producing a communicative action that expresses positive or negative characteristics. Further research beyond the scope of this dissertation is necessary to determine the efficacy of generic phrases generated using this method and phrases taken directly from human-human interaction.

Once a set of strings representations for each communicative intent is specified, speech can be generated using a text-to-speech system or by prerecording the complete set of relevant phrases. A random phrase is selected from the set of all phrases for each communicative intent when executing the communication. The phrases for roleallocative feedback can be readily adjusted to be more or less polite using methods from Politeness Theory (Brown, 1987). Baseline phrases were used as collected from a pilot study, with no attempt to make them more or less polite. All experiments assumed that the robot's task control policy is static to avoid confounds. The phrases, stored as text strings, are passed to a text-to-speech engine and played in real-time during the task as triggered by the communication planning system. A speech executive node was developed to allow multiplexing the various types of communicative feedback without interrupting individual phrases. Communicative intents generated as the task state changes over time are queued and played in the order received. This prevents the phrases from being played at the same time and also leads to instances of the robot simultaneously issuing a role-allocative communication followed immediately by a narrative communication, for example. In practice, most phrases are a single sentence and thus relatively short in duration, since the task state changes are dictated by task actions that take a certain amount of time, large backups in the queue were never encountered. At this point relevant non-verbal cues can be integrated. Visual attention behaviors were integrated with a humanoid robot to provide visual cues to the relevant target of a particular phrase. In general, care must be taken integrating gestures if they employ parts of the embodiment that are also employed by the robot for task performance, as the gestures could potentially have a negative effect on the robot's task performance.

For example, if the role allocated at a specific point in time calls for an action a, rather than try to describe the complete role to the person, the robot could ask the person to perform action a at that point in time. Over time the action preferences of the role can be elucidated to the person. This works similarly for narrative feedback where the robot could say it is going to execute some role, r, represented by a policy p. It could periodically describe the actions it is executing at a given time or it could refer to the role in general if this is clear in the task context. If a task typically involves a clear separation of duties between the team members the robot could say "I'll be the painter" rather than the action-level phrase "I'll paint this piece". Combining speech with coverbal behaviors like beat gestures and deixis when reference physical objects

or locations is also important in humanoids as it makes them appear more intelligent and provides other embodied social cues for people to understand (Huang and Mutlu, 2013b).

3.2 Deictic Gesture

Multi-disciplinary research from neuroscience and psychology has demonstrated that human gesture production is tightly coupled with language processing and production (Kelly et al., 2009, Mayberry and Jaques, 2000). There is also evidence that gestures are adapted by a speaker to account for the relative position of a listener and can, in some instances, substitute for speech functions. This substitution effect has been demonstrated in deictic gestures by studying performance on a target disambiguation task, where it is found that deictic speech combined with deictic gesture offered no additional performance gain compared to one or the other used separately (Bangerter and Oppenheimer, 2006, Louwerse and Bangerter, 2005).

These findings have important implications for the field of human-robot interaction. Robots interacting with people in a shared physical environment should be able to complement their use of verbal feedback during the course of an interaction with other gesture types. To make this possible, it is necessary to gain an empirical understanding of how to map well-studied human gestures to robots of varying capabilities and embodiments. Specifically, we are interested in identifying variables for proper production of robot gestures, to maximize the likelihood of some desired interpretation by a person. In general, this is difficult for the same reasons that processing natural language is difficult; many gestures are context-dependent and rely on accurately estimating a mental model of the scope of attention and possible intentions for peoples' actions given only low-level perceptual input. Deictic gestures, however, are largely consistent in their mapping to linguistic constructs, such as "that" and "there", and serve to focus the attention of observers to a specific object or location in the environment, or perhaps to indicate an intended effect involving such an object e.g., "I will pick up that." These characteristics, while simplifying their interpretation and production, also make the gestures useful for referring to objects and for grounding attention. Intentional analysis and timing are still challenging problems, except in the context of performing a specific pre-determined task. Both recognition (Cipolla and Hollinghurst, 1996, Kortenkamp et al., 1996, Nickel and Stiefelhagen, 2007, Pook and Ballard, 1996, Wong and Gutwin, 2010) and production (Hato et al., 2010, Marjanovic et al., 1996, Sugiyama et al., 2006) of deictic gestures have been studied in human-human, human-computer, and human-robot interaction settings. This work adds to this field as a step toward obtaining an empirically grounded HRI model of deictic gestural accuracy between people and robots, with implications for the design of robot embodiments and control systems that perform situated distal pointing.

A study of the literature on human deictic pointing behavior suggests a number of possible variables that could potentially affect the robustness of referent disambiguation, i.e., accurately pointing at a desired target, when having a robot employ deictic gestures. The physical embodiment of the robot constrains the appearance of the part of the robot used for pointing. Pointing with a blunt or irregularly shaped object has lower resolution than a sharp, pointed object with a clear vector interpretation (Bangerter, 2004, Matarić and Pomplun, 1998). It has also been shown that people take into account the position and orientation of their audience when staging a gesture in order to improve the audience's interpretation (Ozyurek, 2002); this consideration of audience perspective is even a core concept in hand-drawn (Lasseter, 1987) and computer animation (Thomas et al., 1995). Assuming mobility, a robot could also relocate or reorient itself relative to its audience or to the referent target to improve viewer's interpretation accuracy.



Figure 3.1: The PR2 mobile manipulator attempting to produce a "step back" gesture. Producing gestures on a robot platform that was not designed with this purpose in mind is a difficult problem and currently requires skilled hand animation.

Other considerations beyond physical variables include parameterizing the timing and appearance of the gesture itself. People in most cases attempt to minimize effort when performing deixis and thus do not tend to make maximally expressive gestures unless necessary, such as when an object is far away (Bangerter, 2004). Most robots, however, point without accounting for distance to the target by using a more or less constant arm extension or head gaze that is reoriented appropriately (Mutlu et al., 2009). Finally, since the gesture is grounded with respect to a specific referent in the environment, the robot must be able to correctly segment and localize visually salient objects in the environment at a similar granularity to the people with whom it is interacting. Biologically inspired methods to assess and map visually salient features and objects in an environment exist (Itti et al., 1998, Walther et al., 2010), as do models of human visual attention selection (Desimone and Duncan, 1995), but the role of visual saliency during deictic reference by a robot is largely uninvestigated.

3.3 Other Gesture Modalities

Besides deixis, other gesture modalities could be used to issue feedback from the robot to the person during the course of a joint activity. The other types of gestures commonly used by people in interactions include iconic, metaphoric, emblematic and beat gestures (McNeill, 2008, Mehrabian, 1977). Iconic gestures are closely related to specific linguistic terms and used to illustrate concrete concepts, such as moving one's hand up and down in a hammering motion to indicate a hammer. Metaphoric gestures are similarly used to illustrate abstract concepts such as presenting an idea with an open hand. Emblematic gestures have well-defined accepted meanings but may be culturallydependent such as a thumbs-up gesture, which is interpreted to mean "all is well" in the United States, despite having the opposite meaning in other cultures. Finally, beat gestures are produced coverbally and used in a rhythmic manner while speaking and are associated with various conversational cues.

Since people working in teams benefit from shared visual information (Kraut et al., 2003), enabling the robot to use these additional gestures might improve teamwork. Unfortunately, producing these types of gestures on a robot can be quite difficult and is highly-dependent on the robot's embodiment. Robots designed with a specific purpose besides producing coherent social gesture, such as the Willow Garage PR2 in Figure 3.3 designed to be a mobile manipulation platform, may lack the degrees of freedom necessary to produce certain gestures in a way that will be readily understood. Besides production difficulties and issues with cross-person and cross-cultural interpretation, there is also the problem of correctly selecting the gestural repertoire of the robot with respect to the task. As an example, since iconic gestures are directly related to physical concepts, such as the work items or other objects in the environment, they are dependent on the task domain. Similarly, although emblematic and metaphoric gestures have meanings independent of the task-domain, their use is dependent on the concepts the



Figure 3.2: Bandit pointing with its (a) head; (b) straight-arm; (c), bent-arm; and (d) and head+arm.

robot is trying to convey and thus they may not be applicable in all scenarios. For these reasons, deictic gestures were selected as a promising means for coverbal gesture production that are useful in a wide variety of task settings. We performed a detailed analysis of robot deictic gesture production to verify the stability of the interpretation by people under a number of conditions.

3.4 Evaluating Human Interpretation of Robot Deictic Gestures

We performed a factorized IRB-approved experiment over three robot pointing modalities: the head with 2 degrees-of-freedom (DOF), the arm with 7 DOF, and both together (i.e., head+arm) with two saliency conditions: a blank (or non-salient) environment and an environment with several highly and equally visually salient targets, using an uppertorso humanoid robot (Bandit). For an overview of the robot pointing gestures consult Figure 3.4. The results, particularly the modality condition, may be specific to Bandit, but a similar, smaller, test with a human performing the pointing gestures was also conducted for comparison.

3.4.1 Study Design

In the experiments, the participant was seated directly facing Bandit from a distance of 6 feet (1.8 meters). The robot and the participant are separated by a transparent, acrylic screen measuring 12 feet by 8 feet (2.4 by 3.6 meters) (see Figure 3.4.1). The screen, thus, covered a horizontal field of view from approximately -60 to 60 degrees and a vertical field of view -45 to 60 degrees. Two screen conditions were tested, one in which the screen was blank (the non-salient condition), and another in which salient objects were placed on the screen at target locations. This was to test to see if the participants exhibited a tendency bias perceived referents towards salient objects. The robot performed a series of deictic gestures and the participant was asked to estimate the referent location on the screen. All gestures were static and held indefinitely until the participant estimated a location, at which point the robot returns to a home location (looking straight forward with its hands and its sides) before performing the next gesture. Participants were given a laser pointer to mark their estimated location for each gesture. These locations were recorded using a laser rangefinder, which was placed facing upwards at the base of the screen. For each gesture, an experimenter placed a fiducial marker over the indicated location, which was subsequently localized using the rangefinder data within approximately 1 cm. The entire experiment was controlled via a single Nintendo Wiimote, with which the experimenter could record marked locations and advance the robot to point to the next referent target.



Figure 3.3: The experimental setup showing a participant indicating the perceived location of the Bandit robot's referent target with a laser pointer among a set of salient targets.

The face-to-face nature of the experiment, as shown in Figure 3.4.1, was chosen intentionally although other work in gesture perception (Bangerter and Oppenheimer, 2006) has tested human deictic pointing accuracy when the pointer and the observer are situated more or less side-by-side, observing a scene. In our work, and in most humanrobot interaction settings the robot is typically facing the participant, rather than sideby-side, which, in terms of proxemics, is more likely to occur when coordinated motion and interaction are concurrent (e.g., the robot and human walking together while an interaction is taking place (Desimone and Duncan, 1995)). For this reason, our design tests the face-to-face scenario, since it is more applicable to the types of proxemic configurations we have tended to encounter during our prior work

3.4.2 Participants

A total of 30 runs of the experiment were conducted, with 17 (12 female, 7 male) participating in the non-salient condition and 12 (7 male, 5 female) participants in the salient condition. In total, approximately 3600 points were estimated and recorded (Table 3.1 for a breakdown by condition). The conditions are close to equally weighted with the exception of the non-salient arm-only and head+arm conditions, which was done initially to allow for comparison of the cross-body versus away-from-body arm gestures. Participants were recruited from on-campus sources and all were undergraduate or graduate students at the University of Southern California from various majors. The participants were roughly age- and sex-matched with an average age of 20. The data collected for each run included a log file recording the desired target on the screen (i.e., the desired location the robot should have pointed to), the actual target on the screen for each modality (i.e., the location the robot actually pointed to), and the perceived point as indicated by the participant and recorded by the laser rangefinder. We also captured timing data for each point and video of the sessions taken from a camera mounted behind and above the robot.

Data Counts Per Experimental ConditionHeadArmHead+ArmNon-salient5651175298Salient800801811

Table 3.1: Samples collected in the deixis experiment for each condition; arm is overrepresented in the non-salient condition to compare away-from-body and cross-body

3.4.3 Hypotheses and Outcome Measures

Given the large number of possible variables when considering how to best reference, via gesture, a point in the environment with a particular robot embodiment, we conducted an initial pilot study to test distance and angle to target, distance and angle to audience, and modality for a face-to-face interaction between a person and our upper-torso humanoid robot, Bandit. No strong correlations emerged in early testing, so we narrowed the set of conditions and hypotheses. The conditions tested include head (Figure 3.4a), away-from-body, straight arm (Figure 3.4b), cross-body, bent arm (Figure 3.4c), and combined head and arm (Figure 3.4d).

3.4.3.1 Modality

- **Hypothesis 1:** The straight-arm modality will lead to more accurate perception since, when fully extended, it is the most expressive and easily interpreted as a vector from the robot to the screen.
- Hypothesis 2: Away-from-body, straight arm (Figure 3.4b) gestures will be easier to interpret than cross-body bent-arm (Figure 3.4c) gestures, since it is staged in front of the robot's body rather then laterally (Lasseter, 1987, Thomas et al., 1995).
- Hypothesis 3: The head modality will have a higher error rate since Bandit's head does not have movable eyes, leading to an ambiguous point of reference.
- **Hypothesis 4:** Using both modalities together will reduce error relative to a single modality, since participants have two gestures on which to base the estimate.

3.4.3.2 Saliency

- Hypothesis 5: In the salient condition, the salient objects will affect people's interpretations of the points, given that people exhibit a tendency to to the salient objects (Itti et al., 1998), thus reducing error for points whose targets were on or near markers.
- Hypothesis 6: There will not be a significant difference in the performance of each pointing modality when comparing the salient and non-salient conditions.

3.4.4 Results

A two-way analysis of variance (ANOVA) of the resulting data was performed with modality and saliency as the independent factors. Both were found to have significant effects on the angular error between perceived and desired target points as well as on perceived and actual target points (Table 3.2). Additionally, the interaction effects between the modality and saliency factors were found not to be significant. Mean angular error computed from the perspective of the person and confidence intervals are shown in the graphs in Figure 3.4 and Figure 3.5. We used angular error as a metric to effectively normalize for different distances to target. For comparison purposes, human perceptual error when estimating human pointing gestures (arm or eye gaze) has been measured to be approximately 2-3 degrees for people up to 2.7 meters apart (Bangerter and Oppenheimer, 2006). We conducted post-hoc analysis using the Tukey's honestly significant differences (HSD), which revealed that mean error tends to be about 1.5 degrees higher for arm points with p < 0.01, and that using both modalities tends to outperform the arm modality in most cases with the means differing by 1.8 degrees in the salient case and 2.0 degrees in the non-salient case with significance of p < 0.01. The arm alone, however, performed equally poorly in both saliency conditions.

F and P Values for 2-Way ANOVA		
Condition	Perceived-Desired	Perceived-Actual
Saliency	$F = 4.53, p < 0.03^*$	$F = 11.4, p < 0.01^*$
Modality	$F = 9.6, p < 0.01^*$	$F = 5.0, p < 0.01^*$

Table 3.2: Results of deixis experiment ANOVA by condition.

To compare cross-body bent-arm versus away-from-body straight-arm gesture, we looked at arm points in the non-salient case. Partitioning them into two sets, depending on the side of the screen they were on, resulted in 450 cross-body points and 440 straightarm points. Conducting a one-way ANOVA with arm gesture type as the dependent



Figure 3.4: (a) Mean angular pointing error and (b) Mean angular error between perceived and actual targets.



Figure 3.5: (a) Mean angular error between perceived and desired targets and (b) Mean angular error between perceived and desired targets for cross-body and straight-arm points.



Figure 3.6: (a) Mean time from start of pose to marking of estimated point and (b) Mean error by saliency condition, human pointer.

variable, we find there is a significant difference between the straight-arm case (M = 5.6 degrees, SD = 3.7) and the bent-arm case (M=10.4 degrees, SD=7.9) with p < 0.001 (Figure 3.6a). We also obtained similar results when conducting a full 3-way ANOVA with the other two conditions, although there were interactions between some of the factors in addition to significant main effects, likely due to the high variance in arm accuracy.

To assess whether the accuracy of the other modalities varied with angle to target, we first fit a linear regression model with angular error as the dependent variable and the desired target as the independent variable. The resulting model did not perform well upon cross validation suggesting that the error was nonlinear in nature. To cope with the nonlinearity, we then binned points by angle to target in 9 uniform intervals covering the extent of the screen. We then performed an n-way ANOVA with target x and y coordinates as a conditional factor, and found that the head-only and head+arm



Figure 3.7: (a) Mean angular error with respect to horizontal target position as seen from the participant's perspective and (b) Mean angular error with respect to vertical target position as seen from the participant's perspective.

conditions were significantly better (p < 0.001) in the center of the screen with error increasing about halfway to the edge before leveling out. These effects were largely symmetric for the head-only and head+arm conditions. The arm-only condition, as described above, was asymmetric and was significantly worse in almost all cases except for the middle of the left side of the screen corresponding to away-from-body straightarm points. Finally, all modalities tended to result in more erroneous estimates in the lower extreme of the screen. Figure 3.7a depicts a smoothed cubic fit to the entire dataset for each modality with respect to target x-coordinates and Figure 3.7b with respect to the target y-coordinates. Each is pictured from the participant's perspective, meaning cross-arm gestures correspond to the right side of the horizontal graph. The average time taken to estimate each point was nearly a second faster in the non-salient case (M=6.6, SD=2.6) than in the salient case (M=5.8, SD=2.3). This effect was found to be significant with p < 0.001, while there was no significant difference for the modality conditions (see 3.6a).

For the human pointing phase of the experiment, we collected a total of 70 data points from two participants. While this is not a considerable amount of data and should include more participants before drawing conclusions, we did find a significant (p < 0.09) effect when comparing the two saliency conditions. For the salient condition (M=2.23, SD=2.46), the error was small enough that both the pointer and the observer were able to "hit" all the salient targets, while the non-salient condition (M=3.74, SD=2.74) resulted in a 60% increase in error, as noted in Figure 9. When plotting the error versus intended referent x-coordinate, note that points directed at the center of the screen tend to result in lower perception error than points directed between the center and the periphery. Overall, the error in estimating human-produced points appears to have a similar profile to that of the robot-produced points; however, more investigation is necessary.

In the responses to the survey, participants in the non-salient condition estimated that their points were within an average of 28 centimeters (11 inches); this is very close to the mean error of 27 centimeters we found in practice. There was no significant difference between participants' estimated error when comparing across the two conditions. Pointing with the head-only and with the head+arm were preferred by the majority of the participants, with only 4 (or 16%) stating a preference for the arm modality. When asked if there was a noticeable difference in straight-arm points versus the bent-arm, 65% said there was with the remainder not seeing a difference. Ten out of 12 (or 83%) of the participants in the salient condition said that the markers would have an effect on their estimate of the referent target.

3.4.5 Discussion

The mean error, as computed (using the perceived point and the desired target points), tells us how close to a desired target (either a randomly chosen one in the non-salient case or one of the markers in the salient case) the robot was actually able to indicate. The performance of the head-only in the salient condition is improved by approximately 1 degree, the arm-only is not appreciably different, and head+arm exhibits only a modest gain. This suggests that the snap-to-target effect that we expected to see when salient objects were introduced is modest at best, resulting in a best-case improvement of approximately 1 degree. This is also seen when we consider the mean perceived-actual error and find that nearly every condition is slightly more accurate. This suggests that participants estimate the referent to be closer to the actual point the robot is physically indicating than the nearest salient object. This could be potentially useful because it allows us to consider pointing without having to assess scene saliency beforehand. That is, if we have not specified the referential target of a point *a priori*, through some other means, such as verbal communication or previous activity, people tend to evaluate the point in an ad hoc manner by taking their best guess. When disambiguating referents, if there are unknown salient objects in the environment, we can anticipate their effects on the perception of a given gesture to be small enough in most cases that a precise point to our actual target should suffice to communicate the referent.

As we hypothesized, the modalities did result in different pointing accuracy profiles. When considering modalities, pointing with the head+arm does appear to perform appreciably better than either the arm-only or the head-only, in most cases. One possible explanation of this is that it more closely emulates typical human pointing, in which people tend to align arm gestures with their dominant eye (Bangerter and Oppenheimer, 2006), or that multiple modalities provide more diverse cues that indicate the referential target resulting in better priming of the viewer to interpret the gesture. The poor performance of the arm in the salient condition was also somewhat unexpected. This might be due to its higher actual error compared to the head, which we discovered to be related to a small dead-band created in the Bandit firmware that can become compounded over each joint in the arm, resulting in poor overall performance. Another source of the error could be use of the cross-body arm gesture, which, while equally weighted, resulted in nearly twice the perceptual error as compared to the away-frombody arm. This might be a result of the reduced length of the arm, which forces people to estimate the vector based on only the forearm versus the entire arm as in the awayfrom-body case. Another explanation is that the gesture is staged against the body, that is, with minimal silhouette and is, thus, more difficult for people to see. In either case, roughly one-third of the participants did not notice a difference in the arm gestures while their performance was, in fact, affected. This illustrates the impact that gesture and embodiment design can have on interpretation, and underscores the need to validate gestural meaning with people.

When considering the horizontal and vertical target position analyses, we see that people are best at estimating points directly between the participant and the robot. Performance then drops off when the target is located laterally, above, or below. This effect could be due to a field-of-view restriction, preventing the viewer from seeing both the robot's gesture and the target at the same time in high acuity, foveal vision. Estimating these points then requires the viewer to saccade their head between the two points of interest. We believe the slight improvement at the far periphery for some of the modalities is due to the fact that we informed participants that the points would be on the screen, thus, creating a bound for points near the screen edges.

The results of our smaller-scale investigation of human pointing did find that the salient condition resulted in approximately 1.5 degrees less error than in the non-salient

condition, which is consistent with our finding using the robot pointer. Also, the 2degree perceptual accuracy that we found when testing a human pointer seems to agree with prior study of people conducted in relevant literature. It is also worthy of note that, although the deictic pointing performance of the robot is several times worse than what we saw in the human experiment or would expect from literature, we can use the estimate of our resolving power (i.e., the minimum angle between referents that we could hope to convey) to inform controller design and ensure that the robot re-positions itself or gets close enough to prevent these effects. The salient condition also resulted in a 16% increase (or approximately 1 second) in time needed to estimate the gesture. This is intuitive, as the participants were presented with more stimuli in the form of the salient objects and, thus, take some extra time to ground the point, possibly checking to see if it is coincident with any objects first. This information could be useful in developing methods for effective timing control.

3.5 Summary

In this chapter several modalities for physically instantiating planned coordinating abstractions using a physical robot were discussed with emphasis on speech and deictic gesture use. Strategies for handling timing and observability issues using a system of queuing and filler word insertion were presented. A study of the accuracy of human perception of robot-produced deictic gestures including arm-only, head-only and combined head and arm gestures was presented demonstrating that robot deictic gestures can be used to disambiguate targets but are likely less accurate than their human-produced counterparts. The difference in performance of the various modalities motivates modeling and accounting for pointing inaccuracy during gesture production to ensure the clarity of the robot's intended referential target to on-lookers.

Chapter 4

Approach to Coordinating Social Communication

In this chapter the task representation is presented as well as the approach for modeling user activity using the notion of roles and planning social communication from the robot to the person during task performance including narrative, role-allocative, and empathetic communication.

The main contribution of this work is in the area of planning and execution of robot coordinating communication actions during collaborative task execution. Coordinating communication in a collaborative task can take many forms, depending on the modalities used (speech, embodied gestures, pre-determined commands, etc.), the time-scale over which the coordination occurs (every atomic action, major sub-tasks, or extended task trajectories in the case of repeated/recurring tasks), and the specificity of the communication with the task. This dissertation focuses on using explicit verbal communication to issue relevant coordinating social communication. This is in contrast to more implicit signaling mechanisms, such as adjustments to the robot's task performance or anticipatory action.

The goal of the approach is to plan coordinating social behavior that will improve the situational awareness of the robot's teammate during the course of a joint activity. Through this improved situational awareness, the approach seeks to guide the person to better decision-making and increased performance by reducing the idle time that results from the person's uncertainty of the robot's beliefs and intentions. More specifically, by providing verbal feedback based on the team's performance of the task the work aims to improve the situational awareness of the user. This is complementary to existing approaches that adjust the task performance of the robot to allow it to better interleave its actions with those of a teammate. In addition to adjusting the robot's behavior in response to the person's action, a methodology for providing explicit verbal feedback can also play an important role in mediating human-robot joint activities by supporting human-like coordination methods such as speech. This approach to communication opens up some interesting possibilities in adjustable autonomy that allow for a robot to guide an interaction by explicitly telling a user what to do or to provide implicit feedback, allowing a user more flexibility in action selection. These types of interactions are applicable in teaching or coaching scenarios in which the robot guides the person to a desired goal through a combination of physical task performance and verbal feedback as well as more traditional collaborative work scenarios.

As noted in in Section 1.2.2, social science literature indicates many types of human communication behavior that a collaborative robot could make use of, including attentional cues to indicate an area of focus, staging actions to maximize shared visual information, gestural and speech cues indicating intentional goals or instructions, and coaching actions such as feedback, encouragement, and empathetic displays to build team rapport. Based on a review of the relevant social science literature covering human-human collaborations, observations of person-person task collaborations in our experimental setting, and the capabilities of the robot systems currently available, this work focuses on using speech as the primary form of coordinating communication. The approach involves the combination of three robot communication components: 1) the robot's self-narration of its activities, 2) role allocation suggestions for the user, and 3) empathetic displays when positive and negative events occur. This combination provides a balance of information aimed at improving the human collaborator's situational awareness. The second component, offering suggestions to the user, is particularly interesting because it informs the user that the robot is monitoring the user's progress and evaluating the world from their perspective and allows the robot to potentially influence the joint decision making of the team.

In order to achieve these goals, the task control and communication approach consists of three constituent components (shown in the system diagram in Figure 4.1):

The task control component of the approach is responsible for enabling the robot to actively participate in performance of the task (shown in dark green in Figure 4.1). Few restrictions are placed on this system, as the representation of roles used to generate communication is compatible with tasks represented as Markov decision processes.

The following sections describe how the approach can be applied to a broad range of cases in which Markov-based task representations are used, such as full multi-agent MDPs with a joint policy, hierarchical models, and partially-observable environments.

Since the work is focused on producing social behavior, it aims to be compatible with a range of MDP-based task models and policy representations. Even if the robot's behavior is suboptimal for a given task, the approach should at least generate communication that is consistent with the robot's behavior and assists the user in making coordinated decisions during the task. Users' evaluations of the entire robot system are dependent on what the robot does and what the robot says. Since these components are not easily separable, the evaluation of the work is conducted using a control component that competently performs parts of the task and is designed to assess the effect of the



Figure 4.1: A brief over of the major components of the approach including the user policy recognition system, the communication planner with three types of communicative intents and the communication realizer.

coordinating communication through user studies. The robot's controller is assumed to monitor sensor input and decides which of a set of discrete actions is best for the robot to take over time. These discrete actions are then executed by the robot by employing one or more lower-level controllers.

4.1 Human-Robot Task Representation

The approach is comprised on three main components: the user role recognition component, the communication planning component, and the communication realizer. These components require as input a stream of recognized user actions and information about the current task state. Each component is described in detail below.

Human activity modeling and recognition: This component of the approach is responsible for tracking and recognizing intentional human behavior over time (Figure 4.1). It is necessary to explicitly model the person's behavior so that effective communicative actions can be planned in expectation of future robot and user action. This component takes the output of whatever sensing systems are available to the robot and recognizes the performance of each of a set of discrete actions, given *a priori*, that are assumed to be available to the person. There are existing approaches to segmenting, recognizing, and modeling human activity in various contexts as surveyed in Section 4.2. In this work a simplified methodology was developed to serve these purposes in the experimental task settings employed for evaluation, although the work can also be integrated with and improved by state-of-the-art recognition and task planners. The user action recognition methodology used in the experimental evaluations of the approach are relevant to task settings where objects can be reliably tracked and are suitably separated and distinguishable with high accuracy.

There are existing approaches to segmenting, recognizing, and modeling human activity in various contexts. In most real-world environments there are manageable subsets of objects or features that will play a role during performance of the task. Thus constructing an environmental map involves tracking a subset of key objects throughout the collaboration. This could be accomplished with a tailored detection system for a specific task domain or could be accomplished via some unsupervised activity segmentation and recognition process, wherein the system observes someone performing the task while tracking salient features and infers what is important for each segment of the activity (Hamid et al., 2009, Kirsh, 1995). In cluttered or partially observable environments multiple robots can be used to provide better monitoring of the task space and the robots' human teammates.

As high-fidelity person detection and tracking methods become more capable and corresponding sensing technologies, such as the Microsoft Kinect, improve in performance, tracking a person's actions in a room or building-type environment becomes increasingly viable. As the approach relies on the recognition of distinct task actions, it is expected that the work could be integrated with, and improved by, a state-of-the-art human action recognition system. **Communication planner and executive:** This component of the approach is responsible for planning various types of coordinating social communication during the course of the task (Figure 4.1). This takes as input a stream of recognized human actions as well as information from the robot task execution component, including the policy that the robot is currently executing, and generates three types of communicative intents: 1) narrative feedback about the robot's action, 2) role-allocation suggestions for the robot's teammate, and 3) empathetic feedback when large positive and negative transitions occur. These communicative intents are symbolic; for instance, a role allocative intent might be to assign the user to a specific role. These must then be realized in the form of actual robot behavior including speech and gesture. This process is performed by a communication realizer, which receives streams of communicative intents and plays them sequentially from a queue to prevent incomplete playback.

The coordination problem is formulated as follows. First, the robot represents the task as a Markov decision process (MDP) to plan its actions. The task model is defined as follows: $M_t = S_{task}, A_{task}, T, R$, where S_{task} is the finite state of the environment, A_{task} is the set of task actions the robot can execute, T is a function giving a probability distribution over states for executing a given action in a given state, and R is a reward for each state. In most tasks the state will consist of a set of random variables $S = \{X_1, X_2, X_3, ..., X_n\}$ that each take on a set of discrete values and describe relevant world state information for conducting the task. These types of models were selected as the basis for the approach because they have been studied extensively (Kaelbling et al., 1996, Shani et al., 2013) and provide many useful properties for planning communication related to task performance.

This task formulation has also been successfully used previously to conduct a humanrobot task collaboration without communication (Nikolaidis and Shah, 2013). The disadvantage to this approach is that the user's activity is not modeled directly. In this case, the transition function captures the robot's uncertainty due to changes the human collaborator might make in the environment in addition to the traditional sources in single-agent activity such as sensor or motor noise or other uncertainty. This may be sufficient for some human-robot collaboration settings, where the user ends up following a stationary policy in the robot's model. Since the user does not share the same underlying representation of the environment as the robot, the user's behavior may not appear consistent and predictable to the robot. Another disadvantage of not having a direct model of the user's behavior is that uncertainty in the transition function comes from both environmental and human sources. This combination makes it an inadequate basis for planning communication actions that specifically relate to reasoning about what the user is doing.

The work in this dissertation also assumes that the robot is an active participant in the task and has a policy mapping the robot's states to actions $\pi_{robot} : S \to A$. Since the primary goal is to produce coordinating social behavior from the robot that clearly conveys the robot's goals, few restrictions are placed on the policy used. This policy could be manually developed or learned from data and is not required to be optimal, although it should result in the robot providing the user with some assistance during the task, i.e., the robot should be an active collaborator. In later sections, application of the approach in scenarios where the robot is permitted to learn the task and communication at the same time are detailed. Given a policy and the model of the task, the goal of the approach is to generate coordinating social communication that explicates the robot's behavior and allocates actions to the person that will be beneficial to the team. A complete overview of the components required for a complete instantiation of the approach is show in Figure 4.1. A description of the approach assuming a fully-observable environment follows with a demonstration of how these constraints can be relaxed in Section 4.5.

55



Figure 4.2: A system diagram showing the interaction of the components of the approach with human activity tracking in green, robot task planning in pink, robot communication planning in blue, and external input and output modules in orange and purple, respectively.

4.2 Representing and Recognizing Human Activity

In order to formalize useful communication from the robot to the person, first a definition formal definition of the human construct of *roles* needs to be established within the task model and a mechanism for planning coordinating social communication actions. The goal is for the robot to provide verbal feedback to its collaborating partner: toward that end the approach focuses on communicating intended action via the mechanism of roles. People use explicit and implicit roles in team activity and other organized behavior. For example, if a group of people is tasked with assembling a wooden box, they might assign a sawing role, a hammering role, and a painting role that different people assume during the course of the activity. Although no formal definition of human role use exists and the mechanisms by which people partition a given activity into a set of roles are typically described informally (Smith et al., 2001), the following observations can be made given the types of task environments under consideration:

- First, roles are used by people to limit the scope of their responsibility during an activity. Intuitively, a person playing a particular role will have a set of actions expected of them in certain situations that might be completely different from a person playing a different role.
- Second, these responsibilities may consist of more than one discrete action, typically involving work objects.
- Third, people can change roles during the course of an interaction.

Given these requirements for the role representation, the roles that the person can assume during the task modeled as a set of policies, $\Pi_{roles} = \pi_1, \pi_2, ..., \pi_n$. This set of policies is domain-dependent and not assumed to be optimal or otherwise sufficient, when executed alone, for solving the task. Rather, the policies in the set are a means of quantifying the user's likely behavior during task performance, and of grouping related actions according to the roles people perform them in for a given task. A distribution, D, is maintained over Π , based on the likelihood that the user is executing each policy and updated by observing state transitions over time, inferring user actions, and reweighting policies that agree. For ease of explanation, these roles are assumed to be defined over the same state and action space employed by the robot (see Section 4.5 for relaxation of this requirement).

The process of generating coordinating communication proceeds as follows: 1) the user's role usage is inferred from recognized actions over time, 2) given the robot's policy, the inferred user role, and task state information communicative intents for each of the three types of communicative feedback are produced, and 3) verbal feedback and corresponding embodied gesture is executed to convey the communicative intent to the user.

4.2.1 User Policy Recognition

To accurately infer the role, $r \in \Pi$, that best describes the user's action selection preferences, it is assumed the robot has access to some means of user action recognition. Specifically, it needs to know when an agent (person or robot) takes an action $a \in A$ at time t. Human action segmentation and recognition is a challenging open problem with relevant existing work from computer vision Hamid et al. (2009). A heuristic action recognition system that makes use of this representation of the state of the task as a set of discrete features:

$$S = \{X_1, X_2, X_3, \dots X_n\}$$

. In order to execute its policy in a fully observable environment, the robot must take sensor input and determine the value of each random variable so that it can determine the current state. Given this capability, some of the X_i can be directly associated with actions that the agents can take, while others are consequences of the environment. For example, if a switch must be pressed during the course of a task, then we might have an action $a \in A_{task}$ called *FlipSwitch* and an associated random variable $X_{switch} = \langle Off, On \rangle$, describing the state of the switch. The heuristic action recognition system employed in this work determines the subset of:

$$X_{recognizable} \subseteq \{X_1, X_2, X_3, \dots X_n\}$$

that are associated with state changes that are attributable only to the person or the robot. On each update, if the current inferred value of one of these variables is not equal to its previous value, $X_i^{(t-1)} \neq X_i^t$, it is assumed that one of the agents made the change. The action is then attributed to the agent closest to the object at the current time, distinguishing actions performed by the robot from those performed by the user.

This heuristic method works well in simple scenarios with straightforward mechanics, where the physical objects and locations where specific subtasks are performed are distinct, and when the sets of random variables are carefully selected, but does not transfer well to all spatiotemporal actions that a person and robot might undertake. The overall approach could benefit from integration with a more principled state-ofthe-art action recognition method; this is beyond the scope of this work.

Given the stream of recognized actions, the system next attempts to recognize the role the user is employing at the current time. This process is detailed in Algorithm 1. We first discard all recognized actions as being taken by the robot, leaving only the person's actions. For each $a_{rec} \in A_{task}$, we perform a Bayesian update of a multinomial distribution over the set of role policies:

$$p(R = r_i | A = a_{rec}) = \alpha * p(A = a_{rec} | R = r_i) * p(R = r_i)$$

59

The likelihood of the action being taken given a certain role r_i is selected based on whether the policy would have executed the action in the previous state, with a specified weight, as follows:

$$p(A = a_{rec} | R = r_i) = \begin{cases} p_{agree}, & \text{if } r_i(s_{prev}) \equiv a_{rec} \\ \\ p_{disagree}, & \text{otherwise} \end{cases}$$

This likelihood function accounts for the compatibility of a policy from the set of user policies with the evidence provided by the user action recognition component. The constants p_{agree} and $p_{disagree}$ can be obtained empirically and depend on the accuracy of the user action recognition system used. After normalization, the inferred user role is selected by $\pi_{user} = argmax_rp(r|a_{rec})$. This estimate of the user's role is thus based on their action selections at previous times during the task and is dependent on individual actions taking around the same amount of time. If this is not the case, the likelihood of the action given the role could be expanded to include more than one past state, with an appropriate discount factor. In our validation, the component maintains a multinomial distribution over the set of user roles based on the likelihood that the user is executing each policy. On each update, policies are reweighted based on their agreement or disagreement with the recognized action.

Algorithm 4.1 User role recognition algorithm

Require: Recognized human action $a_{rec} \in A_{task}$, previous state $s_{prev} \in S_{task}$
for all roles $r = \{\pi_1 \pi_n\}$ do
if $r(s_{prev}) = a_{rec}$ then
$p(r a_{rec}) = 0.95$
else
$p(r a_{rec}) = 0.05$
end if
$p_r \leftarrow p_r * p(r a_{rec})$
end for
return $\arg \max_r p(r)$
4.3 Planning Robot Feedback

Based on this recognized action and the robot's own policy, three types of verbal feedback are generated: *self-narrative, role-allocative,* and *empathetic.* A wide range of verbal communication could be produced to improve team coordination as reviewed in Section 1.2.2. These feedback types were selected based on a review of the social science literature and a textual analysis of transcriptions of human-human verbal communication in a pilot collaboration experiment described in Section 5.3 to provide a balance of information to the person about what the robot is doing and what the robot expects them to do and are well-supported by the task model. This mimics the mix of implicit and explicit feedback is also used in human teams (Shah and Breazeal, 2010). Additional types of communication, such as providing salient information by highlighting important observations, which could occur at varying levels of granularity, as well as qualitative commentary about aspects of the task environment were also used by people but are not supported by the communication model.

4.3.1 Self-Narrative Feedback

The simplest form of feedback to produce is *self-narrative feedback*, where the robot tells the user what it is going to do, e.g., "I'll take care of X". This type of feedback is considered implicit in that it does not require any action or acknowledgment from the user but does imply that the next action should match well with the action that the robot has narrated. Other work in human factors has demonstrated that people use this type of feedback during coordination (Shah and Breazeal, 2010). It is also useful as it provides information on what the robot is going to do via verbal channels allowing the user to select actions they believe will work well with what the robot is about to do next. This is particularly helpful in scenarios with non-humanoid robot embodiments, as these typically give non-human-like cues that can make the robot's next action less

intelligible to a user. One challenge with employing self-narrative feedback lies in longterm interactions where repeated narration may become boring or annoying to a user. Another related concern is providing users with different desired levels of verbosity with the correct amount of self-narrative feedback. In longer interactions the amount of self-narrative feedback can be scaled back, or if the robot has learned the task reward structure, it can selectively issue self-narrative feedback in situations where the user's misinterpretation of the robot would lead to a greater reduced reward.

A communicative intent for each policy is available to the robot and defined as follows:

$$C_{narrative} = RobotDoing\langle r_i \rangle, \forall r_i \in \Pi$$

to select a single active communicative intent based on the role the robot is executing. To convey these communicative intents to the person, the robot has two options. It could say a phrase indicating its role assignment, i.e., telling a user what role it will execute during the task, or it could implicitly communicate its role assignment by providing feedback about which actions it is executing at a given point in time, leaving the person to infer the robot's role (policy) from this sequence of action selections. The former case requires a mapping of roles to a set of verbal feedback phrases, V, that refer to the correct role, $f_{narrate-role} : \prod_{robot} \to V$, while the latter case requires a mapping of every action to one or more verbal feedback phrases $f_{narrate-actions} : A_{task} \to V$. Ideally, V will contain multiple options for each policy or action so that the robot does not use the same redundant phrasing, which could lead to it being perceived as less intelligent, an issue to be avoided in HRI.

The purpose of this type of self-narrative implicit feedback is to make the robot's policy clear to the user and by extension its action selection mechanism given the current state, assuming the user is able to understand the robot's communication. This provides implicit guidance to the user's action selection by making the robot's course of action



Figure 4.3: An example of the effect of implicit feedback on model-based action selection. For simplicity, the user and robots actions are combined in the transitions. If the robot clearly conveys that it will take action a_{r2} then the user can safely select action a_{h2} without fear of reaching S_3 .

clear. In Figure 4.3.1, a scenario is depicted where the robot providing implicit feedback can alter the user's action selection. In scenarios, where the robot has learned the task and has value information for each state these scenarios can be detected automatically and used to generate a self-narrative communicative intent when the decision made by the user would change the resulting state value given combined robot and human action. In a single-agent model this can also be simulated by considering all possible sequences of user and robot actions and transition probabilities if available. Besides influencing a user's action selection periodic implicit feedback places less burden on a user than an explicit request since it is informational in nature and also serves to improve situational awareness of the user.

4.3.2 Role-allocative Feedback

Role-allocative feedback is an explicit request by the robot to the user to play a particular role during the task, i.e., that the person should change their behavior to a specific, desired policy that better coordinates with the robot's planned course of action. If the user is already performing the role the robot allocates this serves to reinforce the user's selected course of action, whereas if they have selected a different role this serves to guide them to a pattern of behavior that is better coordinated with the robot's selected policy. As an explicit form of communication it places an additional expectation that the user do what the robot says. In experimentation it has been found that people in the synthetic tasks are very likely to follow the robot's instructions but many possible confounds exist including the effects of novelty of the robot, the familiarity of the user with the task, and the importance or significance of the task beyond the experimental setting.

As an explicit request users have the option of complying or attempting to do a role other than the one the robot has requested. This leads to the need for adjustable autonomy in which the robot can both communicate and adjust its task performance in response to the user and vice versa. This type of scenario may also be amenable to extended dialogue with the robot in which the person and robot deliberate to determine who should do what aspects of the task. In this dissertation it is assumed that the robot's task policy is static over time and instead focus on generating role-allocative feedback that is sensible if complied with by the user.

To generate role-allocative feedback for the user, such as "you take care of X", the expected next action of the user given by our recognized policy π_{user} is compared to the robot's next action as given by the robot's policy, π_{robot} . This comparison is dependent on the type of model employed. A heuristic approach can be employed when there is no specific information about joint action allocation i.e., pairs of optimal action for each

state. The heuristic operates on a single-agent task model and is based on the principle of avoiding conflicting action. This approach assumes a list of discrete actions in A_{task} that would conflict in some way if executed by the robot and person at the same time. Given the policies of each agent, a check is made to see if they are expected to execute the same action given the current state. If this action is in the set of conflicting actions, an alternate role from Π_{roles} is selected for the user to do that is not conflicting, as detailed in Algorithm 2.

This heuristic approach can be improved upon if the robot has a more complete model of the task with access to a reward signal. The role-allocation suggestions for the robot should only make suggestions to the user when they are executing a policy that is clearly suboptimal and it should suggest actions that progress the task further, instead of simply avoiding actions that are in conflict. To formulate intelligent user role suggestions, the communication planner must have information about the best action for the person to take a given point in time. A more formal approach can be used to compute role-allocation suggestions, by employing a technique such as q-learning on the task model and joint robot and user action set. This yields a set of q values defined over $S \times [A_{human}, A_{robot}]$ where we can compare the user's probable next action based on the recognized user policy and select an alternate policy that matches the learned optimal action for the current state. If the user's recognized role matches the intended target allocation, a reinforcing action can be optionally issued indicating they have selected the right action according to the robot's model of the task. This assumes that the user and robot actions are both defined in the same state space and that both the user and the robot are taking action at similar time-scales.

The approaches described above result in a communicative intent that assigns the user to a given role:

 $C_{role-allocative} = PersonRole\langle r_i \rangle, \forall r_i \in \Pi_{roles}$

To convey this role to the user the system must ultimately generate some verbal feedback that conveys the requested allocation to the user in terms they can understand. The phrases for role-allocative feedback have similar requirements to those in the narrative feedback case. Specifically, we require a function mapping roles (or alternatively actions) to a set of verbal feedback phrases that the robot can say $f_{role-allocative} : \Pi \to V$. It should be noted that, at this step, the robot could proactively perform the better action instead of recommending that the user do it. For our initial experiments, we do not consider this option, but assume the robot's policy is stationary.

Algorithm	4.2	Alternative	role-	allocation	suggestions	

Require: Person and robot policies π_{user}, π_{robot} , current state s	\in
$S_{task}, ConflictingAction(a_{user}, a_{robot})$	
$a_{user} \leftarrow \pi_{user}(s)$	
$a_{robot} \leftarrow \pi_{robot}(s)$	
if $a_{user} = a_{robot}$ and ConflictingAction (a_{robot}, a_{user}) then	
for all roles $r = \{\pi_1 \pi_n\}$ do	
alternatives $\leftarrow \{\}$	
if $(r(s) \neq a_{robot})$ or	
not ConflictingAction $(r(s), a_{robot})$ then	
append r to alternatives	
end if	
end for	
return alternatives	
else	
$\mathbf{return} \{\}$	
end if	

4.3.3 Empathetic Feedback

Empathetic feedback, in which the robot makes a positive or negative statement in response to favorable or unfavorable changes in the task, is another form of implicit communication that serves to alert the user to the robot's interpretation of the current situation. People employ this form of communication to build rapport and situational awareness by alerting others to potential problems or other outcomes that might change further task planning. In the human-robot case it serves a similar purpose and can also be thought of as a reward signal issued by the robot, letting the user know the robot's assessment of the team's performance.

One approach to generating empathetic feedback, employed in the work is to monitor for specific state transitions in the task model that typically trigger people to issue an empathetic display. When such a transition occurs, we generate an appropriate communicative intent in the form of either a positive or negative expression (e.g., "Oh no!" or "Great!"). This approach maps specific transitions to communicative intents, e.g., $f(s, s_{prev}) = c \in C_{empathetic}$. By mapping specific state transitions to instance of empathetic feedback, these may also include more specific expressions of empathy, such as a description of what exactly went right or wrong. This more specific feedback, may affect users' subjective opinions of the intelligence of the robot as it identifies not only that a fault was detected but also what went wrong, leading people to conclude the robot has a higher-level understanding of the task.

Another approach that can be applied in certain scenarios where the robot is actively learning the task, and consequently has access to a reward signal, is to monitor the execution of the task and issue empathetic feedback when a fault or large change occurs. Existing work on monitoring execution and fault detection in these types of models exists (Pettersson, 2005). It is important that the robot match people's policies for producing empathetic feedback, requiring an appropriate threshold for triggering a positive or negative empathetic display. This can be calibrated by observing people perform the task and annotating such displays. Further investigation beyond the scope of this dissertation is needed to determine the efficacy of this approach and the impact of human-like empathetic displays in coordinating task collaborations.

4.4 Communication Executive

After a communicative intent is planned, the robot must then execute a social communication gesture using a combination of speech and embodied gesture to convey the intent to the user reliably. To generate an audio file, either pre-recorded or via a text-to-speech system, the exact content of the robot's verbalization that should be spoken aloud to the person is needed. To avoid repetitive use of the same phrase, which could lead to the robot being perceived as less intelligent (Torrey et al., 2013), a set of equivalent phrases is required for each action that can be played at random and ensure some variability in the robot's speech usage. If the robot has an expressive embodiment these speech actions can be combined with co-verbal gesture to make the robot easier to understand. Production of nonverbal behavior is a challenging problem that has been studied extensively in the embodied conversational agents community (Lee and Marsella, 2006). Existing robots are not currently capable of producing humanlike nonverbal behavior as they may have a reduced number of degrees of freedom, especially with respect to facial expression, or non-humanoid form factors. In Section 3.1 the approach for translating an abstract communicative intent into embodied speech is presented, followed by analysis and a study of particular gesture types and their use in human-robot collaboration.

4.5 Extensions of the Approach

In this section relaxations of the assumptions of the approach are assessed and applicability in different task domains is discussed. The full extension of the approach to these scenarios is beyond the scope of the dissertation.

- Observability: In partially-observable environments the robot can use of a POMDP, for the task model resulting in several changes to the approach. In that case, we would instead maintain a distribution over states and perform the requisite belief updates. All steps involving the directly observed state must be extended to use the belief state or take the most likely underlying state.
- Multiagent task models: If we are fortunate to have a multiagent MDP formulation of the task, the approach algorithm can be modified to look up the requested action for the user in the multiagent model. Given this action a communicative intent can be produced using the action-conveying method mentioned in Section 4.3 or by producing a communicative intent for allocating the user to a role in Π_{roles} that agrees with the joint policy, see Algorithm 4.3. The resulting communicative intent can be handled identically to the role-allocative intents in Section 4.3. If the user is already expected to take the correct action, the system can optionally provide some positive reinforcement.

Algorithm 4.3 Role-allocative feedback - single-agent MDP - non-conflicting action

Require: Person and robot policies $\pi_{person}, \pi_{robot}$, current state $s \in S_{task}$, joint action	on
$\langle a_{robot}, a_{person} \rangle$	
if $a_{person} = \pi_{person}(s)$ then	
return positive feedback	
else	
return role r in Π_{roles} where $r(s) = a_{person}$	
end if	

• Unmodeled and unequal human-robot task capability: The first concern here is that the human may have capability that we did not incorporate into the model, i.e., they can do something that the robot did not know about. This case is not extremely problematic as we assume the robot will still take the right actions based on the resulting state or observations after the user complete the unmodeled action. Since we do not include the action in the task model the robot will never ask the robot to perform it, which is not especially damaging to the interaction. The second possibility is that a particular person cannot perform one of the tasks that we assumed they could. In this case, if unmodified the approach might repeatedly attempt to allocate the action to the user, despite the person's inability to perform the task. In this case, it would be best to track user compliance over time and adapt the role-allocative communication of the robot to prevent this from happening.

The second concern with unequal human and robot capability occurs when the robot action set and the person's action set are disjoint. In this case our method for selecting non-conflicting actions fails as the robot and person can never take the same action. This case can be handled either by using a multi-agent model with joint policy to suggest the proper action to allocate to the person, or by modifying the non-conflicting action comparison to check a list of tuples containing all elements in the cross product of the robot action set and person's action set that could potentially conflict.

4.6 Summary

This chapter discussed the approach to generating coordinating social communication by generating three types of communicative intents including narrative feedback, roleallocative feedback and empathetic feedback. The requirements and a set of algorithms for producing these intents were discussed as well as design decisions in specifying the set of verbal phrases that the robot might use to convey each communicative intent. Finally, consideration was given to various extensions to the approach to environments with partial observability, different task structures, and areas in which the robot and person have unequal capabilities. In the next chapter, a pilot experiment informing the specific verbalizations employed by the robot is described as well as evaluations of the approach in two different task scenarios.

Chapter 5

Evaluation of Coordinating Social Communication

This chapter presents evaluations of the coordinating communication approach as instatiated on two pairwise human-robot collaborative systems and evaluated in two different task settings. An augmented reality task simulation environment developed to allow rapid-prototyping of different dynamic collaborative tasks using overhead projectors and a depth camera-based tracking system is described. A pilot experiment using this task simulation system and two user studies applying the coordinating communication planning approach to two different task scenarios are presented. The pilot study, informs the selection of roles and phrases used by the robot. The first user study evaluates the humanrobot task collaboration with and without communication with a convenience population, demonstrating an improvement in quantitative task performance. Finally, results from a cooking task with an older adult user population in an eldercare facility are presented.

In order to fully test a system for planning coordinating communication in the context of human-robot task scenarios, the first requirement is a robot and accompanying low-level control systems that can reliably perform tasks in an environment co-located

with a person. Because manipulation of real objects is a challenging problem and offthe-shelf software for performing these tasks is typically slow and requires tuning for specific environmental factors, using these behaviors to prototype real-world tasks is time-consuming and greatly reduces the range of possible scenarios that can be tested. Additionally, since the robot must monitor its teammate's activity over time, a collaborative task involving real object manipulation presents difficulties in balancing sensing and perception of the person and the objects the robot manipulates. There is relevant work in modeling, recognizing, and tracking human activity (Crick and Scassellati, 2008, Hamid et al., 2009, Kirsh, 1995) that can be applied in collaborative scenarios to offer more robust user action recognition. When applying this approach in real-world tasks, the constraints of the task, robot, and objects used influence the choice of robot system applied to the task. Although some robots are currently capable of performing some tasks reliably, making a robot reliably perform a task even seemingly simple manipulation is a difficult problem. To circumvent these restrictions and make the problem of creating a reliable collaborative robot that can reliably perform tasks, an augmented reality task simulation engine was developed. Using this simulation engine we developed a dynamic pseudo-herding task to be performed by human-human and human-robot teams and a simulated cooking game to be played between a person and robot on a large touchscreen. We applied the coordinating communication approach to these systems, allowing the robot to perform the task and issue coordinating communication. Evaluation was conducted with two target user populations through a series of user studies designed to evaluate the efficacy of the approach by robot's communication to improve team performance.



Figure 5.1: Diagram of the augmented reality task simulation environment with a person and Pioneer 2 robot.

5.1 Augmented Reality Task Environment

The goal of the augmented reality task environment was to develop a challenging, dynamic cooperative set of tasks involving multiple people and robots, thus necessitating the effective use of coordinating communication. The task simulator modeled the behavior of multiple virtual agents and virtual objects over time. The agents were projected onto the floor of the room via an overhead projection system (Figure 5.1). Simultaneously, the system made use of environmental sensing to track people and on-board sensing to localize robots in the shared space. By updating the simulated positions of agents and objects, robots and people in the room were able to interact with the virtual world through the same physical action, consisting of navigation around the room. The environment was calibrated by assuming a rigid transformation from the virtual space to the physical space of the room, allowing for generation of augmented sensor data, such as laser scans, from the point of view of a physical robot in the space. A calibration procedure for the projectors allowed for clicking four overlapping points and computing

a homography from the physical floor of the room to the projected view. Using OpenGL the projected output of each projector is warped, resulting in sub-pixel accurate alignment. Two Microsoft Kinects RGB-D sensors were mounted in the ceiling of the room pointing downward diagonally. These were extrinsically calibrated with respect to each other by detecting a standard calibration target in an overlapping region of the frame. The depth sensors were calibrated to the virtual floor of the room by placing physical markers at the corners of the projected area, allowing for accurate alignment. A point cloud processing pipeline was developed for tracking objects in the room since the built-in skeletal tracking did not work well for the selected mounting locations. To track moving objects such as people in the room the generated point cloud from each Kinect was downsampled and merged into a single cloud. Next, planes including the wall and the floor were removed from the image leaving only points at least some parameterized distance above the ground. These points were then clustered and the resulting clusters were tracked over time using an unscented Kalman filter assuming the number of objects to track is known a priori. The robot used in this scenario was the Pioneer 2AT as shown in Figure 5.2a. This platform was selected due to its competent navigation of the environment using an on-board Hokuyo laser rangefinder for localization and a pre-built map of the room. Additionally, the small size of the Pioneer was beneficial as it reduced the possibility of harm in an accidental collision with the person also moving around the environment.

This experimental setup had several advantages for conducting tasks compared with using physical objects and/or confederate experimenters. First, it allowed for many repeatable, dynamic agents that move autonomously in a directed or pseudo-random manner. It also allowed for calibrating the velocity, shape, and behavior dynamics of the simulated agents to make a given task easier or more difficult as desired. Finally, since it does not depend on physical objects, several different tasks could be conducted



Figure 5.2: A person and a Pioneer 2 AT robot collaborating on the pseudo-herding augmented reality task with the virtual elements shown projected on the floor of the room. In a) the lock is unlocked and only one sheep is in the pen (light blue) while in b) the robot has just finished the game by collecting the last sheep.

quickly by switching the task controller. The simulator abstracts away physical tasks that are not the focus of this work, such as object manipulation, allowing for the study of collaborative behavior. This is consistent with the notion of *research tasks* described by Martin et al. (1998) as a systematic abstraction of a real-world task. Despite this abstraction, the setup preserves some of the complexity of real-world environments, such as partial observability, noise, and occlusion by augmenting the robot's sensor data. This represents a reasonable trade-off of environmental realism for repeatable HRI experiments and expanded task possibilities.

5.2 Pseudo-herding Task

The first task scenario implemented in the augmented reality environment was a pseudoherding task in which a group of "sheep" appear from the boundary of the room over time with a specified arrival rate in agents/minute added to the scene. The sheep moved about exhibiting Brownian motion with simple avoidance rules and must be collected

and brought to a central holding pen. Only one sheep could be collected by an agent at a time, hence it is *pseudo*-herding. To collect a sheep the agent must move to a position overlapping the sheep. People's tracked locations were depicted by projecting a ring around the last tracked position. After a sheep is collected it follows the person until it is dragged inside the circular pen, at which point it drops off and stays inside. In addition to the pen there were two timed game-play elements that count down on a timer: a lock and a light. When the lock timer reached zero the pen unlocked and any sheep that were captured in the pen escape and began to wander the room again, effectively reverting any progress made to that point. When the light timer reached zero the visibility of the projection was modified to simulate the lights going off in the room. A stencil buffer was used to draw only a circular area around each agent with the rest of the projection blacked out. This did not affect the state of the sheep herding but prevented users from easily finding any more wandering sheep. Each of the timer elements were drawn with sprites indicating their current state and featured a progress bar and color change to clearly indicate when they were about to expire. The participants could reset the amount of time left on the timers by moving on top of the timer object of their choosing. The task was designed to be dynamic since the sheep move about unpredictably avoiding the users; it was divisible and loosely-coupled, since sheep could be collected by a single individual; and it supports different completion strategies such as allocating one person to maintain the timed elements or dividing the room into different spatial zones. Almost all elements of the simulation including the number, speed, arrival rate, etc. of the sheep and countdown interval of the timers were configurable. A separate mode designed to elicit time-induced stress in participants was developed in which the sheep have a life meter that counts down over time. When the life meter reaches zero for a specific sheep, that sheep would disappear making it no longer collectible and reducing the number of sheep the team can successfully collect. The stated goal of the game, as introduced to participants, was to collect all the sheep into the pen area as quickly as possible.

5.2.1 Application of the Approach in the Pseudo-Herding Task

Applying the approach to the pseudo-herding task required development of the task model, robot policies, person policies, and verbal phrases specific to the herding task. Since the environment is fully observable, the task is modeled as a Markov decision process $M_{herding} = (S_{herding}, A_{herding}, R, T)$. The state is comprised of a set of independent random variables $Sherding = \{X_{lock}, X_{light}, X_{sheep}\}$. For the herding task a state space consisting of three features is defined, $X_{lock} = \langle high, low, unlocked \rangle$, $X_{light} = \langle high, low, off \rangle$, and $X_{sheep} = \langle true, false \rangle$. Because navigation is used as a proxy for manipulation in the augmented reality task environment, navigation actions for the robot to go to each of the work items (lock, light, pen) were developed. Since all the sheep are equally valuable, a single navigation action is used for sheep retrieval in which the robot gets the sheep closest to itself. Thus the possible actions are $A_{herding} = (GotoLock, GotoLight, GotoPen, CollectSheep)$. The robot's task behavior during the task consists of a set of policies mapping states to action $\pi : s \to a$. Note that in the augmented reality environment the robot and person have the same capability and thus share the same set of actions, $A_{robot} = A_{person} = A_{herding}$.

To apply the communication approach, a set of policies that covers the set of roles that a team of people might use to complete the task, is specified, as noted in Section 4.2. Fitting with the assumption that roles are correlated with distinct work objects, roles associated with each work object as well as all possible pairwise combinations are used: $P = \{p_{lock}, p_{light}, p_{sheep}, p_{lock-light}, p_{lock-sheep}, p_{light-lock}, p_{light-sheep}\}.$ In the single object policies the action in every state is the navigation action associated with the object i.e., $\forall s \in S, p_{lock}(s) = GotoLock$. For the policies with multiple objects, action preferences are given by the following object priorities $o_{lock} > o_{light} > o_{sheep}$, thus the p_{lock_light} policy will make the robot go to the light if the lock is in the high state and will go to the lock in all other states. No three-object roles were defined since that would consist of a single agent trying to do all parts of the task. Given the equal capability of the robot and person for this particular task, this set of policies is used to both control the robot and to model the person's behavior. This mapping of roles to policies was informed by a pilot experiment in which two people performed the task together.

5.3 Study: Communication and Role-Usage in the Pseudo-Herding Task

To inform the definition of the roles and referential phrases used to assign them in the pseudo-herding task, a pilot study was conducted with two-agent teams collaborating on the pseudo-herding task in the augmented reality (AR) environment including allperson teams and human-robot teams. There were twelve total sheep in the herding task with speeds and other parameters set to make the task challenging but not too difficult. In addition to two-person teams four sessions were conducted with a person interacting with a Pioneer 2AT robot, that was controlled by Wizard of Oz (WoZ) (Figure 5.3a), in which an experimenter triggered actions from $A_{herding}$ by hand while observing a person do the task and listening to an audio recording of their speech from a headset microphone. The operator was instructed to trigger the actions to best match the needs and requests of the person performing the task. These sessions were designed to see how the participants instructed the robot with verbal commands and how the task was partitioned between the collaborators. To verify that the selected actions for the WoZ interface would be compatible with the desired level of instruction given by the users,



Figure 5.3: Example of the Wizard of Oz interface used by the experimenter to control the Pioneer 2 robot during the pilot study. The control interface in a) contains buttons to trigger each of the autonomous actions from the robot's action set and a cancel button to stop all motion, as well as text displays for completion status and error output. In b) the visualization of the robot, person, and simulated objects potions is depicted.



Figure 5.4: An overhead diagram of the experimental setup used when constructing the action set of the robot. Users were presented various diagrams and asked to instruct one of the participants what to do next.

a survey-based data collection was taken by administering surveys using an overhead diagrammatic rendering of the experimental setup (Figure 5.3).

The action set for the robot was developed based on the feedback of 10 participants viewing 10 slides with various task states presented, for example, a scenario with a lock unlocked and another with a locked lock but one sheep left to be collected. In the actual pilot experiment all users, except one, instructed the robot at a similar level of detail, by requesting that it perform the various object-related subtasks. One participant attempted to provide low-level navigation commands verbally by requesting the robot move forward a particular distance and then turn left or right. This particular user's transcripts were discarded as an outlier, since the robot was only able to perform higher-level sub-tasks autonomously.

5.3.1 Study Design

For the human-human experiments, participants were recruited in two-person teams. They were introduced to the task by standing outside the room and watching the experimenter perform each of the following task elements: herding a sheep, triggering the lock, and triggering the light. They were also shown a demonstration of the consequences of letting the lock timer reach zero. They were then told that the goal of the game is to herd all of the sheep as quickly as possible. After they entered the room it was verified that they were both being correctly tracked before they were given a countdown, and asked to complete the task. For each group the task was performed once with the stress condition, in which the sheep fade away and disappear after a set length of time, and once without, with the order counter-balanced.

For the person-robot experiments, participants were introduced as above and told that the robot could understand their speech but could not speak itself. Furthermore, they were told that the robot knew how to do the task. The robot was operated via Wizard of Oz (WoZ) by an experimenter listening to the participant's audio feed and triggering the autonomous robot actions, $A_{herding}$ via a simple button user interface.

Data collected included a record of the underlying task state, the tracked locations of the agents (people and robot, if applicable), an audio recording of each person from a wireless headset microphone, and video of the interaction from overhead and side-view cameras. From these data a series of post hoc metrics were computed including idle time, distance traveled, number of sheep collected, and number of triggers of the timer elements. The recorded audio from each person was also manually transcribed post hoc and quantitative measures of task performance were automatically annotated using the recorded task state.

5.3.2 Participants

The experiment was conducted with six two-person teams and five person-robot teams with a Pioneer 2AT robot. In total 17 undergraduate students were recruited from the University of Southern California to perform the task. All teams were able to understand and perform the task.

5.3.3 Hypotheses and Outcome Measures

The following hypotheses were formulated based on an understanding of the task and human-human collaboration:

- **H1:** In the human-human collaborations, role use will be correlated with the work objects (lock, light, pen, and sheep). This fits with the design goals for the modeling of roles as discussed in Section 4.2.
- **H2**: In the human-human collaborations, people will use other gestural cues, including eye gaze, to conduct the collaboration.
- H3: In the human-robot collaborations, people will instruct the robot to do the task by assigning roles correlated with the work objects.
- H4: The two person teams will have higher speech usage than the human-robot teams, as people collaborating with the robot will expect it to have a more limited verbal capability directly applicable to the task..

5.3.4 Results

In the pilot experiment, it was found that people engaged with the game with some running around the room and others strongly reacting to negative events, such as the sheep getting out. Confirming hypothesis **H1**, in the person-person condition people largely used speech to either convey information about role allocation, with phrases like "get the light" and "get the lock" and their variants as some of the most spoken phrases. Another common use of speech was to express positive or negative reactions to the team's progress through the task via empathetic responses, motivating our use of this type of feedback by the robot. People's use of self-referential feedback to convey their intended actions and their willingness to ask or allocate an activity to a partner in times of need are competencies that the approach seeks to duplicate via robot feedback. Contrary to hypothesis **H2**, gesture usage overall was very low with few participants choosing to point to objects, however gaze and attentional cues appear to be crucial to the task as people appear to monitor where others are looking and plan accordingly. Interestingly, all but one of the two-person teams experienced at least one incidence of the lock timer lapsing and having to restart the task.

In the human-robot case it was found that users spoke less overall and only used speech to allocate tasks to the robot, thus neglecting to provide self-narrative feedback about their own intentions and confirming hypothesis **H4**. One user attempted to teleoperate the robot verbally using phrases like "robot, go forward", but all others used similar intuitive phrases that corresponded well to the roles used in the human-human case confirming hypothesis **H3**. Some users assigned the robot to a static position where it triggered the lock at the beginning of the task and then proceeded to perform the rest of the task on their own, presumably due to a desire to minimize their interaction with a difficult-to-use robot as well as to remove the need to continually monitor and assign it tasks. Although this is a valid strategy to complete the task, it does not take full advantage of the robot's capability and resulted in a less equal partitioning of the task as noted by extracted task performance metrics and distance traveled.



Figure 5.5: Bar chart showing sheep herding allocation across agents for person-person and person-robot teams. Note the less equitable allocation of sheep herding in the human-robot teams compared to the human-human teams.

5.3.5 Discussion

When considering usage of role-based communication a common decomposition of the task was found across teams; although some roles described atomic actions involving work objects, such as triggering the lock or the light, others, such as herding a sheep, conflated a series of compound actions (finding a stray sheep, catching it, and bringing it back to the pen). This informal notion of roles by people performing the task validates our selection of a policy-based representation for flexible role definitions that support both atomic actions and sequences of actions. Most of the positive and negative qualitative feedback people employed was in reaction to one-time events in the course of the game as opposed to evaluations of their own or other's performance over time. Each two-person team experienced at least one event in which the timed lock lapsed and all sheep escaped the pen. Whether this is due to the learning curve of the task or breakdown in communication is not clear but in most cases it elicited negative comments without assigning blame, such as "Oh no! The sheep escaped".

In the human-robot condition there was room for improvement in the task in the form of significant robot idle time as well as instances of the lock lapsing and reseting task progress. An autonomous robot, even one with a static policy, that proactively takes action during the task, could yield better task performance and would be welcomed by users. The disparity in speech use between the human-human and human-robot case suggests that the robot could be made a better teammate and more interesting interaction partner by the use of coherent verbal feedback.

5.4 Study 2: Coordinating Social Communication in Pseudo-Herding Task

This study aims to evaluate the application of the coordinating communication system in the augmented reality environment on the pseudo-herding task with a convenience population recruited from on-campus sources.

5.4.1 Study Design

Although the method presented for generating coordinating verbal communication is designed to be applicable in a wide range of task settings, the following user study is designed to evaluate the approach in a co-located human-robot task collaboration with a convenience population as compared to a non-communicating robot. The approach was evaluated by measuring the time to complete the pseudo-herding task in the augmented reality environment as well as collecting subjective measures of the robot's performance via surveys. The communicating robot was implemented as described above and compared to a control, a non-communicating robot, i.e., a robot that performed the task in an identical manner without issuing any verbal feedback. This control was selected since few existing methods for high-level task planning integrate human feedback, making a task-only robot a likely scenario.

A within-subjects experiment was conducted, in which each participant performed the task twice, once with the communicating robot, and once with the noncommunicating one, with order of presentation counterbalanced across subjects. A within-subjects design was selected since the overall interaction is quite short (usually about 2-5 minutes) and our pilot experiments found the times to finish the task varied a great deal among people. Also, since the task is unfamiliar to all users due to its unique setup, there is the potential for each person to learn the nuances of the task over time, and improve performance. The robot's task behavior was identical for each condition with the on-board speakers either enabled or disabled for the talking and silent robots, respectively.

5.4.2 Task Description and Setup

The task setup for this experiment was the pseudo-herding task with 12 sheep and speeds and timings set the same as the pilot to make the task challenging but not too difficult. The static robot policy used for the task was the light and sheep favoring policy, resulting in the robot alternating between turning on the light and herding sheep (collecting them and returning them to the pen area).

Participants watched an experimenter explain and perform each of the task elements, including: herding a sheep and triggering the lock and light. Participants were told that the goal of the game was to herd all of the sheep as quickly as possible. The random seeds used to generate the sheep behavior were initialized identically for each participant to keep the simulated sheep behavior as consistent as possible between runs. The task simulator was configured with 12 sheep. The speeds of the sheep and timing of the lock and light elements were tested in a pilot experiment. The task thus calibrated was possible but somewhat challenging for a single person to perform alone. The static robot policy used for the task was favored light- and sheep-collection, resulting in the robot's first priority being turning on the light and collecting and returning sheep to the pen area while the light is on. The only difference in robot behavior in each condition was use of communication, which was disabled without an effect on task planning on the non-communicating robot.

5.4.3 Participants

A total of sixteen participants (5 female, 11 male) were recruited from campus sources. Participants' ages ranged from 17-29, with a mean of 21. Most had completed at least some undergraduate college education. One participant was an outlier, not able to complete the task due to a misunderstanding of the rules; this participant's data were discarded. The order of presentation of the conditions was pseudo-randomized with half seeing the communicating robot first and the other half seeing the non-communicating robot first.

5.4.4 Hypotheses and Outcome Measures

The following data were collected: a record of the underlying task simulator state, the tracked locations of the agents (people, robot, sheep and virtual objects), an audio recording of each person from a wireless headset microphone, and video of the interaction from overhead and side-view cameras. After completion of the task participants were asked to take a survey asking demographic questions as well as a series of 27 questions about the robot as a teammate, rated on a 7-point Likert scale from "strongly agree" to "strongly disagree". The following hypotheses were formulated based on our understanding of the task and human-human collaboration:

- **H1**: Participants will prefer the communicating robot over the noncommunicating robot.
- **H2**: The talking robot will be rated as a better teammate by participants despite identical task control.
- H3: The communicating robot will decrease the incidence of conflicting actions taken and also increase the situational awareness of the person, leading to improved task performance.
- **H4**: The communicating robot will be able to make its static role and assessment of the task clear through the use of verbal feedback.

5.4.5 Results

For the objective results, the time to complete the task to the nearest second was recorded for each trial by watching the recorded video, marking the start of the task as signaled by the experimenter, and the end as noted by the freezing of the projected display upon retrieval of the final sheep. Trials with the communicating robot were completed on average 17 seconds faster (M = 100.5 seconds, SD = 21 seconds) than the ones with the non- communicating robot (M = 117.6 seconds, SD = 40 seconds), p < 0.03. As our conditions were repeated within subjects, we then performed a post hoc analysis by fitting a linear mixed effects regression and comparing to a baseline single mean model. In comparing the two models, we find incorporating the communication factor has a significant impact on the duration (p < 0.04), confirming Hypothesis 2.

In survey results participants reported their preference for the communicating robot.



Figure 5.6: Time to completion and allocation counts of agent actions in communicating and silent conditions.

Question	0 (strongly disagree) - 6 (strongly agree)		
	М	SD	
The things the robot said made sense.	5.27	1.6	
The talking robot was a better teammate than the silent robot.	5	1.4	
The task was more fun with the robot than if I had done it alone.	5.4	0.8	
The talking robot was more fun.	5.3	1.2	
The robot's talking helped me understand what it was going to do next.	4.9	1.9	
The things the robot said helped me de- cide what to do.	3.8	2.0	
I tried to do what the robot told me to do.	3.6	2.4	

Table 5.1: Post-experiment survey results

Table 5.2: Survey results on the pseudo-herding evaluation comparing communicating and non-communicating robots. Scale is 0 for strongly disagree, 6 for strongly agree.

In a free response question participants were asked to describe the strategy of the robot. These responses were graded for accuracy where it was found that all but one user correctly named the two aspects of the task that were performed by the robot.

5.4.6 Discussion

In the pseudo-herding task, in the quantifiable terms of state value, the most costly mistake that can be made is to let the timed lock lapse, releasing all collected sheep and forcing the team to start over. We intentionally selected a static robot policy in which the lock was ignored by the robot, requiring the participant to manage it instead. Across all 30 trials, the lock was allowed to lapse in 2 of 15 trials in each condition. Because these costly errors were roughly evenly distributed across conditions and were relatively rare with (M = 0.4, SD = 1.1), it is unlikely that the objective performance improvements are due to an avoidance of mistakes; it is more likely due to the clarification of responsibilities provided by the verbal feedback and minimizing the user's context switching between lock and light (Figure 4).

Distinguishing the impact of each type of feedback is challenging as they are issued while the robot and the person are moving about the room, making subtle effects difficult to discern from planned, task-directed actions. The most common feedback issued was self-narrative as the robot narrated its actions whenever they differed from what it was doing previously. Since the role-allocative feedback is very specific, with post hoc analysis we can compare the robot's suggested action with the action ultimately undertaken by the person. Over the course of the 15 trials in which the robot used verbal communication, and therefore was capable of assigning roles, the robot made a total of 12 role-allocation suggestions to users. To assess users' adherence to the robot's requests, we assessed whether a matching action for the participant was recognized within a short time interval. In total, participants adhered to all 12 requests within 10 seconds and often much faster. This adherence rate is likely due to the low cost of completing the robot's requests as well as the novelty of the interaction. Still, anecdotally, there are multiple instances of participants abandoning a prior plan and altering course in response to the robot's suggestion.

The low overall number of role-allocation suggestions is likely due to most participants allowing the robot to take full responsibility for the light. Also, since the robot attempts to track the person's active policy using a stream of recognized actions (currently from the task simulation), it only accounts for completed actions when determining the user's most likely policy. This has the effect of ignoring additional information provided by the person's trajectory, i.e., motion in the direction of a given object, and relies on the user reaching a steady state of task performance where they perform the same types of actions repeatedly. Anecdotally the refresh rate of the various communication models often appears to result in a narrative phrase and role-allocative phrase happening immediately after one another, such as "I'll go get the light. Can you take care of the lock?", often resulting in the user expressing agreement. The survey data reveal that all but two people agreed that the talking robot was a better teammate (see Table 1 for a summary). The communicating robot had similarly high scores for being more fun and for its verbal feedback making sense during the interaction. Questions asking how participants adjusted their behavior in response to the communicating robot are less consistently agreed with, perhaps due to people taking credit for their own actions and the team's success or failure.

5.5 Study 3: Coordinating Social Communication with Older Adults

The communication system was also evaluated outside of the augmented reality environment on a simulated cooking task with an elderly user population using the Bandit upper-torso humanoid robot mounted. This study was designed to demonstrate the application and efficacy of the coordinating communication approach in a task involving guiding a target user population to desired outcomes and evaluating the effect of implicit compared to explicit communication on user compliance and subjective outcome measures.

5.6 Simulated Cooking Task

In addition to the validation of the task in the augmented reality environment, the approach was applied to a second task developed to demonstrate the applicability of the approach in another task setting. The second task consisted of a simulated cooking task, in which users collaborated with an upper-torso humanoid robot (Bandit II). The cooking task is conducted entirely on a touch screen interface (Figure 5.6) and requires users to fill a set of orders using a set of food items. Each order consists of combinations of six food types. Users must fulfill these orders by moving ingredients through two stages of cooking, preparation and cooking, before the item is placed on a matching picture of the item in the order list. The task was conducted with an elderly user population at an off-site eldercare facility.

The cooking task was created to be an engaging activity that is thematically-relevant to the population, featuring element of activities of daily living i.e., food preparation. The robot also has the capability to prepare food items, making the task cognitively engaging as it requires selecting food items that complement the items the robot is

Cooking with Bandit				
Order #1	Order #2	Order #3	Order #4	Order #5
		× b		*
Cook				Trash Item
Preparation				Trash Item
Ingredients				

Figure 5.7: Screen capture of the interface used for the virtual cooking task. Users are asked to fill all the orders in the top row by moving food items from the green ingredients box through two stages (preparation and cooking) until all items are fulfilled as indicated by green check marks.

making. The task was piloted with users at the facility and modified to ensure that the users were able to do the task using the interface and did not find it too difficult.

The task simulator was implemented using the Robot Operating System (ROS) as a state machine monitoring the progress of each food item through each step as well as the completed items requested as orders. The interface was developed as a web-based interface connected to ROS with roslibjs, which dispatched changes to the task state requested by users by tapping elements of the user interface. The robot control system separately monitors the state of unfulfilled items remaining on each order and the user's food preparation actions, which are given as input into the communication approach, allowing the robot to issue verbal feedback during the task.

5.6.1 Application of the Approach in the Cooking Task

Applying the approach to the cooking task required development of the task model, robot policies, person policies, and verbal phrases specific to the herding task. As the task state is simulated and fully-observable, the task is modeled as a Markov decision process $M_{cooking} = (S_{cooking}, A_{cooking}, R, T)$. The state is comprised of a set of independent random variables representing the level of need for each food item $Scooking = \{X_{steak}, X_{chicken}, X_{salad}..\}$. Each of these features is defined, $X_{food} =$ $\langle high, med, low \rangle$. Because the task is conducted using the touch screen interface, the robot's performance of the task is completed virtually, with the control system simulating the robot's ability to transition the food items through each stage with an appropriate amount of waiting between each step to make the robot's performance similar to a person. The actions available to the robot consist of preparing each food item $A_{cooking} = (MakeSteak, MakeChicken, MakeSalad, ...)$. The robot's task behavior during the task consists of a set of policies mapping states to action $\pi : s \to a$ that consist of the robot monitoring and preparing combinations of three types of food items. Note that in the augmented reality environment the robot and person have the same capability and thus share the same set of actions, $A_{robot} = A_{person} = A_{cooking}$.

To apply the communication approach to this task, a set of policies that covers the set of roles that people might use to complete the task, is specified, as noted in Section 4.2. Fitting with the assumption that roles are correlated with distinct work objects, roles were created for the robot to monitor and fulfill combinations of the six selected food items. A static policy, in which the robot monitored three of the six food items and fulfilled all requests for those food types from left to right. This specialization based on food types was motivated by piloting task performance with four participants recruited from a convenience population performing the task. Other types of specialization could include partitioning orders spatially, with one person taking care of the left half and the other doing the right half. Since the aim of the approach is to generate useful feedback to guide user behavior, rather than alter the robot's behavior to best match the person's completion strategy, a specific and static completion policy was selected, requiring the robot to use social communication strategies to inform the user of its action choices and guide them to complementary roles. Given the equal capability of the robot and person for this particular task, this set of policies is used to both control the robot and to model the person's behavior. This mapping of roles to policies was informed by a pilot experiment in which two people performed the task together.

Phrases for the task were collected via a data collection with four participants from a convenience population. Participants were shown the user interface depicting a set of desired orders and a description of the food type-based specialization, e.g., that they were responsible for each of three selected food types or they were telling a partner to specialize similarly. Phrases for conveying responsibility for each type of food item were chosen and decomposed into the part of the phrase indicating responsibility such as "I'll make the steaks" and the specific phrase used to describe each food item. From these


Figure 5.8: In (a) an overhead diagram of the experimental setup for the cooking task is shown. The participant (depicted in red) is seated at a table in front of the touch screen. The robot (grey) is placed across the table at an angle. In (b) The experimental setup during the instruction phase, where participants were taught in a group how to perform the task to ensure a common learning of the task. The robot is seated diagonally across from the robot.

two components, phrases were generated using the NeoSpeech text-to-speech engine allowing the robot to convey responsibility for each of the food types.

5.6.2 Task Description and Setup

The cooking task setup consisted of the participant seated at a table in front of a large 27-inch touch screen computer with the upper-torso humanoid Bandit robot seated diagonally across from them (Figure 5.8a). The large screen was selected for use with the elderly user population to alleviate potential eyesight issues as much as possible. The user was introduced and familiarized with the task by the experimenter. The user was introduced to the task by the experimenter, first in a group session with the other participants, to ensure everyone had the same understanding of the task, and then once again before their interaction with the robot. Users were shown how to use the interface to complete the task. The participants were told that the robot would assist them in completing the task, although the progress of the robot is not depicted on the interface.

The robot displayed an idle behavior during the interaction in which it moved its head direction to saccade between the person, the screen and locations behind the screen. It also opened and closed its mouth when issuing speech, with the mouth opening at the beginning of a phrase and closing afterword to return to a smile. This was done as the mouth motors are noisy and not fast enough to mimic humanlike mouth movements during speech. A short introductory speech was issued before starting, the task as well as an acknowledgment of finishing the task.

5.6.3 Study Design

The user study was designed to evaluate the effectiveness of the coordinating communication system at enabling the team to successfully and efficiently perform the task, evaluate the use of the approach with a target population for assistive service robots, and to identify differences in performance and user preference when using implicit and explicit communication modalities. The pseudo-herding task evaluated the combination of all three forms of feedback as compared to a silent robot. As a complement to those results, this study was designed to identify relative differences in the usage of the two most common feedback types: self-narrative and role-allocative types. The study design was within-subjects with two conditions, one with the robot only issuing implicit communication, i.e., self-narrative feedback and one with the robot only issuing explicit communication, i.e., role-allocative feedback.

5.6.4 Participants

This experimental validation was conducted with three female participants recruited through a partnership with be.group, an organization of senior living communities in the Los Angeles area. Participants were recruited by a flyer distributed at the facility. Six participants responded to the flyer with three participants ultimately participating; the others were absent due to unrelated health concerns. Each participant was shown both the robot issuing only role-allocative communication (explicit feedback) and the robot only issuing self-narrative communication (implicit feedback). The order of these conditions was counterbalanced by alternating the order of presentation between participants.

5.6.5 Hypotheses and Outcome Measures

H1: Participants will prefer the robot that uses implicit communication over the robot that uses explicit communication.

H2: Team performance will be improved when the robot uses implicit feedback as compared to explicit feedback.

The following data were collected: a record of the underlying simulator state, the selections users made using the interface, and video of the interaction from a side-view camera. A survey to evaluate the user's opinion of the robot as a teammate, rated on a Likert scale, was administered between performances of the task and at the end of the experiment. An open-ended interview was also conducted after the conclusion of the experiment.

5.6.6 Discussion

Due to the small number of participants, and large differences in task completion time, the objective measures collected during the task showed no statistically significant differences between the explicit and implicit conditions, offering no support for H2 and suggesting a larger sample size is needed. Task completion time was highly-variable between participants as different users showed differing levels of familiarity with touch screen interfaces and ability to tap precise points on the screen. In the survey, all users expressed a preference for the robot that used implicit communication and described the robot that used explicit communication to be "bossy" and "demanding". This evidence supports H1 and is aligned with the results from the survey results from the pseudo-herding experiment, in which users expressed lower agreement with phrases such as, "the robot helped me decide what to do". It should be noted that the implicit feedback, which provides information about what sub-tasks the robot is performing, requires the user to understand what the robot has said and then select complementary actions, while the explicit feedback only requires that users follow the robot's instructions. Although both types of feedback, as produced by the approach, were sufficient for successful completion of the task the use of one or other would appear to have an impact on user preference. This pattern is also demonstrated in prior work in human factors that demonstrates that people generally tend to use implicit modes of communication during task collaboration (Shah and Breazeal, 2010). While further research is necessary, the balance of implicit and explicit feedback used by a robot appears to have implications for the perception of the robot as a collaborator and may require only the situational application of explicit, role-allocative feedback to guide user behavior. In the pseudo-herding evaluation, the role-allocative feedback is issued only when a potential conflict is detected, which reduces the relative proportion of feedback in that evaluation compared to the cooking task.

5.7 Summary

In this chapter we discussed the development of an augmented reality task simulation environment that allows for human-robot teams to collaborate on a task involving simulated agents and objects. The implementation of a pseudo-herding task in this environment as well as the application of the communication approach to this task was discussed. A pilot study to inform the definition of roles and verbal feedback used by people in this task setting was conducted with human-human and human-robot teams, with the robot controlled by Wizard of Oz. An evaluation of the approach in the pseudoherding task was conducted in a within subjects experiment demonstrating objective improvement in team performance and subjective improvement in user satisfaction with the robot collaborator as compared to a silent robot with identical task controller. A second evaluation of the approach with an elderly an older adult user population was conducted at an eldercare facility, studying the effect of implicit and explicit role-based communication on team performance and user opinion.

Chapter 6

Summary and Conclusions

This dissertation addresses a number of research problems related to generating appropriate coordinating social communication from an autonomous robot to a person during the course of a pairwise task collaboration with the aim of allowing the robot to provide helpful and guiding feedback to a human teammate. An approach to planning coordinating social communication based on a formalism of the human notion of role allocation is presented as well as an analysis of the limitations and applicability of the system to a variety of task scenarios. The implementation details of the communication approach as applied to two sample task environments were presented. User studies were conducted in each of these task environments with two different user populations for the purposes of evaluating the efficacy and generality of the approach.

6.1 Major Contributions, Findings, and Insights

Enabling an autonomous robot to produce relevant communicative social behavior during a pairwise task collaboration poses many challenges. The communication produced by the robot must be relevant and contextually accurate given the current state of the collaborative task and it must convey this information in a way that the person is able to understand and use to improve their situational awareness. The following section summarizes the main contributions of this dissertation.

- 1. A novel representation for modeling user roles in coordinated activities capable of capturing the action selection preferences that people typically exhibit when allocating responsibilities during joint activity and compatible with traditional planning methods.
- 2. An approach for planning coordinating social communication actions in a situated pairwise task collaboration between a robot and a person aimed at providing useful feedback to the user and enabling the robot to actively participate in the joint planning. The planning approach incorporates activity recognition, planning under uncertainty, and principles from social science of human joint activity to produce narrative, role-allocative, and empathetic feedback in real-time during dynamic task activity.
- 3. A framework for executing social coordinating communication during robot task performance, making use of the given social communicative capabilities of the robot in use and potentially including speech and embodied gesture. The robot's communication actions are derived from human-human communication in similar task situations and applies relevant design principles from the social sciences on joint attention and politeness theory to ensure that the robot provides coherent feedback.

The following are secondary contributions:

1. A robot communication executive that produces coordinated speech and gesture during a human-robot task collaboration and also supports speech interrupts with stop words for state changes.

- 2. User studies demonstrating the ability of the approach to improve both team performance and subjective measures of the robot as a teammate, as compared to a silent robot with an identical task control system.
- 3. A user study of the perceptual accuracy of robot head and arm deictic gesture production and control system for reducing grounding errors accounting for the noise introduced due to interpersonal differences in perception.
- 4. An augmented reality task simulator created using overhead projectors and depth sensors (Microsoft Kinect (Microsoft, 2010)) that allows for rapid-prototyping of different dynamic task scenarios supporting any combination of people, robots, virtual agents, and virtual objects.

The results of the user studies as presented in Chapter 5 demonstrated improvements in team task performance with a robot using communicative feedback as compared to a silent robot despite identical task control systems. This demonstrates the potential for effective mediation of human-robot task collaborations through humanlike social communication modalities. A separate evaluation of the approach with an older adult user population demonstrated that the balance of implicit and explicit communication modalities is an important consideration and has an effect on users' subjective evaluations of the robot as a teammate.

6.2 Open Problems and Future Work

A selection of open problems and potential future extensions to the work presented in this dissertation and applicable relaxations of the human-robot collaboration assumptions made in the approach described in this work.

6.2.1 Extension to multi-robot and multi-person teams

In the pairwise scenario addressed in this work, all social communication is produced by one party and consumed by the other. In a scenario with more than two agents there is the possibility that one of the agents misses a communication that other agents received. Also, if addressing a group it is necessary to clarify who one is addressing if providing an explicit instruction but not necessary if providing an informational prompt or implicit cue. Managing these complexities and also dealing with networked communication between robots that the user cannot overhear makes this an area with many potential research questions that could be applicable in scenarios where one or more users interact with multiple robots.

6.2.2 User role specification and generalization across tasks

The policies used to represent roles in the evaluations presented in this dissertation have so far been hand-designed based on data from two people collaborating on a task, as they are domain-dependent. The limitations of this representation are that roles must be specified as domain knowledge when applying the approach to a new task. This is often straightforward if the task is separated by two person teams in a standard manner. In the tasks presented in this work the partitions of different actions into roles often involve separation of the task into common work areas or work objects. A method for extracting roles from observational data or teaching a robot the roles present in task through a guided interaction remains unexplored.

6.2.3 Communication personalization and adaptation

There is some evidence that users will change their communication tendencies when under time-induced stress (Shah and Breazeal, 2010). Users with different personalities may also prefer different amounts of communication from the robot and as the user gets increasingly familiar with the task at hand, it is likely they will require less verbal guidance, which could ultimately get annoying under repeated performance of the task. Adapting the amount, type, and execution of the robot's behavior over time in response to these factors remains an open but important problem that could lead also lead to improved human-robot task collaborations. The use of methods from dialogue management is also potentially applicable as this would enable bi-directional communication from the user to the robot and vice-versa, a scenario that this work does not currently address.

There is a growing body of work in the area of enabling robots to be better collaborators and partners when working with people. The work in this dissertation addresses a range of scenarios featuring a person and robot interacting in a co-located environment to achieve a shared goal. As robots become increasingly pervasive in home and work environments, they will require a means of coordinating their behavior with people of various backgrounds and levels of familiarity with robotics systems. The results of the user studies presented in this work also demonstrate that providing communicative feedback to people interacting with a robot can have significant effects on user behavior during a structured activity even when all else is held constant. The capability for a robot to interact and assist users via humanlike communication modalities is a promising field of future inquiry as it opens up the possibility for robots to reinforce, assist, and guide user behavior in complex task environments, providing beneficial improvements in a user's situational awareness, engagement and ability to safely coordinate behavior with autonomous robot teammates.

Bibliography

- A. Bangerter. Using pointing and describing to achieve joint focus of attention in dialogue. *Psychological Science*, 15(6):415, 2004.
- A. Bangerter and D. M. Oppenheimer. Accuracy in detecting referents of pointing gestures unaccompanied by language. *Gesture*, 6(1):85–102, 2006.
- S. Baron-Cohen. The evolution of a theory of mind. The descent of mind: Psychological perspectives on hominid evolution, pages 261–277, 1999.
- J. M. Beer, C. Smarr, T. L. Chen, A. Prakash, T. L. Mitzner, C. C. Kemp, and W. A. Rogers. The domesticated robot: design guidelines for assisting older adults to age in place. In *Human-Robot Interaction (HRI), 2012 7th ACM/IEEE International Conference on*, pages 335–342. IEEE, 2012.
- D. S. Bernstein, S. Zilberstein, and N. Immerman. The complexity of decentralized control of markov decision processes. In *Proceedings of the Sixteenth conference on Uncertainty in artificial intelligence*, pages 32–37. Morgan Kaufmann Publishers Inc., 2000.
- D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein. The complexity of decentralized control of markov decision processes. *Mathematics of operations research*, 27 (4):819–840, 2002.
- C. Boutilier. Planning, learning and coordination in multiagent decision processes. In Proceedings of the 6th conference on Theoretical aspects of rationality and knowledge, pages 195–210. Morgan Kaufmann Publishers Inc., 1996.
- C. Breazeal, C. D. Kidd, A. L. Thomaz, G. Hoffman, and M. Berlin. Effects of nonverbal communication on efficiency and robustness in human-robot teamwork. In *Intelligent Robots and Systems, 2005.(IROS 2005). 2005 IEEE/RSJ International Conference* on, pages 708–713. IEEE, 2005.
- C. Breazeal, J. Gray, and M. Berlin. An embodied cognition approach to mindreading skills for socially intelligent robots. *The International Journal of Robotics Research*, 28(5):656, May 2009.

- P. Brown. *Politeness: Some universals in language usage*, volume 4. Cambridge University Press, 1987.
- G. Buccino, F. Binkofski, and L. Riggio. The mirror neuron system and action recognition. *Brain and language*, 89(2):370–376, 2004.
- S. Chernova and M. Veloso. Teaching collaborative multi-robot tasks through demonstration. In *IEEE-RAS International Conference on Humanoid Robots*, Daejeon, Korea, December 2008.
- R. Cipolla and N. Hollinghurst. Human-robot interface by pointing with uncalibrated stereo vision. *Image and Vision Computing*, 14(3):171–178, 1996.
- C. Crick and B. Scassellati. Inferring narrative and intention from playground games. Proceedings of the 7th IEEE International Conference on Development and Learning (ICDL'08), pages 13–18, August 2008. doi: 10.1109/DEVLRN.2008.4640798.
- R. Desimone and J. Duncan. Neural mechanisms of selective visual attention. Annual review of neuroscience, 18(1):193–222, 1995.
- P. Doshi, Y. Zeng, and Q. Chen. Graphical models for interactive pomdps: representations and solutions. Autonomous Agents and Multi-Agent Systems, 18(3):376–416, 2009.
- A. D. Dragan, K. C. Lee, and S. S. Srinivasa. Legibility and predictability of robot motion. In Human-Robot Interaction (HRI), 2013 8th ACM/IEEE International Conference on, pages 301–308. IEEE, 2013.
- J. Fasola and M. J. Matarić. A socially assistive robot exercise coach for the elderly. *Journal of Human-Robot Interaction*, 2(2):3–32, Jun 2013. URL http://robotics.usc.edu/publications/793/.
- D. Feil-Seifer and M. J. Mataric. Defining socially assistive robotics. In *Rehabilitation Robotics*, 2005. ICORR 2005. 9th International Conference on, pages 465–468. IEEE, 2005.
- S. R. Fussell, R. E. Kraut, and J. Siegel. Coordination of communication: Effects of shared visual context on collaborative work. In *Proceedings of the 2000 ACM conference on Computer supported cooperative work*, pages 21–30. ACM, 2000.
- D. Gergle. The value of shared visual information for task-oriented collaboration. PhD thesis, Carnegie Mellon University, 2006.
- B. P. Gerkey and M. J. Matarić. A formal analysis and taxonomy of task allocation in multi-robot systems. *The International Journal of Robotics Research*, 23(9):939–954, 2004.

- M. Gombolay, R. Wilcox, and J. A. Shah. Fast scheduling of multi-robot teams with temporospatial constraints. In *Robotics: Science and Systems*, 2013.
- B. J. Grosz and S. Kraus. The evolution of sharedplans. In Foundations of rational agency, pages 227–262. Springer, 1999.
- B. J. Grosz and C. L. Sidner. Attention, intentions, and the structure of discourse. Computational linguistics, 12(3):175–204, 1986.
- R. Hamid, S. Maddi, A. Johnson, A. Bobick, I. Essa, and C. Isbell. A novel sequence representation for unsupervised analysis of human activities. *Artificial Intelligence*, 173(14):1221–1244, 2009. ISSN 0004-3702.
- J. E. Hanna and M. K. Tanenhaus. Pragmatic effects on reference resolution in a collaborative task: Evidence from eye movements. *Cognitive Science*, 28(1):105–115, 2004.
- Y. Hato, S. Satake, T. Kanda, M. Imai, and N. Hagita. Pointing to space: modeling of deictic interaction referring to regions. In *Proceeding of the 5th ACM/IEEE* international conference on Human-robot interaction, pages 301–308. ACM, 2010.
- F. Heider and M. Simmel. An experimental study of apparent behavior. The American Journal of Psychology, pages 243–259, 1944.
- G. Hoffman and C. Breazeal. Cost-based anticipatory action selection for human–robot fluency. *Robotics, IEEE Transactions on*, 23(5):952–961, 2007.
- G. Hoffman and C. Breazeal. Effects of anticipatory perceptual simulation on practiced human-robot tasks. Autonomous Robots, 28(4):403–423, May 2010.
- C.-M. Huang and B. Mutlu. Modeling and evaluating narrative gestures for humanlike robots. *Proceedings of the 9th R: SS*, 2013a.
- C.-M. Huang and B. Mutlu. The repertoire of robot behavior: Enabling robots to achieve interaction goals through social behavior. *Journal of Human-Robot Interaction*, 2(2): 80–102, 2013b.
- L. Itti, C. Koch, and E. Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 20(11):1254–1259, 1998. ISSN 0162-8828. doi: http://doi.ieeecomputersociety.org/10.1109/34.730558.
- L. P. Kaelbling, M. L. Littman, and A. W. Moore. Reinforcement learning: A survey. arXiv preprint cs/9605103, 1996.

- R. Kelley, C. King, A. Tavakkoli, M. Nicolescu, 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*, 5(2): 1–22, 2008.
- S. Kelly, A. Ozurek, and E. Maris. Two sides of the same coin: Speech and gesture mutually interact to enhance comprehension. *Psychological Science*, 21(2):260–267, 2009.
- B. Keysar, D. J. Barr, J. A. Balin, and J. S. Brauner. Taking perspective in conversation: The role of mutual knowledge in comprehension. *Psychological Science*, 11(1):32–38, 2000.
- D. Kirsh. The intelligent use of space. Artificial intelligence, 73(1):31-68, 1995.
- T. Kollar, S. Tellex, D. Roy, and N. Roy. Toward understanding natural language directions. In *Human-Robot Interaction (HRI)*, 2010 5th ACM/IEEE International Conference on, pages 259–266. IEEE, 2010.
- D. Kortenkamp, E. Huber, and R. Bonasso. Recognizing and interpreting gestures on a mobile robot. In *Proceedings of the National Conference on Artificial Intelligence*, pages 915–921, 1996.
- R. Kraut, S. Fussell, and J. Siegel. Visual information as a conversational resource in collaborative physical tasks. *Human-computer interaction*, 18(1):13–49, 2003.
- C. Lamm, C. D. Batson, and J. Decety. The neural substrate of human empathy: effects of perspective-taking and cognitive appraisal. *Journal of cognitive neuroscience*, 19 (1):42–58, 2007.
- P. A. Lasota, G. F. Rossano, and J. A. Shah. Toward safe close-proximity human-robot interaction with standard industrial robots. In *Automation Science and Engineering* (CASE), 2014 IEEE International Conference on, pages 339–344. IEEE, 2014.
- J. Lasseter. Principles of traditional animation applied to 3d computer animation. In ACM Siggraph Computer Graphics, volume 21, pages 35–44. ACM, 1987.
- J. Lee and S. Marsella. Nonverbal behavior generator for embodied conversational agents. In *Intelligent virtual agents*, pages 243–255. Springer, 2006.
- S. P. Levine, D. A. Bell, L. A. Jaros, R. C. Simpson, Y. Koren, and J. Borenstein. The navchair assistive wheelchair navigation system. *Rehabilitation Engineering*, *IEEE Transactions on*, 7(4):443–451, 1999.
- C. B. Lockridge and S. E. Brennan. Addressees' needs influence speakers' early syntactic choices. *Psychonomic Bulletin & Review*, 9(3):550–557, 2002.

- 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.
- M. Marjanovic, B. Scassellati, and M. Williamson. Self-taught visually guided pointing for a humanoid robot. In From Animals to Animats 4: Proceedings of the Fourth International on Conference Simulation of Adaptive Behavior, pages 35–44, 1996.
- E. Martin, D. R. Lyon, and B. T. Schreiber. Designing synthetic tasks for human factors research: An application to uninhabited air vehicles. In *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, volume 42, pages 123–127. SAGE Publications, 1998.
- M. J. Matarić and M. Pomplun. Fixation behavior in observation and imitation of human movement. *Cognitive Brain Research*, 7(2):191–202, 1998.
- R. Mayberry and J. Jaques. Gesture production during stuttered speech: Insights into the nature of gesture-speech integration, chapter 10, pages 199–214. Cambridge University Press, 2000.
- D. McNeill. Gesture and thought. University of Chicago Press, 2008.
- A. Mehrabian. Nonverbal communication. Transaction Publishers, 1977.

Microsoft. Microsoft kinect, 2010. URL http://www.microsoft.com/en-us/kinectforwindows/.

- B. Mutlu, T. Shiwa, T. Kanda, H. Ishiguro, and N. Hagita. Footing in human-robot conversations: how robots might shape participant roles using gaze cues. In *Proceed*ings of the 4th ACM/IEEE International Conference on Human Robot Interaction, (HRI'09), 2009.
- B. Mutlu, T. Kanda, J. Forlizzi, J. Hodgins, and H. Ishiguro. Conversational gaze mechanisms for humanlike robots. ACM Transactions on Interactive Intelligent Systems (TiiS), 1(2):12, 2012.
- B. Mutlu, A. Terrell, and C.-M. Huang. Coordination mechanisms in human-robot collaboration. In *Proceedings of the Workshop on Collaborative Manipulation*, 8th ACM/IEEE International Conference on Human-Robot Interaction, 2013.
- R. Nair, M. Tambe, M. Yokoo, D. Pynadath, and S. Marsella. Taming decentralized pomdps: Towards efficient policy computation for multiagent settings. In *IJCAI*, pages 705–711, 2003.
- K. Nickel and R. Stiefelhagen. Visual recognition of pointing gestures for human-robot interaction. *Image and Vision Computing*, 25(12):1875–1884, 2007.

- S. Nikolaidis and J. Shah. Human-robot teaming using shared mental models. *ACM/IEEE HRI*, 2012.
- S. Nikolaidis and J. Shah. Human-robot cross-training: computational formulation, modeling and evaluation of a human team training strategy. In *Proceedings of the* 8th ACM/IEEE international conference on Human-robot interaction, pages 33–40. IEEE Press, 2013.
- S. Nikolaidis, R. Ramakrishnan, K. Gu, and J. Shah. Efficient model learning from joint-action demonstrations for human-robot collaborative tasks. In *Proceedings of* the Tenth Annual ACM/IEEE International Conference on Human-Robot Interaction, pages 189–196. ACM, 2015.
- A. Ozyurek. Do speakers design their co-speech gestures for their addresses? the effects of addressee location on representational gestures. *Journal of Memory and Language*, 46(4):688–704, 2002.
- L. Peshkin, K.-E. Kim, N. Meuleau, and L. P. Kaelbling. Learning to cooperate via policy search. arXiv preprint arXiv:1408.1484, 2014.
- O. Pettersson. Execution monitoring in robotics: A survey. *Robotics and Autonomous Systems*, 53(2):73–88, 2005.
- M. J. Pickering and S. Garrod. Toward a mechanistic psychology of dialogue. *Behavioral and brain sciences*, 27(02):169–190, 2004.
- P. Pook and D. Ballard. Deictic human/robot interaction. Robotics and Autonomous Systems, 18(1-2):259–269, 1996.
- D. V. Pynadath and M. Tambe. The communicative multiagent team decision problem: Analyzing teamwork theories and models. *Journal of Artificial Intelligence Research*, pages 389–423, 2002.
- Rethink Robotics. Baxter manufacturing robot, 2012.
- G. Rizzolatti and C. Sinigaglia. *Mirrors in the Brain: How our minds share actions, emotions, and experience.* Oxford University Press, Feb 2008.
- S. Rosenthal, J. Biswas, and M. Veloso. An effective personal mobile robot agent through symbiotic human-robot interaction. In *Proceedings of the 9th International Conference* on Autonomous Agents and Multiagent Systems: volume 1-Volume 1, pages 915–922. International Foundation for Autonomous Agents and Multiagent Systems, 2010.
- M. Roth, R. Simmons, and M. Veloso. What to communicate? execution-time decision in multi-agent pomdps. In *Distributed Autonomous Robotic Systems* 7, pages 177–186. Springer, 2006.

- A. Sauppé and B. Mutlu. The social impact of a robot co-worker in industrial settings.
- A. Sauppé and B. Mutlu. Effective task training strategies for instructional robots. In *Proceedings of the 10th annual Robotics: Science and Systems Conference*, 2014a.
- A. Sauppé and B. Mutlu. How social cues shape task coordination and communication. In Proceedings of the 17th ACM conference on Computer supported cooperative work & social computing, pages 97–108. ACM, 2014b.
- B. Scassellati. Theory of mind for a humanoid robot. Autonomous Robots, 12(1):13–24, 2002. ISSN 0929-5593.
- B. Scassellati. Investigating models of social development using a humanoid robot. In *Proceedings of the International Joint Conference on Neural Networks*, volume 4, pages 2704 – 2709 vol.4, jul. 2003. doi: 10.1109/IJCNN.2003.1223995.
- M. Scopelliti, M. V. Giuliani, and F. Fornara. Robots in a domestic setting: a psychological approach. Universal Access in the Information Society, 4(2):146–155, 2005.
- J. Shah and C. Breazeal. An empirical analysis of team coordination behaviors and action planning with application to human-robot teaming. *Human Factors*, 2010.
- J. Shah, J. Wiken, B. Williams, and C. Breazeal. Improved human-robot team performance using chaski, a human-inspired plan execution system. In *Proceedings of the* 6th international conference on Human-robot interaction, pages 29–36. ACM, 2011.
- G. Shani, J. Pineau, and R. Kaplow. A survey of point-based pomdp solvers. Autonomous Agents and Multi-Agent Systems, 27(1):1–51, 2013.
- C. Sidner, C. Kidd, C. Lee, and N. Lesh. Where to look: a study of human-robot engagement. In *Proceedings of the 9th International Conference on Intelligent User Interfaces*, page 84. ACM, 2004.
- R. B. Smith, R. Hixon, and B. Horan. Supporting flexible roles in a shared space. In *Collaborative Virtual Environments*, pages 160–176. Springer, 2001.
- I. D. Steiner. Group process and productivity (social psychological monograph). 2007.
- O. Sugiyama, T. Kanda, M. Imai, H. Ishiguro, N. Hagita, and Y. Anzai. Humanlike conversation with gestures and verbal cues based on a three-layer attention-drawing model. *Connection science*, 18(4):379–402, 2006.
- A. Tapus, M. Maja, and B. Scassellatti. The grand challenges in socially assistive robotics. *IEEE Robotics and Automation Magazine*, 14(1):N–A, 2007.
- F. Thomas, O. Johnston, and F. Thomas. *The illusion of life: Disney animation*. Hyperion New York, 1995.

- A. Thomaz and C. Breazeal. Reinforcement learning with human teachers: Evidence of feedback and guidance with implications for learning performance. In *Proceedings of the 21st National Conference on Artificial Intelligence*, 2006.
- M. Tomasello et al. Why we cooperate, volume 206. MIT press Cambridge, MA, 2009.
- C. Torrey, S. R. Fussell, and S. Kiesler. How a robot should give advice. In Human-Robot Interaction (HRI), 2013 8th ACM/IEEE International Conference on, pages 275–282. IEEE, 2013.
- J. G. Trafton, N. L. Cassimatis, M. D. Bugajska, D. P. Brock, F. E. Mintz, and A. C. Schultz. Enabling effective human-robot interaction using perspective-taking in robots. Systems, Man and Cybernetics, Part A: Systems and Humans, IEEE Transactions on, 35(4):460–470, 2005.
- J. G. Trafton, M. D. Bugajska, B. R. Fransen, and R. M. Ratwani. Integrating vision and audition within a cognitive architecture to track conversations. In *Proceedings* of the 3rd ACM/IEEE international conference on Human robot interaction, pages 201–208. ACM, 2008.
- T. Ullman, C. Baker, O. Macindoe, O. Evans, N. Goodman, and J. Tenenbaum. Help or hinder: Bayesian models of social goal inference. Advances in Neural Information Processing Systems (NIPS), 22, 2010.
- D. Walther, L. Itti, M. Riesenhuber, T. Poggio, and C. Koch. Attentional selection for object recognition: a gentle way. In *Biologically Motivated Computer Vision*, pages 251–267. Springer, 2010.
- S. Whittaker. Things to talk about when talking about things. *Human–Computer* Interaction, 18(1-2):149–170, 2003.
- R. Wilcox, S. Nikolaidis, and J. Shah. Optimization of temporal dynamics for adaptive human-robot interaction in assembly manufacturing. *Robotics*, page 441, 2013.
- N. Wong and C. Gutwin. Where are you pointing?: the accuracy of deictic pointing in cves. In *Proceedings of the 28th international conference on Human factors in computing systems*, pages 1029–1038. ACM, 2010.