

SOCIALLY ASSISTIVE AND SERVICE ROBOTICS
FOR OLDER ADULTS:
METHODOLOGIES FOR MOTIVATING EXERCISE
AND FOLLOWING SPATIAL LANGUAGE INSTRUCTIONS
IN DISCOURSE

by

Juan P. Fasola

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)

June 2014

Copyright 2014

Juan P. Fasola

Dedication

I would like to dedicate this Dissertation to my family. You have always been there for me, with love and guidance, during my Ph.D. and throughout my whole life. I will be forever grateful for all of your unconditional love and support; this is for you.

Acknowledgements

I would like to thank my advisor, Maja Matarić, for giving me the opportunity to come to the University of Southern California to perform research in her lab alongside great people, and to work on robotics projects that have an impact on the real world with user populations in need. I am very grateful for Maja's guidance and support throughout my Ph.D., and for giving me the freedom to explore the field and work on research projects that I was passionate about.

I would also like to express my gratitude to the other members of my Qualifying and Dissertation Defense Committee: Gaurav Sukhatme, Aaron Hagedorn, Stefan Schaal, and David Traum for their invaluable feedback and guidance regarding my research during the final year of my Ph.D.

I would like to thank all members of the Interaction Lab, for providing me with excellent feedback on my research during lab meetings and practice talks, and for many entertaining discussions in the lab (research related and otherwise).

In particular, I would like to thank Ross Mead, Pierre Johnson, Eric Wade, and Arvind Pereira (from RESL) for the many discussions and adventures we've had while exploring what LA has to offer. Whether it was running into celebrities in Hollywood, lounging in Santa Monica, or eating half price burgers on Rodeo Drive, you guys made my Ph.D. experience so much more fun.

I would also like to acknowledge Dave Feil-Seifer, my desk neighbor of many years and fellow TA, who helped me with many Gentoo and Player/Stage issues before Ubuntu/ROS took over everything, and was the go-to guy for really anything Bandit or Pioneer related; thank you Dave.

I am also very grateful to the undergraduate and MS students who helped me with different aspects of my research work, including: Alex Neiss and Gregory Koch for developing speech recognition tools and integrating the POS tag parser; Aras Akbari for conducting the off-campus SAR exercise coach study; and Farrokh Langroodi for conducting the on-campus SAR exercise coach study, constructing the Bandit gripper mechanism, helping with the spatial language user study, and for always being ready to assist me with last minute demos (thank you Farrokh!).

I especially would like to thank my students Yixin Cai and Han Xiao. In the last year of my Ph.D., Yixin and Han helped me to port my spatial language framework to Gazebo/ROS for realistic 3D simulation and testing, integrated the framework for use with the PR2, and worked diligently to develop the AI algorithms (vision and laser-based object recognition, motion planning, etc.) needed for Bandit to pick-and-place objects reliably with end users. Without their help I would have never finished my Ph.D. on time. Thank you Yixin and Han for all your hard work, I really appreciate it!

I am also grateful to Elena Prieto, who has been by my side with love and support, has helped me to grow as a person, and who has made the last years of my Ph.D. the most memorable.

Lastly, I would like to thank my loving family. Without their love, guidance and support, none of this would have been possible. In particular, I would like to acknowledge my brother, Carlos A. Fasola, for providing linguistics expertise, helpful suggestions, insightful feedback, and many fruitful discussions regarding my research work throughout the years. I dedicate this Dissertation to all of you.

Table of Contents

Dedication	ii
Acknowledgements	iii
List of Figures	ix
List of Tables	xiv
Abstract	xvi
Chapter 1: Introduction	1
1.1 SAR Exercise Coach	5
1.1.1 Motivating Behavior Change	5
1.1.2 SAR Exercise Coach Methodology and User Studies	6
1.2 Spatial Language-Based HRI Methodology	8
1.2.1 Spatial Language Communication in HRI	8
1.2.2 Spatial Language Framework and HRI User Study	10
1.3 Dissertation Contributions	12
1.4 Outline	15
Chapter 2: Background and Related Work	17
2.1 Service Robotics	17
2.1.1 Overview of Service Robotics	17
2.1.2 Service Robots for Older Adults	18
2.2 Socially Assistive Robotics	18
2.2.1 SAR for Older Adults	18
2.2.2 Social Agent Coaches	19
2.2.3 The Effect of Embodiment	20
2.3 Spatial Language-Based HRI	22
2.3.1 Modeling Static Spatial Relations	22
2.3.2 Modeling Dynamic Spatial Relations	23
2.3.3 Robot Control Language	24

2.4	Summary	24
Chapter 3: Socially Assistive Robot Design Principles		26
3.1	Design Principles for a SAR Coach	26
3.1.1	Motivating	27
3.1.2	Fluid and Highly Interactive	28
3.1.3	Personable	28
3.1.4	Intelligent	29
3.1.5	Task-Driven	30
3.2	Summary	31
Chapter 4: Socially Assistive Robot Exercise Coach		33
4.1	Interaction Scenario	33
4.2	Robot Platform	35
4.3	Exercise Games	35
4.3.1	Workout Game	36
4.3.2	Sequence Game	36
4.3.3	Imitation Game	37
4.3.4	Memory Game	37
4.4	Software Architecture	38
4.5	Visual User Activity Recognition Algorithm	41
4.6	Robot Behaviors	45
4.6.1	Interaction Flow and Behavior Management	45
4.6.2	Feedback Procedure	48
4.7	Summary	49
Chapter 5: Intrinsic Motivation User Studies with Older Adults		50
5.1	Motivation Study 1: Praise and Relational Discourse	50
5.1.1	Study Design	51
5.1.2	Participant Statistics	52
5.1.3	Measures	52
5.1.4	Hypotheses	56
5.1.5	Results	57
5.1.6	Discussion	62
5.2	Motivation Study 2: User Choice and Self-Determination	63
5.2.1	Study Design	63
5.2.2	Participant Statistics	64
5.2.3	Measures	65
5.2.4	Hypotheses	66
5.2.5	Results	67
5.2.6	Discussion	71
5.3	Summary	73

Chapter 6: Embodiment and SAR Evaluation User Study	74
6.1 Study Design	74
6.1.1 Robot Platforms	75
6.1.2 Between-Subjects Design	76
6.2 Participant Statistics	76
6.3 Measures	77
6.3.1 Evaluation of Interaction	77
6.3.2 Evaluation of Robot	77
6.3.3 User Performance Measures	79
6.3.4 Relation to Design Principles	79
6.4 Hypotheses	81
6.5 Results	81
6.5.1 Embodiment Comparison Results	81
6.5.2 SAR System Evaluation Results	85
6.5.3 Study Expansion with Young Adults	93
6.6 Summary	93
Chapter 7: Spatial Language Understanding for Human-Robot Interaction	95
7.1 Approach	95
7.1.1 Semantic Fields	95
7.1.2 Modeling Dynamic Spatial Relations	97
7.2 Robot Software Architecture	100
7.2.1 Syntactic Parser	100
7.2.2 Grounding Noun Phrases	101
7.2.3 Semantic Interpreter	104
7.2.4 Planning	105
7.3 Modeling DSR Representations of “To”, “Through”, and “Around”	106
7.3.1 “To” Representation	106
7.3.2 “Through” Representation	108
7.3.3 “Around” Representation	110
7.4 Generating Paths for Dynamic Spatial Relations	112
7.4.1 “To” Path Generation	113
7.4.2 “Through” Path Generation	113
7.4.3 “Around” Path Generation	114
7.5 Parsing Spatial Language Instructions with Figure Objects	116
7.5.1 Extension for Figure Objects	116
7.5.2 Pruning Multiple Parses of a Single Instruction	117
7.6 Object Pick-and-Place Movement Planning	120
7.6.1 Object Pick Up Planning with Grasp Fields	120
7.6.2 Object Placement Planning with Semantic Fields	122
7.6.3 Pragmatic Fields for Object Placement Planning	122
7.7 Summary	124

Chapter 8: Evaluation of Spatial Language-Based HRI Methodology	125
8.1 Semantic Inference Accuracy	125
8.2 Instruction Following Results	127
8.3 Instruction Sequence Following Results	130
8.4 Speech Recognition Results	134
8.5 Instruction Following Results with SLAM Maps	135
8.6 Object Pick-and-Place Task Results	139
8.7 Summary	141
Chapter 9: Spatial Language Discourse with Pragmatics for HRI	143
9.1 Interpreting Instructions in Discourse	144
9.1.1 Probabilistic Extraction of Instruction Sequences	144
9.1.2 Reference Resolution	147
9.2 Pragmatics for Physically Embodied Interaction with People	151
9.3 Generalized Transfer to Robot Systems	153
9.4 Evaluation Results	155
9.5 Summary	159
Chapter 10: Spatial Language-Based HRI User Study with Older Adults	162
10.1 Study Design	162
10.1.1 Virtual Robot Condition	163
10.1.2 Physical Robot Condition	168
10.2 Participant Statistics	172
10.3 Measures	172
10.3.1 Objective Measures	172
10.3.2 Subjective Measures	174
10.4 Results	175
10.4.1 Virtual Robot Condition Results	175
10.4.2 Physical Robot Condition Results	177
10.4.3 Results of Framework Modification for Prepositional Phrase Attachments	180
10.4.4 Spatial Language Usage Statistics	185
10.4.5 Subjective Evaluation Results	189
10.5 Summary	192
Chapter 11: Summary	193
11.1 SAR Coach for Motivating Therapeutic Behavior	194
11.2 Spatial Language-Based HRI Framework	195
11.3 Limitations and Future Work	196
Bibliography	197

List of Figures

4.1	(a) The setup for the one-on-one interaction between user and robot coach; (b) Bandit humanoid torso robot.	34
4.2	Diagram of system architecture, including the six system modules: vision and world model, speech, user communication, behaviors, action output, and database management.	39
4.3	(a), (c), (d) Example face and arm angle detection results; (b) segmented image.	42
4.4	Diagram of the behavior module’s finite state machine (FSM). The FSM illustrates the flow of the interaction during the exercise sessions. The white arrow emanating from the Break Prompt state denotes a user decision.	47
5.1	(a) Plot of participant evaluations of the interaction, in terms of enjoyableness and usefulness, for both study conditions; (b) plot of participant evaluations of the robot (as a companion, exercise coach, and level of social presence) for both study conditions. Note: significant differences are marked by asterisks (*).	58
5.2	Graphs of: (a) the participants’ preferences of study condition; (b) the participants’ ratings in response to survey questions on their perception of the robot’s intelligence, helpfulness, their mood during sessions, and how important the sessions were to them; (c) the participants’ preferences of exercise game.	70
6.1	(a) Physical robot; (b) virtual robot computer simulation; (c) virtual robot on the screen, with camera.	75

6.2	(a) Plot of participant evaluations of the interaction of the SAR exercise system in terms of enjoyableness and usefulness; (b) plot of participant evaluations of the robot coach of the SAR exercise system in terms of helpfulness, intelligence, social presence, and as a companion. Note: Significant differences ($p < .05$) in comparison to neutral rating distribution are marked by asterisks (*).	86
6.3	Participant performance results across all four sessions of interaction for both study groups, showing (a) Average gesture completion time (Workout game); (b) feedback percentage (Workout game). Note: Means are least squares means from the ANOVA; error bars are standard error. The statistical difference between the physical and virtual robot in the first session is significant for completion time ($p < .01$), and marginally significant for feedback percentage ($p = .07$), see text.	92
7.1	Semantic fields for static prepositions (a) near; (b) away from; (c) between.	96
7.2	Semantic field for “along” (a) near subfield; (b) direction subfield ($90^\circ = \text{red}$, $0^\circ = \text{blue}$); (c) combined field.	97
7.3	Two example alterations to semantic field due to user feedback statement “Move a little away from the wall”. (a) Field moved entirely; (b) Field mean shifted.	98
7.4	Robot software architecture system modules.	100
7.5	(a) Parse tree for “Go to the table by the kitchen”; (b) Semantic field for ‘near kitchen’ with candidate tables.	102
7.6	Two example paths for “to the kitchen”. (a) Path value = 1.8×10^{-10} ; (b) Path value = 1.0. Note: $\sigma = \text{robot width} \times 2.5$.	107
7.7	Topology of hallway in 2D home environment showing three (configuration space) entrance boundaries.	108
7.8	Path traveling around a dining table reference object, showing start and end orientations, with resulting <i>circumcentric</i> semantic field value = $ 223^\circ /360^\circ = 0.61944$.	111
7.9	Two paths for “Go around the bed” with/without enforcing visibility region. (a) Path value = 0.857; (b) Path value = 0.895.	115
7.10	Apartment environment with cup locations and <i>on</i> (bookcase) semantic field shown.	119

7.11	(a) Proximity-based grasp field for teddy bear object; (b) Orientation-based grasp field for cup with handles; (c) Task solution for “Take the cup to the kitchen” with grasp field and semantic field shown at the pick (1) and place (2) locations, respectively.	121
7.12	(a) Pragmatic field indicating suitable surfaces for object placement (weight values in grayscale); (b) Task solution for “Take the cup to the kitchen” incorporating pragmatic constraints, with combined semantic/pragmatic field shown at object placement location.	123
8.1	Executed paths and semantic fields (command = blue, constraint = red) for test runs (a) run 1; (b) run 2; (c) run 3; (d) run 4.	128
8.2	Semantic field values along execution paths in test runs (a) run 1; (b) run 2; (c) run 3; (d) run 4.	129
8.3	DSR path generation results for entire instruction sequence with and without user-specified constraints. (a) Test Run #1 (no constraints); (b) Test Run #2 (constraints). Note: path endpoints for each instruction are labeled with the instruction number.	132
8.4	Semantic field values along task execution paths in Test Run #1. (a) <i>circumcentric</i> field value along solution path for instruction 1; (b) <i>at</i> field value along solution path for instruction 5.	134
8.5	Executed paths and semantic fields (command = blue, constraint = red) for test runs (a) run 5; (b) run 6; (c) run 7; (d) run 8.	137
8.6	Executed paths and semantic fields (command = blue, constraint = red) for test runs (a) run 9; (b) run 10; (c) run 11; (d) run 12.	138
8.7	Robot execution paths for test runs. (a) Run #1 (no constraints); (b) Run #2 (constraints).	140
8.8	PR2 robot executing task for “Put the Coke can on the coffee table” in 3D household environment using Gazebo simulator/ROS framework. . .	141
9.1	Reference resolution example for the instruction “toss it in the sink” expressed by the user during spatial language discourse. Parsed NPs of the natural language input are shown in brackets with their corresponding unique grounding ID numbers as subscripts.	150

9.2	(a) SLAM map of laboratory space with pragmatic field for robot safety shown; (b) Example robot approach behavior with combined semantic/pragmatic field shown for at/person safety.	152
9.3	Dynamic obstacle avoidance for instruction “Go to the dinner table”. (a) Planned path (green) and actual path (red); (b) Visualization of obstacles detected in robots local map and global plan after robot re-planning.	154
9.4	Combined semantic/pragmatic field and execution result for the task “Stand between the printer desk and the whiteboard”.	156
9.5	(From left to right) Planned (green) and executed paths (red), cost map used for navigation planning with AMCL particles corresponding to robot position estimates, and photograph of PR2 robot just before task termination for test runs 1-2. (a),(b),(c) run 1; (d),(e),(f) run 2.	158
9.6	Test run 3 results. (a) Planned path (green) and actual path (red) with semantic/pragmatic fields calculated for hand-off behavior; (b) PR2 robot handing bottle (grounded object referent) to intended person (ground referent for “him”) during task execution.	160
10.1	(a) Virtual robot condition setup; (b) 2D computer simulated household environment, with example robot task execution path shown for instruction “Pick up the medicine in the guest bathroom”.	164
10.2	Annotated environment maps provided to the user during the session. (a) Room names; (b) appliances/furniture item names.	166
10.3	Example target objects and placement locations for scenarios 3 and 4 of virtual robot condition. (a) target object left of stove; (b) target object by kitchen sink; (c) target location on coffee table; (d) target location left of kitchen sink.	167
10.4	Physical robot condition. (a) View of interaction setup with labeled tables representing typical household areas (from left to right: coffee table, bedroom, dinner table, and kitchen); (b) Bandit, the physical robot platform; (c) Example household items used in the study (from left to right: plant, milk, medicine, bottle, cereal).	169

10.5	Scenarios 1 and 2 of Physical Robot condition. (a) Object identification and placement scenario, with target object (bottle) and target location (kitchen) marked by green sticky notes; (b), (c) Example task photographs provided to the user displaying task goal states with object target locations (the relative location of objects with respect to one another is important).	171
10.6	Plot of number of dialogue rounds engaged in by the user for each task scenario presented during the Virtual Robot condition as an illustration of the varying level of difficulty among the tasks.	177
10.7	Between-subjects results comparing original framework to new framework capable of interpreting PPAs. (a) Task success rate and round success rate; (b) Average number of clarification queries per dialogue round (Note: Significant results marked by asterisks (*)).	183
10.8	Within-subjects results comparing original framework to new framework capable of interpreting PPAs. (a) Task success rate and round success rate; (b) Average number of clarification queries per dialogue round (Note: Significant results marked by asterisks (*)).	184
10.9	Plot of the task success rate (with repeated attempts) achieved in the Virtual Robot condition by each user of the study ($n = 19$), showing both participant groups (original framework vs. new PPA-capable framework) and with the success rates listed in increasing order.	185
10.10	Word count histograms for spatial language encountered in all $N = 1239$ instructions provided by participants during both study conditions. (a) Verb counts; (b) Path preposition counts; (c) Static preposition counts for those expressed within noun phrases as spatial relations.	190
10.11	Participants evaluation results. (a) Evaluation of interaction and of service robot; (b) Evaluation of interaction with respect to USE questionnaire items.	191

List of Tables

5.1	Participant Responses to Direct Comparison Survey Items	60
5.2	Participant Exercise Performance Statistics	62
6.1	Summary of the relations between the five design principles and their related study measures (evaluation and/or performance based)	80
6.2	Results of between-subjects data comparison for all $n = 33$ older adult participants showing means and standard deviations (in parentheses)	82
6.3	User exercise performance statistics for all $n = 33$ older adult participants engaging with the SAR exercise system, showing means and standard deviations (in parentheses)	88
6.4	Means, standard deviations, and intercorrelations among dependent measures	90
6.5	Analysis of variance testing the fixed effect tests of condition, sessions, and measure on performance	91
7.1	Semantic field values of candidate groundings for NP “the table”	102
7.2	DSR Representations for “Through”	108
7.3	DSR Representations for “Around”	111
7.4	Grammar Constituency Rules for English Directives using Spatial Language	117
7.5	Possible Flags Raised during Parse Pruning	118
8.1	Inference Accuracy of Semantic Interpretation Module	126
8.2	Semantic Inference Results for Instructions and Constraints of Test Runs	129

8.3	Instruction Sequence Given in Test Runs	131
8.4	Results of Semantic Inference and Pragmatics for Test Run Instructions	131
8.5	Speech Recognition Module Accuracy	135
8.6	Instructions Given in Test Runs 5-12	136
8.7	Instruction Sequence for Test Runs with Corresponding Constraints . .	140
9.1	Grammar for Spatial Language Directives	144
9.2	Probabilistic Instruction Sequence Extraction Procedure Example . . .	147
9.3	Accuracy of Final Robot Positions in Spatial Task Experiment	156
9.4	Instructions Given in Test Runs 1-3 with Inference Results for Instruction Sequences	157
10.1	Results of Interaction with Participants ($N = 19$) in Virtual Robot Condition	176
10.2	Results of Interaction with Participants ($N = 19$) in Physical Robot Condition	178
10.3	Example Instructions Given by Participants in the Physical Robot Condition with Inference Results for Instruction Sequences	179
10.4	Spatial Language Statistics of Verb and Path Preposition Usage in $N = 1239$ Total Instructions Given by Participants	186
10.5	Spatial Language Statistics of Static Preposition Usage within Noun Phrases of $N = 1239$ Total Instructions Given by Participants	189

Abstract

The growing population of aging adults is increasing the demand for healthcare services worldwide. Socially assistive robotics (SAR) and service robotics have the potential to aid in addressing the needs of the growing elderly population by promoting health benefits, independent living, and improved quality of life. For such robots to become ubiquitous in real-world human environments, they will need to interact with and learn from non-expert users in a manner that is both natural and practical for the users. In particular, such robots will need to be capable of understanding natural language instructions in order to learn new tasks and receive guidance and feedback on task execution.

Research into SAR and service robotics-based solutions for non-expert users, and in particular older adults, that spans varied assistive tasks generally falls within one of two distinct areas: 1) *robot-guided interaction*, and 2) *user-guided interaction*. This dissertation contributes to both of these research areas.

To address *robot-guided interaction*, this dissertation presents the design methodology, implementation and evaluation details of a novel SAR approach to motivate and engage elderly users in simple physical exercise. The approach incorporates insights from psychology research into intrinsic motivation and contributes five clear design principles for SAR-based therapeutic interventions. To evaluate the approach and its effectiveness in gaining user acceptance and motivating physical exercise, it was implemented as an

integrated system and three user studies were conducted with older adults, to investigate: 1) the effect of praise and relational discourse in the system towards increasing user motivation; 2) the role of user autonomy and choice within the interaction; and 3) the effect of embodiment in the system by comparing user evaluations of similar physically and virtually embodied SAR exercise coaches in addition to evaluating the overall SAR system.

To address *user-guided interactions*, specifically with non-expert users through the use of natural language instructions, this dissertation presents a novel methodology that allows service robots to interpret and follow spatial language instructions, with and without user-specified natural language constraints and/or unvoiced pragmatic constraints. This work contributes a general computational framework for the representation of dynamic spatial relations, with both local and global properties. The methodology also contributes a probabilistic approach in the inference of instruction semantics; a general approach for interpreting object pick-and-place tasks; and a novel probabilistic algorithm for the automatic extraction of contextually and semantically valid instruction sequences from unconstrained spatial language discourse, including those containing anaphoric reference expressions. The spatial language interpretation methodology was evaluated in simulation, on two different physical robot platforms, and in a user study conducted with older adults for validation with target users.

The successful acceptance of the presented SAR and service robotics approach by elderly users, as evidenced by the high participant evaluations of the system and consistent task performance in all of our user studies, validates the approach, design, algorithms, and effectiveness of our SAR and service robotics methodologies, and illustrates the potential of such technology to help older adults achieve beneficial health outcomes and improve quality of life.

Chapter 1

Introduction

The growing population of aging adults is increasing the demand for healthcare services worldwide. By the year 2050, the number of people over the age of 85 will increase threefold [Centers for Disease Control and Prevention, 2003], while the shortfall of nurses and caregivers is already an issue [American Association of Colleges of Nursing, 2010; American Health Care Association, 2008; Buerhaus, 2008]. Regular physical exercise has been shown to be effective at maintaining and improving the overall health of elderly individuals [Baum et al., 2003; Dawe and Moore-Orr, 1995; McMurdo and Rennie, 1993; Thomas and Hageman, 2003]. In addition, physical fitness is associated with higher functioning in the executive control processes [Colcombe and Kramer, 2003] and is correlated with less atrophy of frontal cortex regions [Colcombe et al., 2004] and with improved reaction times [Spirduso and Clifford, 1978]. Social interaction, and specifically high perceived interpersonal social support, has also been shown to have a positive impact on general mental and physical wellbeing [Moak and Agrawal, 2010], in addition to reducing the likelihood of depression [George et al., 1989; Paykel, 1994; Stansfeld et al., 1997; Stice et al., 2004]. Thus, the availability of physical exercise therapy, social interaction, and companionship will be critical for the growing elderly

population; *socially assistive robotics* (SAR) and *service robotics* have the potential to help address this need. Our work addresses challenges in both of these important and related research fields.

A socially assistive robot is a system that employs hands-off *human-robot interaction* (HRI) strategies, including the use of speech, facial expressions, and communicative gestures, to provide assistance in accordance with a particular healthcare or assistive context [Feil-Seifer and Matarić, 2005]. SAR systems equipped with motivational, social, and therapeutic capabilities have the potential to facilitate elderly individuals to live independently in their own homes, to enhance their quality of life, and to improve their overall health.

In addition to addressing the healthcare needs of the growing elderly population, SAR technology has the potential to help many populations in need of healthcare and/or therapeutic assistance. For example, previous SAR research by our group includes systems that were developed for and tested with stroke patients [Matarić et al., 2007; Tapus et al., 2008], Alzheimer’s patients [Tapus et al., 2009], children with autism spectrum disorders [Feil-Seifer and Matarić, 2010, 2012], as well as healthy young adults [Fasola and Matarić, 2010]. These systems all encouraged target users to engage in therapeutic behaviors through *robot-guided interactions*. In SAR-guided interactions, the robot guides a healthcare intervention in one of two ways: 1) as the coach of a particular healthcare task (i.e., direct guidance), or 2) as a catalyst for promoting therapeutic behaviors (i.e., indirect guidance). In either case, the SAR agent is in charge of steering the interaction towards engaging the users in therapeutic behaviors as well as promoting beneficial health outcomes, which is the ultimate goal of the interaction.

In the context of *service robotics*, which is related to SAR but also encompasses robots capable of hands-on assistive HRI, *user-guided interactions* are those that are initiated and guided by the user towards the robot performing some user-specified task.

In those types of interactions, the goal of the service robot is to successfully carry out the user-defined task.

In designing an HRI methodology to address these and other possible user-defined service tasks, it is important to consider the following characteristics of service robot-based user-guided interactions: 1) the users are non-experts; 2) the tasks may be specific to the user's environment/preferences and are unknown in advance (i.e., the system cannot program all robot task solutions *a priori*); and as a corollary 3) the users may need to teach the robot how to perform the tasks. Along with these characteristics, a guiding principle for the HRI methodology is that the method of interaction/teaching should be natural and practical for non-expert users.

In investigating a natural method of HRI/teaching with target users, the work of Scopelliti et al. [2005] identified natural language speech as the most preferred method of human-robot communication over less familiar methods such as keyboard input, graphical user interfaces, touch screen devices, etc. The practicality and advantages of natural language speech are also especially evident in assistive contexts, where the robots are interacting with people with disabilities, age-related (e.g., reduced mobility, limited eyesight) or otherwise (e.g., individuals post-stroke), as the users may not be able to teach the robot new tasks and/or provide instructions/feedback by demonstration. The HRI methodology for SAR user-guided interactions presented in this dissertation utilizes natural language speech as the method for human-robot communication. In particular, the methodology presented enables service robots with the capability to understand natural language instructions which can be used to learn new tasks and receive guidance and feedback on task execution.

In the context of service robots for older adults, Beer et al. [2012] conducted a study to determine the types of tasks that were most desired by older adults for an in-home

service robot to perform, and found that among the items most reported by the participants were the assistive service tasks for: fetching objects, moving objects, finding objects, and cleaning the house. Given this result that *spatially-oriented tasks* are the most requested by older adults to be completed by service robots (e.g., fetch/move/find objects), and that the most preferred method of human-robot communication was found to be natural language speech, the HRI methodology for service robot user-guided interactions presented in this dissertation is further informed by, among others, the work of Carlson and Hill [2009] which showed that when instructing spatially-oriented tasks, users often utilize preposition-based spatial relations over precise numerical/quantitative descriptions (e.g., “Put the cup *on the table*” vs. “Move the cup *2.3 meters forward*”).

This dissertation aims to contribute novel methodologies for robot-guided and user-guided interactions, with a particular focus in the domain of assistive robots for older adults. In the area of robot-guided interaction, it presents the approach, design methodology, and implementation details of a novel SAR system that aims to motivate and engage elderly users in simple physical exercise. In the area of user-guided interaction, this dissertation presents an approach for enabling autonomous service robots to follow natural language commands from non-expert users, including under user specified constraints, with a particular focus on *spatial language* understanding.

This chapter begins with an overview of this dissertation’s SAR approach, discussing related research in psychology from which it was derived, and the user studies that were conducted with target users to evaluate the approach with a SAR exercise coach. Next, the chapter presents an overview of the approach and evaluation of this dissertation’s spatial language interpretation and HRI methodology. The chapter then follows with a summary of the primary and secondary contributions of this dissertation, and concludes with an outline of each of the dissertation’s chapters.

1.1 SAR Exercise Coach

1.1.1 Motivating Behavior Change

Motivation is a fundamental tool in establishing adherence to a therapy regimen or task scenario and in promoting behavior change. Literature in psychology identifies two forms of motivation: *intrinsic motivation*, which comes from within a person, and *extrinsic motivation*, which comes from sources external to a person [Deci and Ryan, 1985]. Extrinsic motivation, though effective for short-term task compliance, has been shown to be less effective than intrinsic motivation for long-term task compliance and behavior change [Dienstbier and Leak, 1976].

Intrinsic motivation, however, can be, and often is, affected by external factors. In a task scenario, the instructor can impact the user's intrinsic motivation through verbal feedback. Praise, for example, is considered a form of positive feedback and has the potential to increase the user's intrinsic motivation for performing the task, whereas criticism, a form of negative feedback, tends to negatively impact the user's intrinsic motivation [Vallerand, 1983; Vallerand and Reid, 1984]. The effect of positive feedback, however, is closely tied to the user's own perceived competence at the task. Once the user believes he is competent at the task, additional praise no longer affects his intrinsic motivation.

Indirect competition, wherein the user is challenged to compete against an ideal outcome, has also been shown to increase user enjoyment on an otherwise non-competitive task [Weinberg and Ragan, 1979]. For example, when the user is shown her high score on the task, her intrinsic motivation for the task tends to increase, as she strives to better her previous performance. Thus, in a task scenario, it is important that the task instructor continually report to the user his/her performance scores during the task, for motivational purposes.

Verbal feedback provided to the user by the instructor certainly plays an important role in task-based motivation, but the task itself and how it is presented to the user perhaps plays an even more significant role. Csikszentmihalyi’s research suggests that “when one engages in an optimally challenging activity with respect to one’s capacities, there is a maximal probability for task-involved enjoyment or flow” [Csikszentmihalyi, 1975]. He also states that intrinsically motivated activities are those characterized by enjoyment. Simply put, people are “intrinsically motivated under conditions of optimal challenge” [Deci and Ryan, 1985]. If a task is below the optimal challenge level, it is too easy for the user and results in boredom. Alternatively, if the task is above the optimal challenge level, it is too hard and causes the user to get anxious or frustrated. Therefore, an instructor that oversees user performance in a task scenario must be able to continually adjust the task to meet the appropriate needs of the user in order to increase or maintain intrinsic motivation to perform the task.

Another task characteristic with the potential to influence user enjoyment is the incorporation of direct user input. Studies have shown that tasks that support user autonomy and self-determination lead to increased intrinsic motivation, self-esteem, creativity, and other related variables among the participants [Fisher, 1978], all of which are important for achieving task adherence and long-term behavior change. Self-determination, represented in the task in the form of choice of activity [Zuckerman et al., 1978], choice of difficulty level [Fisher, 1978], and choice of rewards [Margolis and Mynatt, 1986], has been shown to either increase or be less detrimental to intrinsic motivation than similar task conditions not involving choice.

1.1.2 SAR Exercise Coach Methodology and User Studies

This dissertation presents the approach, design methodology, and implementation details of a novel SAR system developed to motivate and engage elderly users in simple

physical exercise. Our SAR approach incorporates insights from psychology research into *intrinsic motivation*, discussed above, and contributes clear design principles developed to maximize the probability of success of SAR-based therapeutic interventions.

To validate our SAR approach and its effectiveness in gaining user acceptance and motivating physical exercise, three user studies were conducted with older adults, to investigate: 1) the effect of praise and relational discourse in the system towards increasing user motivation (13 participants); 2) the role of user autonomy and choice within the interaction (24 participants); and 3) the effect of embodiment in the system by comparing user evaluations of similar physically and virtually embodied SAR exercise coaches in addition to evaluating the overall SAR system (66 participants).

Our user studies with older adults were conducted to evaluate the effectiveness of our robot exercise system across a variety of user performance and evaluation measures. The results of the study validate the system approach and its effectiveness in motivating physical exercise in older adults; the participants consistently engaged in physical exercise throughout the interaction sessions, rated the SAR system interaction highly in terms of enjoyableness and usefulness, and rated the robot coach highly in terms of helpfulness, social attraction, social presence, and companionship.

The third user study conducted also served to investigate the role of physical embodiment, a fundamental topic in HRI, by comparing the effectiveness of a physically embodied robot coach to that of a similar virtually embodied robot coach (a computer simulation of the same robot). Physically embodied agents appear to possess what Lee [2004] refers to as “social presence” to a greater extent than virtually embodied agents do. Social presence mediates how people respond to both embodied and disembodied agents and strongly influences the relative success of the social interaction. It is thus important to note that the embodiment type of a socially assistive agent can influence its effectiveness in social interaction, relationship building, gaining user acceptance and

trust, and ultimately in achieving the desired health outcomes of therapeutic intervention. For these reasons, we explored the role of physical embodiment in our SAR exercise system for the elderly.

The results of the embodiment comparison show a strong preference among the participants for the physically embodied robot over the virtually embodied robot and demonstrate the positive effect that physical embodiment has on participant evaluations of both the interaction and the robot.

1.2 Spatial Language-Based HRI Methodology

1.2.1 Spatial Language Communication in HRI

Service robots designed to interact with non-expert users will benefit greatly from spatial language communication and understanding. For example, consider the following instruction given to a household robot:

- (1) “Go to the kitchen”

If the user says (1), the robot should understand, in principle, what that means. That is, it should understand which task among those within its task/action repertoire the user is referring to. In this example, the robot may not know where the kitchen is located in the user’s specific home environment, but it should be able to understand that (1) expresses a command to physically move to a desired goal location that fits the description of the noun phrase “the kitchen”.

Spatial language plays an important role in instruction-based natural language communication. In (1), the spatial preposition “to” was used in the instruction of the task. The following sentence contrasts minimally in the preposition employed, using “away from” (compound preposition) instead of “to”:

(2) “Go away from the kitchen”

Yet, the meaning of the command specified by (2), or rather the implied goal task/action sequence, is completely different from (1), even though the verb and place noun are the same. The same holds for the spatial prepositions “around”, “through”, “behind”, etc. *Spatial relations*, both *dynamic* and *static*, expressed in language are often expressed by prepositions [Landau and Jackendoff, 1993]. Therefore, the ability for service robots to understand and differentiate between spatial prepositions in spoken language is crucial for their interaction with users to be successful.

Prepositions in English, as well as in many other languages, are identified as a closed class: there are only 80-100 prepositions [Landau and Jackendoff, 1993] (approximate count, as many have multiple meanings) and new words are not being added. The relatively small number of prepositions, combined with their extensive use in spatially-oriented natural language communication across domains, makes the construction of spatial primitives *a priori* based on prepositions for autonomous service robots not only feasible, but also intuitive and beneficial.

Spatial language understanding is especially relevant for interactive robot task learning and task modification. Continuing with the household robot example, the user might teach the robot the complex task “Clean up the room” through natural language and by specifying the subgoals of that task individually, each represented by its own spatial language instruction (e.g., “Put the clothes in the laundry basket”, “Stack the books on top of the desk in the right-hand corner”, “Put all toys under the bed”, etc.). In addition, user modification of known robot tasks can also readily be accomplished with spatial language. For example, the user might modify the task defined by (1) (i.e., robot movement to “the kitchen”) by providing spatial constraints, or rules, for the robot to obey during task execution, such as “Don’t go through the hallway,” or “Move along the wall.” These user-defined constraints do not change the meaning of the underlying

task, but allow the user to interactively modify the manner in which the robot executes the task in the specific instance. Finally, spatial language can be used to provide teacher feedback during task execution, to further correct or guide robot behavior. In the context of our example, as the robot is moving along the wall en route to the kitchen, the user may provide additional feedback by saying “Move a little further away from the wall,” or “Move close to the wall but stay on the paneled floor”. These examples illustrate the importance of spatial language in the instruction and teaching of service robots by non-expert users.

1.2.2 Spatial Language Framework and HRI User Study

This dissertation presents a novel methodology that allows service robots to interpret and follow spatial language instructions, with and without user-specified natural language constraints and/or unvoiced pragmatic constraints. The methodology is generalizable and can be applied across many human-robot interaction domains for a variety of assistive robot behaviors, including user-robot task instruction, teaching, modification, and guidance. In particular, this work contributes a general computational framework for the representation of *dynamic spatial relations* (DSRs), including a novel extension to the *semantic field* model of spatial prepositions, which enables the representation of path prepositions containing both local and global properties. The methodology also contributes a probabilistic approach in the inference of instruction semantics, and in the associated grounding of noun phrases utilizing the proposed computational field representation of spatial relations. The approach allows for robot motion planning and execution of multi-step instruction sequences in real-world continuous domains while providing robustness to sensor noise and environmental uncertainty.

Our service robot software architecture for HRI contains five system modules that enable the interpretation of natural language instructions, from speech or text-based

input, and translation into robot action execution. They are: the syntactic parser, noun phrase (NP) grounding, semantic interpretation, planning, and action modules.

The syntactic parser represents the entry point of our robot architecture, as it is responsible for parsing the user-given natural language instruction into a format that the remaining modules can interpret. The instructions are provided as text strings, either by a speech recognizer (e.g., [Nuance, 2013]) or keyboard-based input. Our system does not attempt to provide a solution for natural language processing in the general case, but instead focuses on well-formed English directives involving spatial language, for which we utilize a specialized grammar. After the syntax of the instruction has been determined, the parse tree is passed on to the grounding module which attempts to associate parsed NPs with known objects in the world. If it is successful, all observations are then passed on to the semantic interpreter for final instruction meaning association.

The semantic interpretation module utilizes a Bayesian approach and infers the semantics of the given instruction probabilistically using a database of learned mappings from input observations to instruction meanings. The input includes: the verb and preposition used in the instruction sentence, and the associated groundings for the specified figure and reference objects as determined by the NP grounding module. The resulting semantic output of the module includes: the command type, the DSR type, and the static spatial relation (if available). The command type is domain-specific, and may include commands such as: robot movement, speech output, learned tasks, etc. While the output specification was designed to represent the instruction of spatial tasks, the inference procedure utilized is general and can easily be modified or expanded to accommodate the requirements of the specific application domain, including the inference of non-spatial tasks.

Once the semantic interpreter has inferred the instruction semantics, the planning module attempts to find a solution for the robot given these command specifications, after which the solution is passed on to the action module for robot task execution.

Evaluation results of the individual system modules, and of the spatial language understanding and HRI framework as a whole, were obtained from 2D/3D simulations of a mobile service robot operating within both manually constructed household environments and real-world environments from robot-generated SLAM maps; the results presented validate the effectiveness of the methodology employed.

In addition, a user study was conducted with older adults to evaluate the effectiveness and feasibility of our spatial language interpretation framework with target users, and to collect data on the types of phrases, responses, and format of natural language instructions given by end users to help inform possible modifications to the spatial language grammar and/or interpretation module of our framework. In particular, the study served to evaluate our approach for interpreting instruction sequences obtained from unconstrained natural language input, and our method of resolving anaphoric reference expressions encountered within user discourse.

1.3 Dissertation Contributions

This dissertation addresses specific challenges in socially assistive and service robotics-based robot-guided and user-guided interactions. In robot-guided interaction, it presents the approach, design methodology, and implementation details of a novel SAR system that aims to motivate and engage elderly users in simple physical exercise. In the area of user-guided interaction, it presents an approach for enabling autonomous service robots to follow natural language commands from non-expert users, including under user specified constraints, with a particular focus on spatial language understanding.

The following are the main contributions of this dissertation:

1. *A set of five general design principles for socially assistive systems* to aid users in a healthcare task by influencing their intrinsic motivation to engage in the therapeutic intervention.
2. *The approach, design methodology, and implementation details of a fully autonomous SAR system* designed to engage elderly users to perform physical exercise while providing real-time feedback, guidance, and encouragement.
3. *The first comprehensive physical vs. virtual robot embodiment comparison study* conducted to observe the effects of embodiment within a SAR-guided healthcare scenario, wherein the SAR agent serves as both an instructor and an active participant in the healthcare task with target users, the results of which demonstrated the positive effect of physical embodiment in a fluid, real-time human-robot interaction based healthcare intervention with elderly users.
4. *A novel representation for DSRs with global properties and DSRs with local properties* that extends the semantic field model of spatial prepositions and facilitates probabilistic reasoning over paths that can be applied for both path classification and path generation scenarios.
5. *The design, system modules, and implementation details of a novel methodology for spatial language-based HRI* that allows service robots to interpret and follow spatial language instructions, with and without explicit figure objects, and with and without user-specified natural language constraints and/or unvoiced pragmatic constraints, including a probabilistic approach for the grounding and interpretation of spatial language instructions.

The following are the secondary contributions of this dissertation:

1. *Validation of the SAR approach* and its effectiveness in gaining user acceptance and motivating elderly users to engage in consistent physical exercise in three separate user studies conducted with the target healthcare population, as well as analyzing the effect of praise and relational discourse, user autonomy, and agent embodiment within the interaction.
2. *Example representations for the DSRs of “to”, “through”, and “around”*, including implementation details of the path generation procedures utilized by our system for these three DSRs, and discussion of relevant pragmatic constraints together with planning methods developed to address these constraints in multi-step robot execution planning of instruction sequences.
3. *A novel probabilistic algorithm for the automatic extraction of contextually and semantically valid instruction sequences from unconstrained spatial language discourse*, including the design and implementation details of a procedure for reference resolution of anaphoric expressions encountered within the user discourse.
4. *Validation of the spatial language-based HRI methodology* from simulation testing in both 2D and 3D, implementation and testing on two separate physical mobile robot platforms (PR2 robot and Bandit humanoid), including in a user study conducted with older adults to evaluate the effectiveness of the approach in interpreting instructions in discourse from target users.

1.4 Outline

The remainder of this dissertation document is organized as follows:

- Chapter 2 discusses related work in the area of service robotics, socially assistive robotics and social agent coaches, as well as work investigating the representation and use of spatial language within a framework for HRI.
- Chapter 3 presents our SAR design principles for socially assistive system-based therapeutic interventions which incorporate insights from psychology research on intrinsic motivation.
- Chapter 4 presents the approach, design methodology, software architecture, implementation details, and interaction scenario of our SAR exercise coach for older adults.
- Chapter 5 describes two user studies conducted with older adults investigating user intrinsic motivation in the SAR exercise system, including the effect of praise and relational discourse in the system, and the role of user autonomy and choice within the interaction.
- Chapter 6 presents our embodiment and SAR evaluation user study conducted with older adults to study the effect of embodiment in the SAR exercise system by comparing user evaluations of similar physically and virtually embodied SAR exercise coaches, in addition to evaluating the overall SAR system with target users.
- Chapter 7 describes the design, system modules, and implementation details of our methodology for spatial language-based HRI that allows service robots to interpret and follow spatial language instructions, with and without user-specified natural language constraints and/or unvoiced pragmatic constraints.

- Chapter 8 presents evaluation results of the individual system modules, and spatial language-based HRI framework as a whole, that were obtained from 2D and 3D simulations of both manually constructed household environments and real-world environments from robot-generated SLAM maps, to validate the effectiveness of the methodology employed.
- Chapter 9 proposes a novel probabilistic algorithm for the automatic extraction of contextually and semantically valid instruction sequences from unconstrained spatial language discourse, and presents the design and implementation details of a procedure for reference resolution of anaphoric expressions encountered within the user discourse.
- Chapter 10 presents the user study we designed and conducted with older adult participants to evaluate the effectiveness and feasibility of our spatial language interpretation framework with end users.
- Chapter 11 summarizes the dissertation.

Chapter 2

Background and Related Work

This chapter surveys prior research related to socially assistive robotics and social agent coaches, as well as work investigating the representation and use of spatial language within a framework for HRI. The chapter begins with a brief overview of service robotics.

2.1 Service Robotics

2.1.1 Overview of Service Robotics

The field of service robotics is extensive and covers a wide range of subfields in robotics; with robots operating in a variety of different environments, all sharing the same goal of providing assistance to human users, typically through some form of HRI, be it social, physical, or otherwise (e.g., GUI/joystick control). Subfields and areas within service robotics include robots developed for rehabilitation [Burgar et al., 2002; Harwin et al., 1988; Kahn et al., 2001], aerial robots [Shen et al., 2012], underwater robots [Pereira et al., 2013], factory robots [Tellex et al., 2011], office robots [Veloso et al., 2012], household robots [Choi et al., 2009; Rybski et al., 2008] and general-purpose service robots capable of learning tasks through HRI and from demonstration [Koenig

et al., 2010; Nicolescu and Matarić, 2005; Pardowitz et al., 2007; Rybski et al., 2008] among others. In this work we focus on the development of socially assistive robots for therapeutic contexts, discussed in Section 2.2, and general-purpose service robots capable of spatial language understanding and instruction following through natural language HRI with non-expert users, discussed in Section 2.3.

2.1.2 Service Robots for Older Adults

The literature that addresses assistive robotics intended for and evaluated by the elderly is limited but growing. Representative work in the field includes robots that focus on providing assistance for functional needs, such as mobility aids and navigational guides. Dubowsky et al. [2000] developed a robotic cane/walker device designed to help individuals by functioning as a mobility aid that provides physical support for walking as well as guidance and health monitoring of a user’s basic vital signs. Montemerlo et al. [2002] designed and pilot-tested a robot that escorts elderly individuals in an assisted living facility, reminds them of their scheduled appointments, and provides informational content such as weather forecasts.

2.2 Socially Assistive Robotics

2.2.1 SAR for Older Adults

Researchers have investigated the use of robots to help address the social and emotional needs of the elderly, including reducing depression and increasing social interaction with peers. Wada et al. [2002] studied the psychological effects of a seal robot, Paro, which was used to engage seniors at a day service center. The study found that Paro, always accompanied by a human handler, was able to consistently improve the moods of elderly participants who spent time petting and engaging with it over the course of a 6-week

period. Kidd et al. [2006] used Paro in another study that found it to be useful as a catalyst for social interaction. They observed that seniors who participated with the robot in a group were more likely to interact socially with each other when the robot was present and powered on, rather than when it was powered off or absent.

Perhaps the robotic system for the elderly most related to our SAR exercise system is the work of Matsusaka et al. [2009], who developed an exercise demonstrator robot, TAIZO, to aid human demonstrators during simple arm exercises with a training group. However, this robot was not autonomous, as it was controlled via key input or voice by the lead human demonstrator, and it did not have sensors by which to perceive the users; hence, it did not provide any real-time feedback, active guidance, or personalized training.

2.2.2 Social Agent Coaches

Social agents that aim to assist individuals in health-related tasks such as physical exercise have also been developed in the human-computer interaction (HCI) community. Bickmore and Picard [2005] developed a computer-based virtual relational agent that served as a daily exercise advisor by engaging the user in conversation and providing educational information about walking for exercise, asking about the user's daily activity levels, tracking user progress over time while giving feedback, and engaging the user in relational dialogue. Kidd and Breazeal [2008] developed a table-top robot to serve as a daily weight-loss advisor; it engaged users through a touch-screen interface, tracked user progress and the state of the user-robot relationship over time, and was tested in a six-week field study with participants at home. French et al. [2008] designed and explored the use of a virtual coach to assist manual wheelchair drivers by providing advice and guidance to help users avoid hazardous forms of locomotion.

These systems are similar to our SAR exercise system in the manner in which they provide feedback (from a social agent) and, with the exception of French’s work, in the activity being monitored (physical exercise). However, our system is clearly distinct in that our robot agent not only provides active guidance, feedback, and task monitoring, but is also directly responsible for instructing and steering the task as well. Hence, our agent is both an administrator and an active participant in the health-related activity, resulting in a unique characteristic of the system: The social interaction between the robot and user is not only useful for maintaining user engagement and influencing intrinsic motivation, but is also necessary to achieving the physical exercise task.

2.2.3 The Effect of Embodiment

Previous studies investigating the role of embodiment within the context of human-agent interaction have demonstrated the potential positive effects that physical embodiment can have on people’s level of engagement and overall perception of the agents with which they are interacting. Wainer et al. [2006, 2007] showed that healthy adult participants engaging in a physical/cognitive task, a Towers of Hanoi table-top game, reported a strong preference for a physically embodied SAR system over similar video-only agents in terms of appeal, perceptiveness, watchfulness, helpfulness, and enjoyableness. Powers et al. [2007] compared interactions between robots and similar computer-simulated agents that engaged participants in a conversation about basic health habits, and found that participants rated the robots as more helpful, more lifelike, and possessing more positive personality traits than the computer-based agents. Bartneck [2003] conducted a study comparing the effectiveness of an emotionally expressive robot, eMuu, with its screen character version in engaging users in a simple negotiation task, and found that participants exerted more effort and received higher task scores when interacting with

the physical eMuu than with the simulated eMuu. Jung and Lee [2004] also demonstrated the positive effects of physical embodiment in relation to interactions with both a Sony Aibo robot and an anthropomorphic dancing robot, April. Bainbridge et al. [2011] found that users in a book-moving task were more likely to fulfill an unusual request and afford more personal space to the agent when interacting with a physically present robot than when interacting with a live video feed of the same robot on a computer screen. Kidd and Breazeal [2008] compared a robotic weight-loss coach (a touch screen with a physical head capable of looking at and speaking to the user) to a similar touch-screen-only device and found that participants interacting with the robotic coach chose to continue with the weight-loss program for twice as long as those interacting with the computer-only device.

While studies such as those mentioned above have previously investigated the effect of physical embodiment in human-agent interaction, most have recruited a participant pool consisting primarily of young adults. However, older adults often respond differently to technology than young adults, as studies have shown [Balakrishnan and Yeow, 2007; Kang and Yoon, 2008], and thus the observed effects do not necessarily generalize across the age span.

Embodiment studies that have targeted the elderly population include the work of Heerink et al. [2010], who investigated the acceptance of assistive social agents by older adults. While their study was similar to our work, the robot used in their evaluation was a table-top robot (the iCat), and was either controlled via a human operator during interaction with elderly users (Wizard of Oz study), or interacted with users through a touch-screen interface. Furthermore, the interaction consisted primarily of short informational or utility interactions (e.g., medication/agenda reminders, weather forecast, companionship), lasting about 5 minutes and often involving only a single session. Lastly, their work did not explore the agent’s role in actively motivating the user

to engage in the task, and instead focused solely on the effect of embodiment. In contrast, our SAR system was designed to engage elderly users in fluid, highly interactive exercise sessions, completely autonomously, while providing active feedback, motivation, and guidance on the task. Also, our user study spanned multiple sessions, each lasting 20 minutes, for the dual purpose of system evaluation and comparison between two different SAR coaching embodiments.

The interactivity of the sessions is important because, as we discuss in the next chapter, the fluidity of the interaction can have a positive influence on the user’s intrinsic motivation to engage in the task and can thereby increase his or her enjoyment level during interaction. It is thus interesting to investigate whether the positive effects introduced by the interaction characteristics might somehow alleviate any potential negative effects resulting from embodiment characteristics. Our study on embodiment (discussed in Chapter 6) addresses this possibility in capturing user evaluations of the system between study groups.

To the best of our knowledge, our embodiment user study is the first to comprehensively demonstrate the positive effect of physical embodiment in a SAR-guided healthcare scenario, wherein the SAR agent serves as both an instructor and active participant in the healthcare task with target users.

2.3 Spatial Language-Based HRI

2.3.1 Modeling Static Spatial Relations

Previous work that has investigated the use of spatial prepositions, and spatial language in general, in an HRI framework, includes the work of Skubic et al. [2004], who demonstrated a robot capable of understanding static spatial relations in natural language instruction. Sandamirskaya et al. [2010] investigated the use of Dynamic Neural

Fields theory in a static spatial language architecture for use in human-robot cooperation tasks on a tabletop workspace. Similarly, the use of computational fields for static relations was implemented in a visually situated dialogue system by Kelleher and Costello [2009]. These works all implemented pre-defined notions of spatial relations, however, researchers have also investigated learning these types of static spatial relations automatically from training data both on- and offline (e.g., [Chao et al., 2011; Hawes et al., 2012; Mohan et al., 2012; Roy, 2002]). Our work aims to extend upon this related work by encoding not only static spatial relations for natural language instruction understanding, but also dynamic spatial relations involving paths, as discussed in the next section.

2.3.2 Modeling Dynamic Spatial Relations

In the context of natural language robot instruction, however, the use of dynamic spatial relations has in fact been explored by recent work. Tellex et al. [2011] developed a probabilistic graphical model to infer task/actions for execution by a forklift robot from natural language commands. Kollar et al. [2010] developed a Bayesian framework for interpreting route directions on a mobile robot, using learned models of dynamic spatial relations such as “past” and “through” from schematic training data. In both of these works there was no explicit definition of the spatial relations used, static or otherwise, and instead they were learned from labeled training data. However, these approaches typically require the programmer to provide an extensive training data set of natural language input for each new application context, without taking advantage of the domain-independent nature of spatial prepositions. Our proposed approach develops novel, pre-defined templates for spatial relations, static and dynamic, that facilitate use and understanding across domains, and whose computational representations enable guided agent execution planning.

2.3.3 Robot Control Language

Researchers have also explored mapping natural language instructions into a formal robot control language using a variety of types of parsers, including those that were constructed manually [Kress-Gazit et al., 2008; Rybski et al., 2008], learned from training data [Matuszek et al., 2012], and learned iteratively through interaction [Cantrell et al., 2011]. Of these, the work of Rybski et al. [2008] and Matuszek et al. [2012] rely on pre-defined agent behaviors as primitives, as opposed to spatial relations, which hinders, if not prohibits, the ability for a user to provide feedback modifications and/or constraints regarding agent execution of a specific primitive behavior. The work of Kress-Gazit et al. [2008] and Cantrell et al. [2011] leave the specification of primitives, from which new behaviors are learned, largely as an open research problem. However, the parsers used in their systems map words to meanings based on dictionary-based rules. Our methodology employs domain-generalizable spatial relations as primitives, and probabilistic reasoning for the grounding and semantic interpretation of phrases, thereby allowing for context-based instruction understanding and user-feedback modifiable robot execution paths.

2.4 Summary

This chapter has provided an overview of the related work in the research areas of service robotics for both robot-guided and user-guided interactions, focusing in particular on those most closely related to the work presented in this dissertation. Specifically, in the area of robot-guided interaction, the chapter included discussion of related work in the area of SAR for older adults, social agent coaches, and studies investigating the effect of embodiment in HRI. In the area of user-guided interaction, the chapter provided discussion of related approaches to representing and interpreting spatial language in an

HRI framework. In both discussions, the differences between our approach, to SAR-based therapy and spatial language-based HRI, and the existing approaches listed were detailed to highlight the contributions of the work presented in this dissertation.

Chapter 3

Socially Assistive Robot Design Principles

Our approach to designing SAR systems to help address the physical exercise needs of the elderly population is motivated by two basic axioms indicating the essential qualities that the SAR coach must possess: 1) the ability to influence the user's intrinsic motivation to perform the task, and 2) the ability to personalize the social interaction to maintain user engagement in the task and build trust in the task-based human-robot relationship. Following the above axioms, we developed five design principles for SAR coaches; all are general and can be applied to a broad range of SAR-based therapeutic interventions.

3.1 Design Principles for a SAR Coach

The five design principles developed stated that the SAR coach should be: 1) motivating, 2) fluid and highly interactive, 3) personable, 4) intelligent, and 5) task-driven. The following discussion elaborates on the importance of each of these qualities in the context

of providing healthcare interventions and details of how each was incorporated into our SAR exercise coach.

3.1.1 Motivating

The coaching style and interaction methodology of our SAR exercise system was guided by psychology research in the area of intrinsic motivation. Intrinsic motivation, in contrast to extrinsic motivation, which is driven by external rewards, comes from within an individual and is based on pleasure derived from engaging in an activity. Specifically, our aim was for the robot to be capable of increasing the user’s intrinsic motivation to perform the therapeutic task (physical exercise, in this case). Intrinsic motivation has been shown to be more effective than extrinsic motivation in achieving long-term user task compliance and behavior change [Dienstbier and Leak, 1976], which is the ultimate goal for any health-related intervention, technology-based or otherwise. The motivational techniques utilized by our system to accomplish this aim were derived from Csikszentmihalyi’s theory of flow [Csikszentmihalyi, 1975], which Deci and Ryan [1985, p. 29] describe as asserting that people are “intrinsically motivated under conditions of optimal challenge.” Toward this end, we focused on providing a variety of challenging exercise games of varying degrees of difficulty. We also focused on the fluidity of the interaction, discussed below, as well as on alternating the games at a regular pace to prevent user boredom and/or frustration.

Our design also included additional motivational techniques besides those based on flow theory. For instance, we incorporated indirect competition into the system design by having the robot periodically report the user’s high score during each of the exercise games. Indirect competition, wherein the user is challenged to compete against an ideal outcome, has been shown to increase user enjoyment in an otherwise non-competitive task [Weinberg and Ragan, 1979]. In addition, self-determination and user

autonomy, which have also been shown to increase intrinsic motivation [Fisher, 1978], were implemented in our system by giving the user control over the exercise routine in one of the exercise games, as discussed in more detail in Section 4.3.3.

As indicated by the first basic axiom above, influencing the user’s motivation was a key component in our SAR agent design; while our first design principle stresses the importance of motivation in the task scenario, all of our design principles were developed to increase user intrinsic motivation during interaction in one way or another.

3.1.2 Fluid and Highly Interactive

A primary goal of our coaching approach was to provide a fluid interaction, which required the robot to both perceive the user’s activity and provide active feedback and guidance in real time, all with the aim of maintaining user engagement in the task. According to Csikszentmihalyi [1993], for any task to achieve a state of flow—or maximal enjoyment—in the user, it must establish a clear set of goals, combined with immediate and appropriate feedback. Toward this end, we developed a real-time vision algorithm to recognize user arm movements (detailed in Section 4.5), as well as coaching behaviors to continually provide the user with appropriate feedback to achieve the goals of the exercise games (discussed in Section 4.6). The result is an exercise coach that is highly interactive and responsive to the user, thus promoting a state of flow.

3.1.3 Personable

The social interaction between the user and the robot is just as important as the task interaction for achieving success in healthcare interventions. Social interaction is the primary means of relationship building, including in therapeutic scenarios. Many social intricacies contribute to the foundation of a meaningful relationship, both in human-computer interaction and in human-robot interaction. These factors include

empathy, humor, references to mutual knowledge, continuity behaviors, politeness, and trust, among others [Bickmore and Picard, 2005]. We place great importance on these relationship-building tools; therefore, we integrated each, in one form or another, into the social interaction component of our robot exercise coach. Primary examples include referring to the user by name; giving praise upon successful user completion of exercise gestures, which has also been shown to increase intrinsic motivation [Vallerand, 1983; Vallerand and Reid, 1984]; displaying empathy after failed gesture attempts; and demonstrating continuity to establish an ongoing relationship with the user. In our system, the robot always uses the user’s name at first greeting, and also when bidding farewell at the end of a session. Having the robot refer to the user by name, along with providing direct feedback specific to the individual user’s performance level and performance history during the games, was an important part of personalizing the interaction. Our SAR exercise system introduced continuity (see Section 4.6) by having the robot refer to previous sessions with the user upon introduction, reference planned future sessions at the end of interaction, and refer to past user exercise performance, such as when reporting previous high scores.

3.1.4 Intelligent

Trust is a key component to the success of any care-provider/user relationship, and one that is closely linked to the intelligence/helpfulness of the care provider as perceived by the user [Tway, 1994]. Tway defines trust as a construct determined by three components: perception of competence, perception of intentions, and capacity for trusting. In the context of SAR-based therapy, the perceived intelligence/competence of the SAR agent by the user is necessarily correlated with the agent’s demonstrated abilities in the therapeutic task (e.g., performance capability, evaluation accuracy, etc.). In our SAR exercise scenario, the effectiveness of the robot’s user activity recognition procedure was

paramount in shaping the user’s perception of the robot’s intelligence. (The algorithm is discussed in Section 4.5.) In particular, if the user were to perceive the robot as slow or ineffective in evaluating performance, this could lead to a decrease in the user’s trust and perception of the robot’s value/usefulness in helping to accomplish the desired healthcare goals, according to Tway. This aspect once again stresses the importance of accurate, real-time sensing in the human-robot interaction.

Intelligent robot social behaviors can also have a positive impact on user engagement. For instance, it has been shown that repetitive discourse tends to have a negative impact on user motivation, whereas increased variability, in discourse and general behaviors, tends to enhance task-based user engagement [Bickmore et al., 2010]. Therefore, in our SAR system, we placed special attention on ensuring variety in the robot’s utterances to minimize the perceived repetitiveness of the robot’s verbal instructions/comments. Toward this end, the robot always draws from a list of phrases that emphasize the same point when speaking to the user, choosing one randomly at run time. For example, there are more than ten different ways in which the robot can praise the user after he or she completes an exercise (e.g., “Awesome!”, “Nice job!”, “Fantastic!”). Variability is also introduced in discourse during game transitions, game introductions, and upon farewell. General interaction behaviors were also designed to reduce predictability; for example, by ensuring that the order of games, game timing, and session schedules were different for each interaction. Further information regarding the methods employed to add variability to our SAR behaviors and interaction sessions is provided in Section 4.6.

3.1.5 Task-Driven

Perhaps the most important property of the interaction is that it be consistent in working to achieve the goals of the healthcare task, or in our case, in motivating exercise

performance to achieve desired overall health benefits. As mentioned above, it is important for the task to be intrinsically motivating and enjoyable and for the coach be personable and intelligent. However, the consistency of the interaction towards accomplishing the healthcare task is fundamental for the user to trust that the system can help to obtain beneficial health outcomes, as this consistency influences the user's perception of the intentions of the SAR agent. Agent intentions, according to Tway [1994], must be seen as mutually beneficial. Without this consistency, users may perceive the robot as simply entertaining, rather than helpful.

Furthermore, it is important that the tasks not only be healthcare-driven, but that they also be successful in achieving the desired therapeutic behavior; in the case of our SAR exercise coach, this means the interaction must elicit consistent physical exercise among the users (measurable through objective quantitative metrics), a result that was achieved in our user studies with older adults as detailed in Chapters 5 and 6.

3.2 Summary

This chapter presented five SAR design principles that were developed in an attempt to capture and address the key components necessary for successful SAR-based therapy. Each of the design principles was developed based on findings from research in psychology, specifically in the areas of intrinsic motivation and trust, and from research in HCI and HRI, specifically regarding the effects of relational agent behaviors on user engagement. Although some previous work has addressed the design of social agents for obtaining long-term engagement and trust, most has focused on agent guidance in task scenarios in which the agent does not monitor or instruct the activity itself (e.g., [Bickmore and Picard, 2005; Kidd and Breazeal, 2008]).

The five SAR design principles presented in this chapter can be applied to a variety of different SAR-based therapeutic domains, including those wherein the SAR agent

provides real-time instruction, activity monitoring, active guidance and feedback, task participation, and continuity behaviors spanning multiple sessions of interaction. The following chapter provides examples of how to incorporate these design principles into practice with a SAR agent, and elaborates on the implementation details of our SAR exercise system developed to engage older adults in physical exercise.

Chapter 4

Socially Assistive Robot Exercise Coach

This chapter presents the design and implementation details of our SAR exercise coach that was developed to engage older adults in physical exercise, and which incorporated the five SAR design principles, discussed in the previous chapter, both in the system and interaction design. The topics of discussion include: the human-robot interaction scenario, the robot platform used, the software architecture and the primary software modules of robot perception and behavior.

4.1 Interaction Scenario

The purpose of the socially assistive robot coach is to monitor, instruct, evaluate, and encourage users to perform simple seated physical exercises. During the exercise sessions, the robot asks the user to perform seated arm gesture exercises. This type of seated exercise, called “chair exercise” or “chair aerobics,” is commonly practiced in senior living facilities. Chair exercises are highly regarded for their accessibility to those with low mobility, for their safety as they reduce the possibility of injury due to falling from improper balance, and for their health benefits, such as improved flexibility, muscle

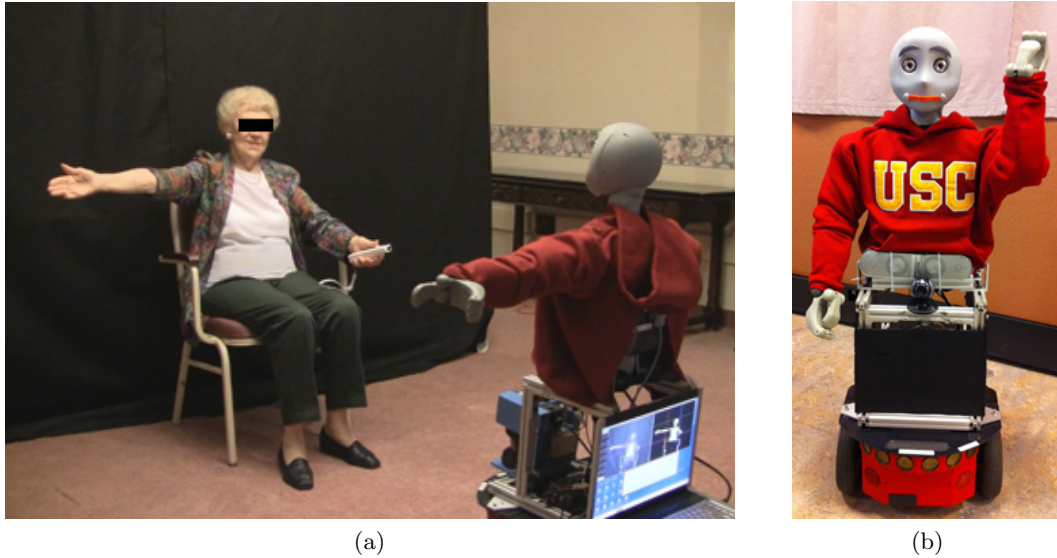


Figure 4.1: (a) The setup for the one-on-one interaction between user and robot coach; (b) Bandit humanoid torso robot.

strength, ability to perform everyday tasks, and even memory recall [Baum et al., 2003; Dawe and Moore-Orr, 1995; McMurdo and Rennie, 1993; Thomas and Hageman, 2003].

The one-on-one interaction scenario consists of the user sitting in a chair facing the robot. In this system, predating the availability of the Microsoft Kinect, the range of the robot’s arm motion in the exercises is planar and restricted to the sides of the body in order to maximize the accuracy of the robot’s real-time visual detection of the user’s arms. A black curtain is used as a backdrop, also to facilitate the robot’s visual perception of the user’s arm movements. An example of the exercise setup and human-robot interaction is shown in Figure 4.1 (a).

The system allows the user to communicate with the robot through the popular Nintendo Wiimote wireless Bluetooth button control interface. There are two buttons available for the user to respond to prompts from the robot, labeled “Yes” and “No,” and

one button for the user to request a rest break at any time during the interaction. The Wiimote is not used for motion sensing of any kind; its role is solely for communication with the robot through button presses. It is also important to note that the robot conducts the exercise sessions, evaluates user performance, and gives the user real-time feedback all completely autonomously, without human operator intervention at any time during the exercise sessions.

4.2 Robot Platform

To address the role of the robot’s physical embodiment, we used Bandit, a biomimetic anthropomorphic robot platform that consists of a humanoid torso (developed with BlueSky Robotics) mounted on a MobileRobots Pioneer 2DX mobile base. The torso contains 19 controllable degrees of freedom (DOF): 6 DOF arms (x2), 1 DOF gripping hands (x2), 2 DOF pan/tilt neck, 1 DOF expressive eyebrows, and a 2 DOF expressive mouth. The robot is shown in Figure 4.1 (b).

A standard USB camera is located at the waist of the robot, and used to capture the user’s arm movements during the exercise interaction, allowing the robot to provide appropriate performance feedback to the user.

The robot’s speech is generated by the commercially available NeoSpeech text-to-speech engine [NeoSpeech, 2009] and a speaker on the robot outputs the synthesized voice to the user. The robot’s lip movements are synchronized with the robot’s speech so that the lips open at the start and close at the end of spoken utterances.

4.3 Exercise Games

Four exercise games are available in our system: the Workout game, the Sequence game, the Imitation game, and the Memory game. Each game presents a different

level of challenge and necessary skill set, in accordance with our system methodology of varying game types for increased user enjoyment and to promote intrinsic motivation in the therapeutic interaction. Following is a brief description of each game.

4.3.1 Workout Game

In the Workout game, the robot fills the role of a traditional exercise instructor by demonstrating the arm exercises with its own arms and asking the user to imitate. The robot shows only one exercise gesture at a time (involving one or both arms). Upon successful completion by the user, the robot generates a different gesture and the process repeats. The robot gives the user feedback in real time, providing corrections when appropriate (e.g., “Raise your left arm and lower your right arm” or “Bend your left forearm inward a little”), and praise in response to each successful imitation (e.g., “Great job!” or “Now you’ve got the hang of it.”). This game has the fastest pace of the four exercise games, as the users generally complete the requested gestures quickly.

4.3.2 Sequence Game

The Sequence game is similar to the Workout game in that the robot demonstrates arm exercises for the user to repeat. However, instead of showing each gesture for the user to perform only once, the robot demonstrates two gestures for the user to repeat in sequence for three repetitions (resulting in six gesture completions per sequence). The robot keeps verbal count of the number of iterations of the sequence performed in order to guide the user, while providing feedback throughout (e.g., after the user completes the first pair, the robot says “One,” and upon completion of the second pair, it says “Two”). This game requires the user to remember the gesture pair while completing the sequence, in addition to promoting periodic movements from the users. While slower-paced than the Workout game, these periodic movements are physically challenging and

can cause the users to physically exert themselves more than in any of the other three games.

4.3.3 Imitation Game

In the Imitation game, the roles of the user and robot are reversed relative to the Workout game: The user becomes the exercise instructor, showing the robot what to do. The robot encourages the user to create his/her own arm gesture exercises and imitates user movements in real time. Since the roles of the interaction are reversed, with the robot relinquishing control of the exercise routine to user, the robot does not provide instructive feedback on the exercises. However, the robot does continue to speak and engage the user by means of encouragement, general commentary, humor, and movement prompts if necessary. As the robot imitates the user's movements in real time, the interaction in this game is the most fluid and interactive of all the exercise games. Furthermore, this is the only game designed to promote user autonomy and creativity in the exercise task. Both of these characteristics are beneficial for increasing the user's intrinsic motivation to engage in the task.

4.3.4 Memory Game

The goal of the Memory game is for the user to try and memorize an ever-longer sequence of arm gesture poses, and thus compete against his/her own high score. The sequence is determined at the start of the game and does not change for the duration of the game. The arm gesture poses used for each position in the sequence are chosen at random at run time, and there is no inherent limit to the sequence length, thereby making the game challenging for users at all skill levels. Each time the user successfully memorizes and performs all shown gestures without help from the robot, the robot shows the user two additional gestures to add to the sequence, and hence the game progresses in

difficulty. The robot helps the user to keep track of the sequence by counting along with each correct gesture and reminding the user of the correct poses, while demonstrating empathy upon user failure (e.g., “Oh, that’s too bad! Here is gesture five again.”). This game is the most cognitively challenging of all the exercise games, and primarily serves to encourage competition with the user’s past performance.

4.4 Software Architecture

The system architecture is comprised of six independent software modules: vision and world model, speech, user communication, behaviors, robot action, and database management. A diagram of the system architecture, showing the connections among the system modules, is provided in Figure 4.2. The following briefly describes the role of each system module:

Vision and World Model Module. This module is responsible for providing information regarding the state of the user to the behavior module for the robot to make task-based decisions during interaction; for example, by providing spoken or demonstrative feedback, prompts for user movement, or by imitating user arm poses. The input for the visual user activity recognition procedure is a monocular USB camera, with a frame resolution of 640x480 pixels. The presence of the user and the locations of the face, hands, and arm angles of the user are captured by the vision process and stored in the world model for subsequent use. The details of the vision algorithm for user state recognition are provided in Section 4.5.

Speech Module. The speech module is responsible for translating the speech requests from the behaviors, provided as text strings, into spoken words (synthesized audio). This module uses the commercially available NeoSpeech text-to-speech (TTS) engine to

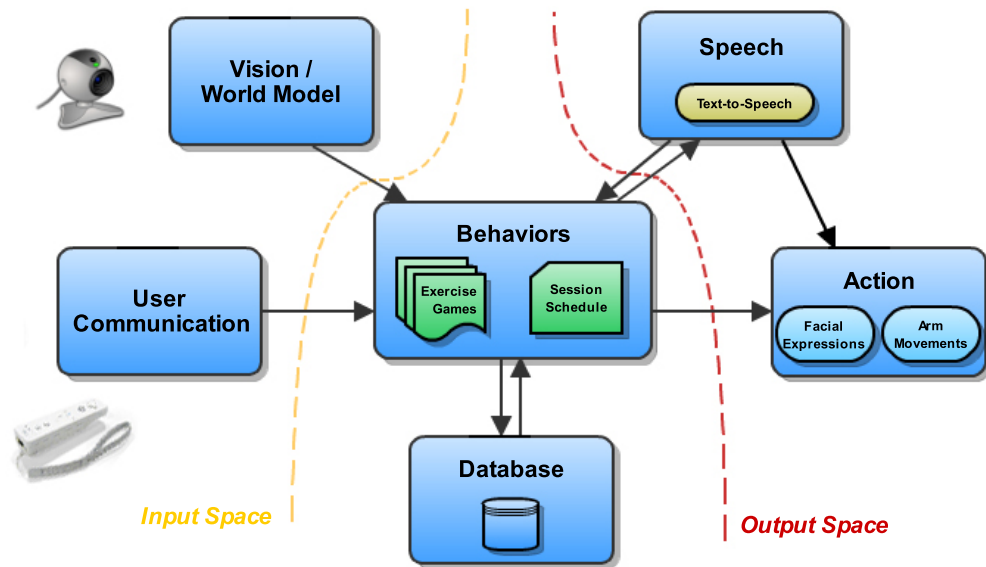


Figure 4.2: Diagram of system architecture, including the six system modules: vision and world model, speech, user communication, behaviors, action output, and database management.

synthesize the given text into communicative speech. The speech module plays the synthesized speech through the robot’s speakers, and is also responsible for synchronizing the lip movements of the robot to match the start and end of spoken utterances. To accomplish this synchronization, the module analyzes the wave file produced by the TTS engine and locates the key start and end points of silence, which correspond to pauses in the speech; for example, due to commas or periods in the input text. Then, at the appropriate times during speech playback, the speech module sends lip movement commands to the action module to simulate the opening and closing of the robot’s mouth during speech. The lip synchronization procedure is used to enhance the naturalness of the interaction to promote user engagement with the robot. The alternative of keeping the robot’s lips closed during speech might appear incongruent to the user and thus have

a disengaging effect. The speech module also communicates speech state information to the behavior module, such as whether or not it is currently speaking, so that the behavior finite state machines can make transitions accordingly (e.g., by waiting for the robot to finish giving a praise comment before demonstrating the next exercise gesture).

User Communication Module. The user communication module is responsible for receiving direct input from the user and relaying the communication attempt to the behavior module. The input in our SAR system is via a Nintendo Wiimote wireless remote control device, which is capable of sending button presses wirelessly via Bluetooth. The simplicity of the input method allows the user to request breaks during interaction, respond to yes/no questions or prompts made by the robot, and also to report key information during the exercise games (e.g., notifying the robot of a gesture completion attempt during the Memory game). Each communication type is represented by a different button specifically labeled on the remote control. This method enables simple two-way communication between the robot and user, is easy to use, and avoids the sensing difficulties often found with alternative communication input methods (e.g., speech recognition).

Behavior Module. The behavior module is responsible for producing all of the coaching behaviors of the system, including steering the interaction and exercise games, providing verbal feedback, demonstrating actions, responding to user input, monitoring user activity and progress in the task, and recording important information to the database. The behavior module therefore communicates with each of the six system modules and represents the system's main loop. The behavior module also keeps track of important session variables, including session time and total duration, current game time, and session schedule.

Action Module. The action module is responsible for sending motor commands to the robot to execute. These include facial expressions (e.g., moving eyebrows and lips) and arm movements. The module was designed to promote robot/agent platform independence, and as a result, all action requests made by connected system modules are only translated to the robot platform-dependent motor commands in the action module, thereby abstracting this knowledge away from all other system modules. This hierarchy allows for seamless integration of different platforms for the social agent, without the need for modification of any of the remaining software modules. In the user study presented in Chapter 6, we explored the use of both a physical humanoid robot and a virtual simulation of the same robot for system evaluation, thus demonstrating the interoperability of the system design.

Database Module. The database module is responsible for storing important user-related information regarding task performance and progress, in addition to session history. Examples include the date and time of previous sessions, exercise performance in each of the exercise games, number of breaks taken, high scores, etc. This information is available for retrieval during interaction by the behavior module, which uses the information to implement continuity behaviors and for reporting user performance and progress. The database information is also useful for obtaining quantitative metrics of user performance for post-session/user study system evaluation.

4.5 Visual User Activity Recognition Algorithm

In order to monitor user performance and provide accurate feedback during the exercise routines, the robot must be able to recognize the user's arm gestures. To accomplish this, we developed a vision algorithm that recognizes the user's arm gestures/poses in real time, with minimal constraints on the surrounding environment and the user.

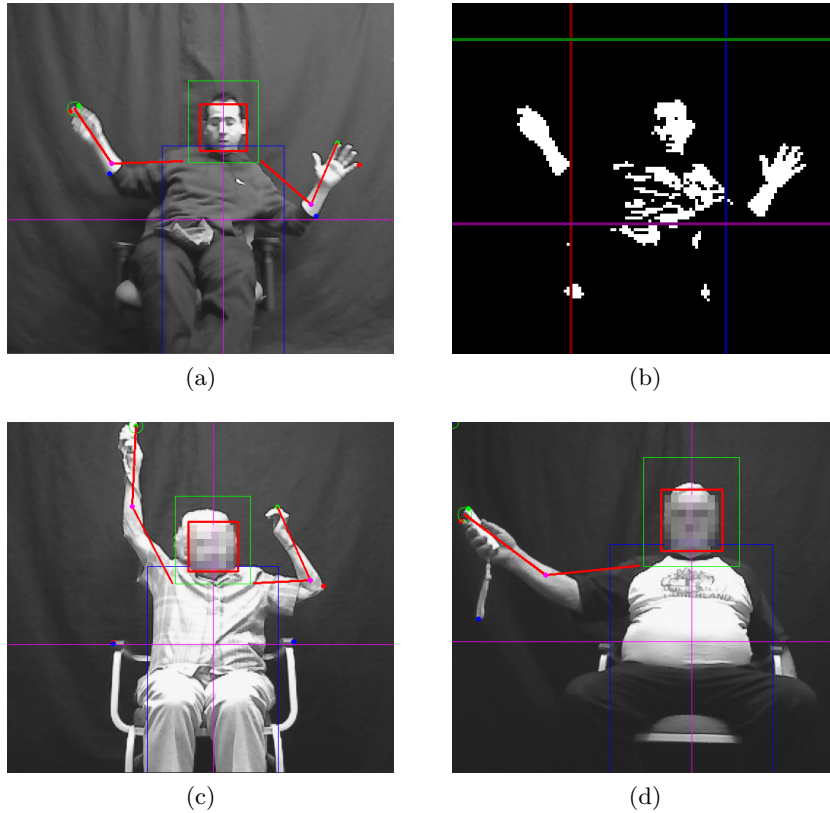


Figure 4.3: (a), (c), (d) Example face and arm angle detection results; (b) segmented image.

Several different approaches have been developed to accomplish tracking of human motions, both in 2D and 3D, including skeletonization methods [Fujiyoshi and Lipton, 1998; Jenkins et al., 2004], gesture recognition using probabilistic methods [Waldherr et al., 1998], and color-based tracking [Wren et al., 1997], among others. We opted to create an arm pose recognition system that takes advantage of our simplified exercise setup in order to achieve real-time results without imposing any markers on the user.

To simplify the visual recognition of the user, a black curtain was used to provide a static and contrasting background for fast segmentation of the user's head and hands,

the most important features of the arm pose recognition task, independent of the user's skin tone. The arm pose recognition algorithm consists of the following four steps:

1) *Create segmented image:* The original grayscale camera frame is converted into a black/white image by applying a single threshold over the image. All pixels below the threshold are set to black, and the rest to white. The threshold is set by the experimenter to appropriately segment out the user from the background, accounting for the specific lighting of the environment. Figure 4.3 (a) shows an original grayscale image captured from the camera, with the segmented image shown in Figure 4.3 (b).

2) *Detect the user's face:* The OpenCV frontal face detector is used to determine the location and size of the user's face. With these values, an estimate for the shoulder positions on both sides of the body is made.

3) *Determine hand locations:* The hand locations of the user are determined by examining the extrema points of the body pixels (maximum/minimum white pixel locations along x and y directions away from body) inside the region above the chest line and to the side of the face in the segmented image (subtracting the head and estimated body regions—represented by green and blue bounding boxes in Figure 4.3). The algorithm applies a simple set of rules, or heuristics, to choose which extrema points correspond to the hand location for a given arm. For example, if the highest white pixel (body pixel) in the segmented image on the left side is further than the approximate shoulder-to-forearm distance, then that is the left hand location.

These extrema point heuristics are valid for this domain, as the silhouette of a person demonstrating planar arm movements (i.e., only two degrees of freedom) will always feature at least one of these extrema points at each of the hand locations, due to the geometric properties of planar two-link manipulators. This planar two-link assumption thus validates the use of extrema point heuristics performed on a binary segmentation of the image (approximated silhouette) to detect the hand locations of the user.

4) *Determine arm angles:* Once the hand location for a given arm is found, the elbow point is estimated, which in turn provides the desired arm angles. The elbow point is estimated using the white pixel (body pixel) that lies furthest from the line connecting the hand position and the shoulder position, while also not exceeding the maximum allowable distance from the hand (to enforce forearm length restriction). Examples of arm angle detection results can be seen in Figure 4.3 (a), (c), and (d).

The vision module only searches for frontal face views, and thus detection rates depend largely on whether or not the user is facing the robot. However, high detection rates are not actually necessary for accurate gesture recognition, as the most recent detected face location is used in the arm detection procedure if a face is not found in the current frame. This substitution works well, since the user's head position generally remains stationary while the user is seated in a chair in front of the robot throughout the interaction.

The arm pose recognition algorithm runs with an average frame rate of 20 fps on a 2.4 GHz Intel Core 2 Duo processor, thus achieving our aim of real-time user activity recognition. While the arm detection procedure was not formally evaluated, the system was confirmed to be robust to different types of clothing, lighting, and user body types, with a notable arm angle estimation accuracy observed during the course of our user studies with both older and young adults.

Although a black curtain was used to aid the visual recognition procedure during our user studies, the algorithm does not rely on equipping the environment, but only on a binary user segmented image. The algorithm has been modified, for example, to receive a motion frame produced by sequential frame differencing in replacement of a color segmented image. The motion frame, once thresholded, serves as the segmentation input in the first step of the algorithm, and the user arm angles can be determined in the same manner as described. Using motion instead of color segmentation, along with

Kalman filtering of the positions of the hands and elbows, provides a more general recognition approach for domain-independent use, even though the overall algorithm remains the same.

The development of the SAR exercise system and visual recognition procedure predated the availability of the Microsoft Kinect. Our future implementations of the system will utilize Kinect-type 3D vision technology and thus do away with the curtain and the planar limits of the motions. Nevertheless, the 2D nature of the exercises was not noted as an issue by any of the participants in our user studies. Furthermore, the real-time vision algorithm described is general and can be applied to a variety of domains, including those wherein use of the Kinect is not appropriate or feasible (e.g., outdoors in sunlight).

4.6 Robot Behaviors

The behavior module represents the main loop of the SAR exercise coach system. It communicates with all of the other system modules and manages the flow of the interaction, including transitioning between the different exercise games, managing the session and game clocks, saving user performance and session statistics to the database, and handling user-requested breaks, among other session-related tasks.

4.6.1 Interaction Flow and Behavior Management

Each of the four exercise games is managed by its own game behavior. In accordance with our principle of encouraging users to be intrinsically motivated to engage in the task, the behavior module transitions between game types every 1-2.5 minutes in order to diminish the possibility of user boredom, which would negatively affect motivation level. The behavior module chooses which game to transition to based on a predefined session schedule, which specifies the order of the games to play and their respective durations.

The session schedule is defined a priori by the experimenter such that each session's order of games is unique. This requirement is enforced to reduce the predictability of the exercise regimen for the user, which might negatively affect his or her perception of the system [Bickmore et al., 2010]. The session schedule, however, though unique for each individual session, follows approximately the same structure throughout all of the sessions. Specifically, the Workout, Sequence, and Memory games are each allotted approximately 25% of the total session time (75% for all three games), with the Imitation game being allotted 15% of the total session time and the remaining 10% being allocated for user breaks, transition behaviors, and the introduction and conclusion of the session. The Imitation game receives a lower allotment compared to the other three games due to the fact that the user is in control of the exercise. Although the Imitation game promotes user autonomy and users often find the game stimulating, in pilot studies we've observed users may get bored if the game goes on for too long. Thus the Imitation game is usually only played for 1-1.5 minutes at a time during the course of the exercise session. If a session schedule is not set, or is incomplete, the robot chooses a game at random to play at regular game change intervals. The exact durations for each game were determined empirically by taking into account multiple factors, including the total duration of the session (20 minutes), the allotted game time ratio, the time required for the robot to familiarize the user with the specific game rules, and the minimum time estimated for the user to appropriately experience the game.

The session schedule allows the robot to guide the interaction and determine which games to play when. An alternative approach would be to allow the user to choose which game to play at each game transition interval, thus perhaps increasing user intrinsic motivation to engage in the task. Section 5.2 discusses a user study we conducted with elderly participants to investigate the role of user choice in the exercise scenario. For the user studies discussed in Section 5.1 and Chapter 6, we predefined the session schedules

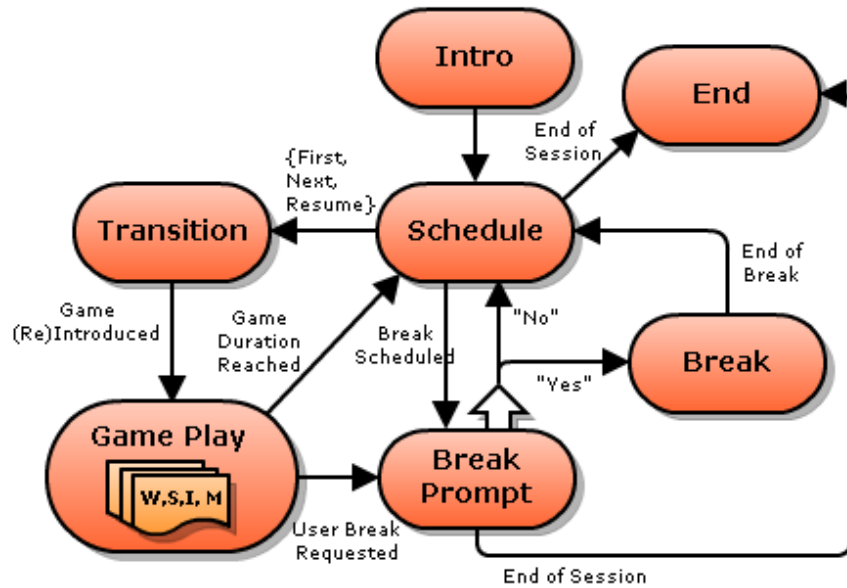


Figure 4.4: Diagram of the behavior module’s finite state machine (FSM). The FSM illustrates the flow of the interaction during the exercise sessions. The white arrow emanating from the Break Prompt state denotes a user decision.

and had the robot determine which games to play to ensure that each participant received similar session interactions to allow for unbiased comparison between them.

Though comprising most of the interaction time, the individual game behaviors are not the only behaviors available within the behavior module. Continuity behaviors, desirable for their relationship-building abilities (as discussed previously), were implemented in both the introduction and conclusion of the exercise sessions. These continuity behaviors are responsible for expressing the time-extended nature of the interaction. Examples of these behaviors include having the robot acknowledge past interactions, as in “It’s nice to see you again, Alice; last time was fun.” These also dictate how the robot re-introduces the games; the robot provides enough information to remind the

user of the rules of the exercise game without re-stating everything over again. This serves to acknowledge the fact that the user has already engaged in the exercise game in a previous session.

The remaining behaviors are related to breaks that provide the user with short periods of rest during the interaction so the exercise session is not too demanding. Breaks are scheduled every 5-6 minutes (3 breaks per session). Whenever a break is scheduled, the robot first confirms with the user if a break is desired, and depending on the response, either proceeds to take a 30-second break or continues with the next exercise game. Thus users who wish to rest are allowed to do so, and those who would rather continue playing have that option as well. This presents another example of user autonomy within the interaction, in addition to promoting a level of optimal challenge for each individual user to achieve maximal enjoyment.

A diagram of the behavior module's finite state machine is given in Figure 4.4 to illustrate the transitions between the different sub-behaviors and the overall flow of the interaction during the exercise sessions.

4.6.2 Feedback Procedure

All four of the game behaviors employ some form of feedback and guidance throughout the interaction. For example, in the Workout game the robot gives the users verbal feedback on how they can alter their arm positions to achieve the correct arm gesture for the current exercise, in the Memory game the robot can remind the users to perform a specific gesture number, and in the Imitation game the robot can prompt the users to move their arms if they happen to stop for more than a brief period. In many cases, the robot needs to provide the same verbal feedback to the user more than once, but must do so without sounding repetitive. As discussed earlier, strict repetitiveness, apart from being a possible annoyance to the user, represents a form of predictability in the nature of

the robot and may negatively affect user engagement in the task [Bickmore et al., 2010], which could ultimately be detrimental to the success of the healthcare intervention. We avoid repetitiveness in all comments made by robot, including feedback, praise, empathy, humor, game rules, transitions, and greetings and farewells.

We aimed for variability in verbal feedback so as to avoid the perception of repetitiveness as much as possible. When sequential verbal feedback comments differ, as in “Stretch out your right arm” followed by “Now raise your left arm above your head,” the variability is inherent. However, when the robot needs to report the same feedback comment in sequence, variability must be introduced into the feedback statement itself. We accomplish this variability by introducing filler words. Example filler words include the user’s name and the words “try” and “just.” Utilizing the user’s name not only serves to provide variability, but also works to gain the user’s attention, in addition to further personalizing the specific feedback comment. The feedback procedure chooses randomly at run time whether to use the user’s name, a filler word, or both when repeating a specific feedback phrase.

4.7 Summary

In this chapter implementation details of our SAR exercise coach system were discussed, including the development of a general visual user activity recognition algorithm capable of detecting planar arm movements in real time, and procedures for providing task-based feedback, encouragement, and guidance while introducing variability to limit potential repetitiveness of the system and maintain user engagement.

Chapter 5

Intrinsic Motivation User Studies with Older Adults

This chapter presents two user studies conducted with older adults and our SAR exercise coach to investigate methodologies for increasing user intrinsic motivation, based on related research in psychology and human-computer interaction. The first study investigated the role of praise and relative discourse in the interaction, while the second study investigated the role of choice and user autonomy in influencing user intrinsic motivation. The chapter presents the details of both studies in turn, including study design, dependent measures, study hypotheses and results.

5.1 Motivation Study 1: Praise and Relational Discourse

We designed and conducted an intrinsic motivation study to investigate the role of praise and relational discourse (politeness, humor, empathy, etc.) in the robot exercise system. Toward that end, the study compared the effectiveness and participant evaluations of two different coaching styles used by our system to motivate elderly users to engage in physical exercise. This section discusses the study methods employed, the subjective

and objective measures that were evaluated, and the outcomes of the study and system evaluation with elderly participants.

5.1.1 Study Design

The study consisted of two conditions, Relational and Non-Relational, to explore the effects of praise and communicative relationship-building techniques on a user’s intrinsic motivation to engage in the exercise task with the SAR coach. The study design was within-subject; participants saw both conditions, one after the other, and the order of appearance of the conditions was counter-balanced among the participants. Each condition lasted 10 minutes, totaling 20 minutes of interaction, with surveys being administered after both sessions to capture participant perceptions of each study condition independently. The following describes the two conditions in greater detail:

1) *Relational Condition*: In this condition the SAR exercise coach employs all of the social interaction and personalization approaches described in Chapter 3. Specifically, the robot always gives the user praise upon correct completion of a given exercise gesture (an example of positive feedback) and provides reassurance in the case of failure (an example of empathy). The robot also displays continuity behaviors (e.g., by referencing past experiences with the user), humor, and refers to the user by name, all with the purpose of encouraging an increase in the user’s intrinsic motivation to engage in the exercise session.

2) *Non-Relational Condition*: In this condition the SAR exercise coach guides the exercise session by providing instructional feedback as needed (e.g., user score, demonstration of gestures, verbal feedback during gesture attempts, etc.), but does not employ explicit relationship building discourse of any kind. Specifically, the robot does not provide positive feedback (e.g., praise) in the case of successful user completion of an exercise gesture, nor does it demonstrate empathy (e.g., reassurance) in the case of user

failure. The SAR coach also does not display continuity behaviors, humor, or refer to the user by name. This condition represents the baseline condition of our SAR exercise system, wherein the robot coach does not employ any explicit motivational techniques to encourage an increase in the user’s intrinsic motivation to engage in the task.

5.1.2 Participant Statistics

We recruited elderly individuals to participate in the study through a partnership with be.group, an organization of senior living communities in Southern California, using flyers and word-of-mouth. Thirteen participants responded and successfully completed both conditions of the study. The sample population consisted of 12 female participants (92%) and 1 male participant (8%). Participants’ ages ranged from 77-92, and the average age was 83 (S.D. = 5.28). Half of the participants ($n = 7$) engaged in the Relational condition in the first session, whereas the other half ($n = 6$) engaged first in the Non-Relational condition.

5.1.3 Measures

Survey data were collected at the end of the first and second sessions in order to analyze participant evaluations of the robot and of the interaction with the exercise system in both conditions. The same evaluation surveys were used for each session to allow for objective comparison between the two conditions.

In addition to these evaluation measures, at the end of the last exercise session we administered one final survey asking the participants to directly compare the two study conditions (labeled “first” and “second”) according to 10 evaluation categories. This survey allowed us to obtain a general sense of the participants’ preferences regarding the different SAR approaches and hence gauge their respective motivational capabilities.

Objective measures were also collected to evaluate user performance and compliance in the exercise task.

The following describes the specific evaluation measures captured in the post-session surveys, and the objective measures captured during the exercise sessions:

1) *Evaluation of Interaction*: Two dependent measures were used to evaluate the interaction with the robot exercise system. The first measure was the *enjoyableness of the interaction*, collected from participant assessments of the interaction according to six adjectives: enjoyable; interesting; fun; satisfying; entertaining; boring; and exciting (Cronbach's $\alpha = .93$). Participants were asked to rate how well each adjective described the interaction on a 10-point scale, anchored by "Describes Very Poorly" (1) and "Describes Very Well" (10). Ratings for the adjective "boring" were inverted to keep consistency with the other adjectives that reflect higher scores as being more positive. The enjoyableness of the interaction was measured to gain insight into the user's motivation level to engage in the task because, as Csikszentmihalyi states, intrinsically motivating activities are characterized by enjoyment [Csikszentmihalyi, 1975].

The second measure was the perceived *value or usefulness of the interaction*. Participants were asked to evaluate how well each of the following four adjectives described the interaction: useful; beneficial; valuable; and helpful (Cronbach's $\alpha = .95$). The same 10-point scale anchored by "Describes Very Poorly" (1) and "Describes Very Well" (10) was used in the evaluation. The perceived usefulness of the system was measured to estimate user acceptance and trust of the system in helping to achieve the desired health goals, which is necessary for the system to be successful in the long-term.

2) *Evaluation of Robot*: The *companionship of the robot* was measured based on participant responses to nine 10-point semantic differential scales concerning the following robot descriptions: bad/good; not loving/loving; not friendly/friendly; not

cuddly/cuddly; cold/warm; unpleasant/ pleasant; cruel/kind; bitter/sweet; and distant/close (Cronbach's $\alpha = .86$). These questions were derived from the Companion Animal Bonding Scale of Poresky et al. [1987]. The companionship of the robot was measured to assess potential user acceptance of the robot as an in-home companion, thereby demonstrating the capability of the system toward uses in independent living.

To assess the perceptions of the capabilities of the system in motivating exercise, we measured participant evaluations of the robot as an *exercise coach*. Participant evaluations of the robot as an exercise coach were gathered from a combination of the participants' reported level of agreement towards two coaching-related statements, and responses to three additional questions. The two statements and three questions were, respectively: I think Bandit is a good exercise coach; I think Bandit is a good motivator of exercise; How likely would you be to recommend Bandit as an exercise partner to your friends? How much would you like to exercise with Bandit in the future? How much have you been motivated to exercise while interacting with Bandit? (Cronbach's $\alpha = .88$). The two statements were rated on a 10-point scale anchored by "Very Strongly Disagree" (1) and "Very Strongly Agree" (10), and the three question items were each measured according to a 10-point scale anchored by "Not at all" (1) and "Very much" (10).

To quantify the effectiveness of the robot's social capabilities, we measured the *social presence of the robot*. Social presence is defined as the feeling that mediates how people respond to social agents; it strongly influences the relative success of a social interaction [Lee, 2004]. In essence, the greater the social presence of the robot, the more likely the interaction is to be successful. The social presence of the robot was measured by a 10-point scale anchored by "Not at all" (1) and "Very much" (10) using questionnaire items established from Jung and Lee [2004] (e.g., While you were exercising with Bandit, how much did you feel as if you were interacting with an intelligent being?) (Cronbach's $\alpha = .82$).

3) *Direct Comparison of Conditions*: The ten evaluation categories assessed by the direct comparison survey, which asked participants to choose between the first or second exercise sessions, were as follows: *enjoy more; more useful; better at motivating exercise; prefer to exercise with; more frustrating; more boring; more interesting; more intelligent; more entertaining; choice from now on*. Analysis of the direct-comparison data serves primarily to support and confirm the results obtained from the within-subjects analysis of the dependent measures across study conditions.

4) *User Performance Measures*: To help assess the effectiveness of the SAR exercise system in motivating exercise among the participants, we collected nine different objective measures during the exercise sessions regarding user performance and compliance in the exercise task. Most of the objective measures were captured during the Workout game, wherein the robot guides the interaction similar to a traditional exercise coach. These measures include the *average time to gesture completion* (from the moment the robot demonstrates the gesture, to successful user completion of the gesture), *number of seconds per exercise completed*, *number of failed exercises*, *number of movement prompts by the robot to the user due to lack of arm movement*, and *feedback percentage*. The feedback percentage measure refers to the fraction of gestures, out of the total given, where the robot needed to provide verbal feedback to the user regarding arm positions in order to help guide the user to correct gesture completion.

We also recorded the *maximum score* over all sessions, *average maximum score* among users, and *average time per gesture attempt* in the Memory game. For the Imitation game, the only measure captured was again the *number of movement prompts* by the robot due to lack of user arm movement.

5.1.4 Hypotheses

Based on the related research on the positive effects of praise and relational discourse on intrinsic motivation (see Chapter 3), seven hypotheses were established for this study:

- *Hypothesis 1:* Participants will evaluate the enjoyableness of their interaction with the relational robot more positively than their interaction with the non-relational robot.
- *Hypothesis 2:* Participants will evaluate the usefulness of their interaction with the relational robot more positively than their interaction with the non-relational robot.
- *Hypothesis 3:* Participants will evaluate the companionship of the relational robot more positively than that of the non-relational robot.
- *Hypothesis 4:* Participants will evaluate the relational robot more positively as an exercise coach than the non-relational robot.
- *Hypothesis 5:* There will be no significant difference between participant evaluations of the social presence of the relational robot and non-relational robot.

The reasoning behind this hypothesis is that people's sense of social presence is largely determined by the embodiment type and perceived intelligence of the social agent, which is assumed to be more or less equal in the two robot conditions.

- *Hypothesis 6:* Participants will report a clear preference for the relational robot over the non-relational robot when asked to compare both exercise sessions directly.
- *Hypothesis 7:* There will be no significant difference in participant exercise performance when interacting with either the relational or non-relational robot.

This hypothesis is based on the assumption that, due to the short-term nature of the study and novelty of the system, performance measures will be approximately equal between robot conditions.

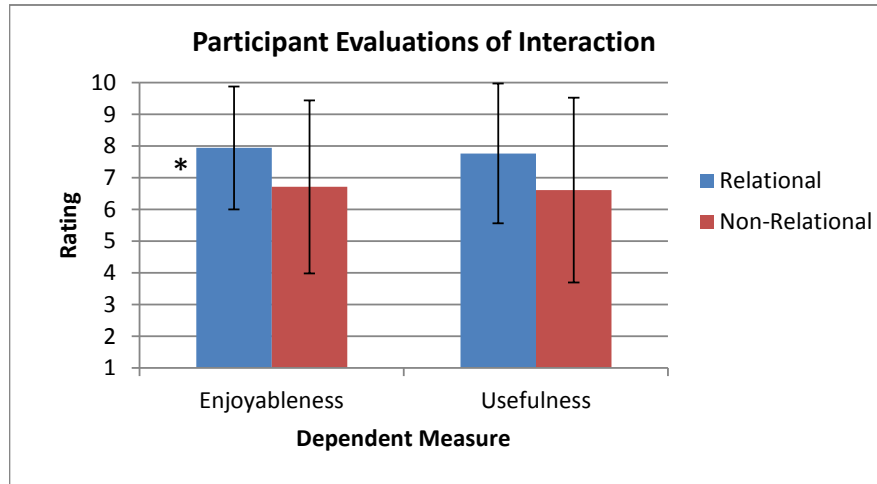
5.1.5 Results

1) Evaluation of Interaction Results

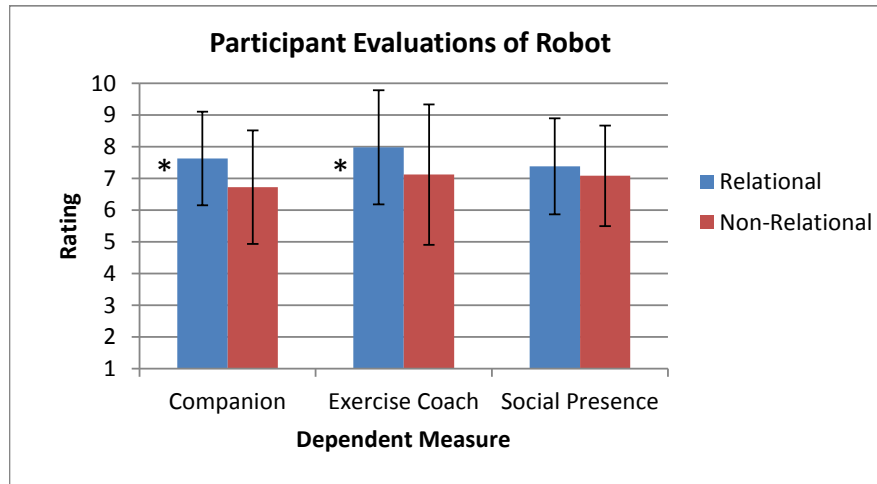
Participants who engaged with the relational robot in their first session, rated the Non-Relational condition on average 22% lower than the Relational condition in terms of enjoyment ($M_R=7.5$ vs. $M_{NR}=5.9$), and 23% lower in terms of usefulness ($M_R=7.5$ vs. $M_{NR}=5.8$). Similarly, the participants who instead engaged with the non-relational robot in their first session also expressed a greater preference for interacting with the relational robot by rating the Relational condition on average 10% higher than the Non-Relational condition in terms of enjoyment ($M_{NR}=7.6$ vs. $M_R=8.4$), and 7% higher in terms of usefulness ($M_{NR}=7.5$ vs. $M_R=8.0$).

Altogether, 85% of the participants (11 of 13) rated the Relational condition higher than the Non-Relational condition in terms of enjoyment, and 77% of the participants (10 of 13) rated the Relational condition higher in terms of usefulness than the Non-Relational condition.

To test for significant differences among the participant evaluations of the study conditions, we performed a Wilcoxon signed-rank test on the data to analyze matched pairs from the sample population's evaluations of both study conditions according to the dependent measures. Supporting Hypothesis 1, the results show that the participants evaluated the interaction with the relational robot as significantly more enjoyable/entertaining than the interaction with the non-relational robot ($W[12] = 4, p < .005$), and as somewhat more valuable/useful than the interaction with the non-relational robot, although not to a significant degree ($W[12] = 15.5, p < 0.10$), hence



(a)



(b)

Figure 5.1: (a) Plot of participant evaluations of the interaction, in terms of enjoyableness and usefulness, for both study conditions; (b) plot of participant evaluations of the robot (as a companion, exercise coach, and level of social presence) for both study conditions. Note: significant differences are marked by asterisks (*).

Hypothesis 2 was not supported by the data. For illustration purposes, Figure 5.1 (a) shows the average participant ratings of the enjoyableness and usefulness of the interaction for both study conditions.

2) Evaluation of Robot Results

Participants who engaged in the Relational condition in their first session rated the non-relational robot on average 11% lower than the relational robot in terms of companionship ($M_R=7.4$ vs. $M_{NR}=6.5$), 11% lower as an exercise coach ($M_R=7.7$ vs. $M_{NR}=6.9$), and 1% lower in terms of social presence ($M_R=7.2$ vs. $M_{NR}=7.1$). Greater positive scores for the relational robot were also reported by the participants who instead engaged first in the Non-Relational condition, having rated the relational robot on average 14% higher than the non-relational robot in terms of companionship ($M_{NR}=6.9$ vs. $M_R=7.9$), 10% higher as an exercise coach ($M_{NR}=7.4$ vs. $M_R=8.2$), and 8% higher in terms of social presence ($M_{NR}=6.9$ vs. $M_R=7.5$).

Altogether, 77% of the participants (10 of 13) rated the relational robot higher than the non-relational robot in terms of companionship, 77% of the participants (10 of 13) rated the relational robot more positively as an exercise coach, and the comparative ratings of social presence between the robot conditions were approximately equal, as 54% of participants (7 of 13) reported higher social presence for the relational robot.

We again analyzed the data to test for significant differences among participant evaluations across the two robot conditions by performing a Wilcoxon signed-rank test. The results show that the participants rated the relational robot as a significantly better companion than the non-relational robot ($W[13] = 14, p < .05$), supporting Hypothesis 3, and as a significantly better exercise coach than the non-relational robot ($W[11] = 7, p < .02$), in support of Hypothesis 4. As expected, there was no significant difference in the participant evaluations of social presence between both robot conditions ($W[12] = 28.5, p > 0.2$), confirming Hypothesis 5, with both robots receiving equally high ratings. The average participant ratings of both robot conditions for all three dependent measures are shown in Figure 5.1 (b).

Table 5.1: Participant Responses to Direct Comparison Survey Items

	Relational	Non-Relational	Both Equal
Enjoy More	10 (77%)	3 (23%)	0 (0%)
More Intelligent	11 (85%)	2 (15%)	0 (0%)
More Useful	11 (85%)	2 (15%)	0 (0%)
Prefer to Exercise with	11 (85%)	2 (15%)	0 (0%)
Better at Motivating	11 (85%)	2 (15%)	0 (0%)
More Frustrating	3 (23%)	10 (77%)	0 (0%)
More Boring	2 (15%)	10 (77%)	1 (8%)
More Interesting	10 (77%)	2 (15%)	1 (8%)
More Entertaining	10 (77%)	2 (15%)	1 (8%)
Choice from now on	11 (85%)	2 (15%)	0 (0%)

3) Direct Comparison Results

At the end of the final exercise session, participants were asked to directly compare both robot conditions with respect to 10 different evaluation categories; results are provided in Table 5.1. It is important to note that the study conditions were labeled as “first session” and “second session” on the survey. These labels would correspond to either the Relational condition or Non-Relational condition, depending on the order of the conditions in which each participant engaged, and were chosen to avoid any potential bias in the survey items.

The results support Hypothesis 6 by demonstrating that, regardless of the order of condition presentation, the participants expressed a strong preference for the relational robot over the non-relational robot. Specifically, the relational robot received 82% of the positive trait votes vs. 16% for the non-relational robot, with the remaining 2% shared equally between them. Other notable results include the high number of participants who rated the relational robot as more enjoyable (10 votes, 77%), better at motivating exercise (11 votes, 85%), more useful (11 votes, 85%), and the robot

they would choose to exercise with in the future (11 votes, 85%). In contrast, the non-relational robot received a high number of votes for being more frustrating (10 votes, 77%) and more boring (10 votes, 77%) than the relational robot.

4) User Exercise Performance Statistics

The collected statistics regarding participant performance in the exercise task were very encouraging as they demonstrated a consistently high level of user exercise performance and compliance with the exercise task. As expected, and in support of Hypothesis 7, there were no significant differences found in participant performance between the two study conditions, with both conditions reporting equally high performance among the participants. For example, the average gesture completion time for participants in the Relational condition was 2.45 seconds (S.D. = 0.65), compared to 2.46 seconds (S.D. = 0.78) for participants in the Non-Relational condition ($W[13] = 37, p > 0.2$). Given the lack of significant difference in user performance between the two conditions, the statistics presented in this section refer to the participant performance across all exercise sessions of the study.

User compliance and performance in the Workout game were high. The average gesture completion time was 2.46 seconds (S.D. = 0.70), and the overall exercise performance averaged 5.21 seconds per exercise (S.D. = 1.0), which also includes time taken for verbal praise, feedback, and score reporting from the robot. The low percentage of necessary corrective feedback, averaging 7.4%, zero failures, and zero movement prompts during the entire study, are all very encouraging results, as they suggest that the participants were consistently motivated to do well on the exercises throughout the interaction.

A summary of all statistics regarding user performance, including those from the Memory and Imitation games, can be found in Table 5.2.

Table 5.2: Participant Exercise Performance Statistics

Performance Measure	Mean (std.)
<i>Workout game:</i>	
Time to Gesture Completion (seconds)	2.46 (0.70)
Seconds per Exercise	5.21 (1.00)
Feedback Percentage	7.4% (4.8%)
Number of Failed Gestures	0
Number of Movement Prompts _W	0
<i>Memory game:</i>	
Maximum Score	6
Average Maximum Score	3.08 (1.12)
Time per Gesture Attempt (seconds)	8.57 (4.11)
<i>Imitation game:</i>	
Number of Movement Prompts _I	0.26 (0.53)

5.1.6 Discussion

The results of the study show a strong user preference for the relational robot over the non-relational robot, demonstrating the positive effects of praise and relational discourse in a healthcare task-oriented human-robot interaction scenario, and supporting all of our hypotheses with the exception of Hypothesis 2, which missed reaching significance by a small margin. Participants rated the relational robot significantly higher than the non-relational robot in terms of enjoyableness, companionship, and as an exercise coach.

Comments made by participants after the study further illustrate the positive response to the relational robot, including “It’s nice to hear your name, it’s personal. I felt more positive reinforcement,” and from another participant “The robot encourages you, compliments you; that goes a long way.” These results provide significant insight into how people respond to socially assistive robots, and confirm the positive influence that praise and relational discourse have on intrinsic motivation. These are of particular

importance for the healthcare domain, where effectiveness in social interaction, relationship building, and gaining user acceptance and trust are all necessary in ultimately achieving the desired health outcomes of the therapeutic interventions.

The effectiveness of the SAR exercise system was also demonstrated by the outcomes of the study. Not only did the participants rate the interaction with our robot coach as highly enjoyable/entertaining, suggesting they were intrinsically motivated to engage in the exercise task, but they also consistently engaged in physical exercise throughout the interaction, as demonstrated by the gathered user performance statistics. These results are very encouraging, as they clearly show that the system was successful in motivating elderly users to engage in physical exercise, thereby confirming its effectiveness and achieving the primary goal of the system.

5.2 Motivation Study 2: User Choice and Self-Determination

As discussed in Chapter 3, allowing the user to gain a sense of self-determination within a task, for example from choice of activity, has been shown to increase or be less detrimental to intrinsic motivation when compared to similar task conditions that do not involve choice [Fisher, 1978; Zuckerman et al., 1978]. To investigate the role of choice and user autonomy in influencing user intrinsic motivation in the robot exercise system, as well as to further test and validate the effectiveness of our system, we conducted a second user study with elderly participants.

5.2.1 Study Design

The study consisted of two conditions, Choice and No Choice, designed to test user preferences regarding choice of activity. The conditions differed only in the manner in

which the three exercise games (Workout, Imitation, Memory) were chosen during the exercise sessions. As in the first study, the design was within-subject; each participant engaged in both conditions one after the other, with the order of appearance counter-balanced among the participants. Each condition lasted 10 minutes, totaling 20 minutes of interaction. The following are descriptions of each condition in greater detail:

1) *Choice Condition*: In this condition the user is given the choice of which game to play at specific points in the interaction. The robot prompts the user to press the “Yes” button upon hearing the desired game, and then calls out the names of each of the three game choices. After the user has made a choice, the chosen game is played for a duration ranging from one to two minutes in length. Then, the robot asks the user if he would like to play a different game. Depending on the user’s response, the robot either continues playing the same game for another one to two minutes, or prompts the user again to choose the game to play next.

2) *No Choice Condition*: In this condition the robot chooses which of the three games to play at the specified game change intervals (every one to two minutes). The robot always changes games, to try to minimize any user frustration, as in this condition the robot is unaware of the user’s game preferences. For simplicity, in this condition the robot always chooses to first play the Workout game, followed by the Imitation and then Memory games, then cycles through them again in the same order.

5.2.2 Participant Statistics

We recruited elderly individuals to participate in the study again through our partnership with the be.group senior living organization. Eleven individuals participated in the first trial of the study, which was subsequently expanded to include thirteen additional participants. Therefore, a total of twenty four participants were recruited and successfully completed both conditions of the study. Half of the participants engaged in the

Choice condition in their first session, whereas the other half engaged first in the No Choice condition. The sample population consisted of 19 female participants (79%) and 5 male participants (21%). Participants' ages ranged from 68-89, and the average age was 77 (S.D. = 5.76).

5.2.3 Measures

As in the first study, survey data were collected at the end of the first and second sessions in order to analyze participant evaluations of the interaction with the exercise system in both conditions. The same evaluation surveys were used for each session to allow for objective comparison between the two conditions.

We administered an additional questionnaire at the end of the last session, asking the participants about their preferences regarding choice in the exercise system, in addition to various other opinion items for further evaluation of the exercise system.

The following describes the specific evaluation measures captured in the post-session surveys:

1) *Evaluation of Interaction*: The two dependent measures used to evaluate the interaction with the robot exercise system were the same as in the previous study, namely the *enjoyableness of the interaction*, and the *perceived value or usefulness of the interaction*. The ratings scales and survey items for each measure also remained the same.

2) *User Preferences Regarding Choice*: Three questionnaire items were used to assess participant preferences and opinions regarding choice in the exercise system (direct user input in choosing the exercise games). The first item asked participants to state their *session preference*, labeled as “first” or “second”, which referred to either the Choice or No Choice conditions, depending on each participant's session ordering. The ordinal labels were again chosen, as in the previous study, to avoid any bias in the survey item.

The second item asked the participants about *user choice*, specifically whether they preferred to choose the exercise games to be played, or whether they preferred to let the robot choose instead. This question is similar to the first item but in more direct terms. Lastly, the third item asked the participants about *added enjoyment due to user choice*, specifically asking whether or not having the ability to choose which game to play added to their enjoyment of the interaction.

3) *Evaluation of SAR System*: The last seven questionnaire items were used to obtain additional feedback on the user perceptions of and feelings towards the SAR exercise system. The first four of these items asked participants to rate, respectively: their perception of the *robot's intelligence*, their perception of the *robot's helpfulness*, the *level of importance* they put on their participation in the exercise sessions with the robot, and their mood in general during the exercise sessions. The rating scales were five-point Likert scales, anchored by “Not at all” (1) and “Very” (5) (e.g., “Not at all intelligent” and “Very intelligent”). The question regarding user mood during the sessions contained a modified scale, where the mood options ranged from “Irritated/frustrated” (1) to “Happy/joyful” (5), with the medium range being “Normal” (3). Participants were also asked to report their *favorite game*, *least favorite game*, and to state their choice of the *robot description* that best fit among four available options: companion, exercise instructor, game conductor, none of these.

5.2.4 Hypotheses

Based on the related research on the positive effects of user choice and autonomy on intrinsic motivation, discussed in Chapter 3, five hypotheses were established for this study:

- *Hypothesis 1:* Participants will evaluate the enjoyableness of their interaction in the Choice condition more positively than their interaction in the No Choice condition.
- *Hypothesis 2:* Participants will evaluate the usefulness of their interaction in the Choice condition more positively than their interaction in the No Choice condition.
- *Hypothesis 3:* Participants will report a clear preference for the Choice condition over the No Choice condition when asked to compare both exercise sessions directly.
- *Hypothesis 4:* Participants will report a clear preference for choosing the exercise games themselves, as opposed to having the robot choose which games to play during the interaction.
- *Hypothesis 5:* Participants will report feeling an increase in the enjoyment of the exercise task when given the opportunity to choose which games to play during the interaction.

5.2.5 Results

1) *Evaluation of Interaction Results*

The evaluation of interaction survey items were introduced after the first trial of the study, therefore the results presented here for the two interaction measures were analyzed from the data gathered solely from the thirteen participants of the expanded study. Nevertheless, all other survey results presented were gathered from all twenty four participants of the study.

Participants who engaged in the Choice condition in the first session rated the No Choice condition on average 7% higher than the Choice condition in terms of enjoyment ($M_C=7.5$ vs. $M_{NC}=8.0$), and 2% lower in terms of usefulness ($M_C=8.7$ vs. $M_{NC}=8.5$).

In slight contrast, the participants who instead engaged in the No Choice condition in their first session rated the Choice condition on average 4% higher than the No Choice condition in terms of enjoyment ($M_{NC}=8.5$ vs. $M_C=8.8$), and 5% higher in terms of usefulness ($M_{NC}=9.0$ vs. $M_C=9.5$).

Altogether, there was no clear participant preference for one condition over the other, as 62% of the participants (8 of 13) rated the No Choice condition higher than the Choice condition in terms of enjoyment, and 62% of the participants (8 of 13) rated the Choice condition higher in terms of usefulness than the No Choice condition.

We performed a Wilcoxon signed-rank test on the data and found no significant differences between participant evaluations of the two study conditions, neither with respect to the enjoyableness ($W[13] = 28.5, p > 0.2$), nor the usefulness of the interaction ($W[13] = 30.5, p > 0.2$). Thus, Hypotheses 1 and 2 were not supported by the data. Nevertheless, participant ratings for the enjoyableness ($M=8.18$, $S.D.=1.67$) and usefulness ($M=8.95$, $S.D.=1.63$) of the interaction across both conditions were very positive, with scores even higher than those seen in the previous study. These high evaluations of the SAR exercise system further illustrate the effectiveness of the system in instructing and motivating elderly users to exercise.

2) User Preferences Regarding Choice Results

The survey results regarding session preference indicated that 42% of the participants (10 of 24) preferred the No Choice condition, 33% of the participants (8 of 24) preferred the Choice condition, and 25% of the participants (6 of 24) expressed no preference for one condition over the other. Figure 5.2 (a) plots the participants' stated preferences of study conditions. The varied participant condition preferences indicate no clear preference for one over the other, and thus Hypotheses 3 was not supported. Concerning user choice in the exercise system, 62% of participants (15 of 24) reported preferring to

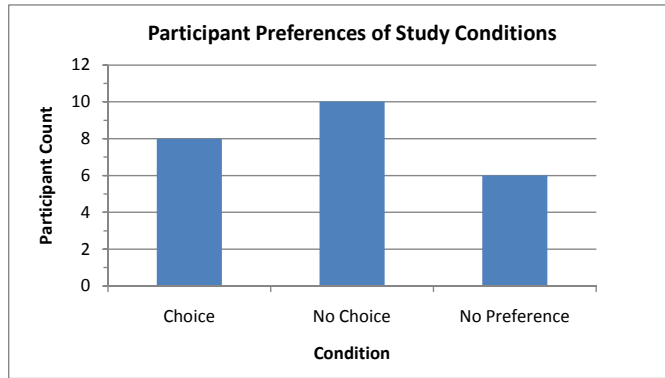
let the robot choose the games to play, with the remaining 38% of participants (9 of 24) preferring to choose the games themselves. The slight preference among participants for having the robot choose countered the reasoning of Hypothesis 4, which was not supported.

It is interesting to note that even though most participants preferred letting the robot decide which games to play, almost all of the participants, 92% (22 of 24) reported increased enjoyment of the task when given the opportunity to choose the exercise game to play. This result supports Hypothesis 5 and is consistent with the literature on the effects of user choice on intrinsic motivation [Fisher, 1978; Zuckerman et al., 1978].

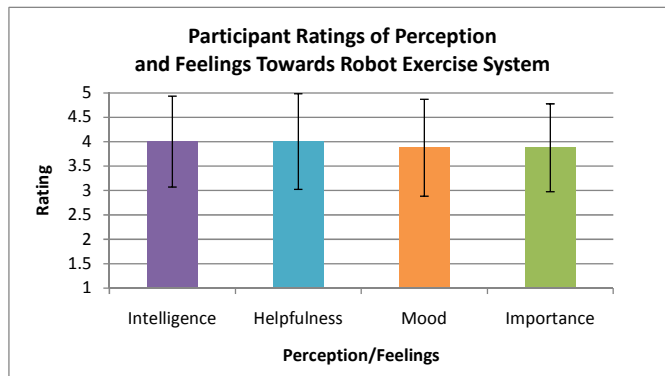
3) Evaluation of SAR System Results

The results of the survey questions regarding participant perceptions and feelings towards the SAR exercise system are very encouraging; the participants rated the robot highly in terms of intelligence (M=4.0, S.D.=0.93) and helpfulness (M=4.0, S.D.=0.97), attributed a moderately high level of importance to the exercise sessions (M=3.87, S.D.=0.89), and reported their mood throughout the sessions to be normal-to-moderately pleased (M=3.87, S.D.=0.99). These results are important because positive user perceptions of the agent's intelligence and helpfulness are a key part of establishing trust in the human-robot relationship. This, along with positive user mood and user-attributed importance to the therapeutic task, are in turn important for establishing and maintaining user intrinsic motivation. These are all key components for achieving long-term success in any socially assistive robot setting. An illustration of the results is shown in Figure 5.2 (b).

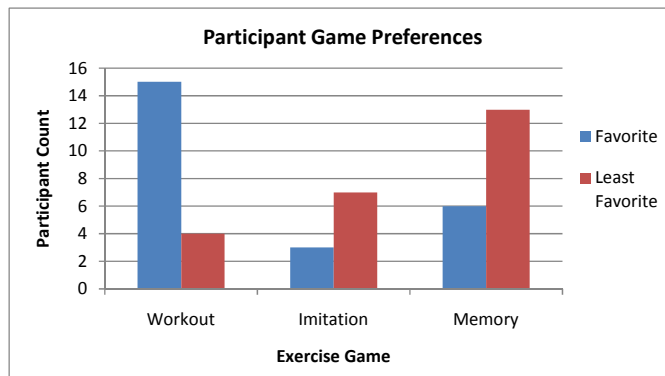
Regarding the exercise games, the participants largely favored (62%, 15 of 24) the Workout game over the others, wherein the robot serves as a traditional exercise coach,



(a)



(b)



(c)

Figure 5.2: Graphs of: (a) the participants' preferences of study condition; (b) the participants' ratings in response to survey questions on their perception of the robot's intelligence, helpfulness, their mood during sessions, and how important the sessions were to them; (c) the participants' preferences of exercise game.

with the Memory game being chosen most often as the participants' least favorite game (54%, 13 of 24). Figure 5.2 (c) summarizes the participants' game preferences.

The description most chosen by the participants as the best fit for the robot was that of an exercise instructor (67%, 16 of 24), not surprisingly, as opposed to that of a game conductor (25%, 6 of 24) or companion (8%, 2 of 24). While all of the descriptions represent characteristics of the robot in one form or another, the primary selection of an exercise coach by the participants illustrates the perception of the robot as an agent that they can trust and that is capable of helping, rather than simply entertaining.

5.2.6 Discussion

The results of the study showed no clear preference for one condition over the other, as the user enjoyment level of the interaction was reported to be equally high for both conditions, with or without user choice of activity. The high participant evaluations regarding the enjoyableness and usefulness of the interaction, the intelligence and helpfulness of the robot, and positive user mood and attributed importance to the exercise sessions, further validate the SAR system's effectiveness in motivating elderly users to engage in physical exercise.

The relatively mixed condition preferences among participants, or rather the lack of clear preference for the Choice condition, seem somewhat counter-intuitive given the positive effect that choice and user autonomy have been shown to have on task-based enjoyment [Fisher, 1978; Zuckerman et al., 1978]. One possible explanation for the mixed preferences may be that, since the robot's role in the interaction was that of an exercise instructor, some participants might have felt it was the robot's duty to determine the exercise regimen, and hence were comfortable relinquishing the choice of exercise games. Another possible explanation may be that the enjoyment derived from choosing the games did not outweigh the enjoyment derived from relaxation due

to the reduced responsibility of not having to choose the games. Both explanations seem plausible, as some of the participants reported preferring the robot to have the “responsibility” of steering the task. A third explanation may be that, given the short-term nature of the study, some participants may have needed more experience with the robot system before they felt confident enough to make task-based decisions themselves.

It is interesting to note that, even though the condition preferences were varied and nearly half of participants preferred letting the robot decide which games to play, all participants at one point or another during the study took advantage of having greater control in the Choice condition. Specifically, when given the option by the robot to change games, all participants at some point either chose to continue playing the same game they were playing, or chose to avoid playing a game they did not want to play. Neither of these cases could occur in the No Choice condition, as the robot was unaware of the user’s current game preferences.

This observation speaks to the value of user preference within the task scenario, suggesting that a hybrid approach that includes both user and robot decision making, personalized and tuned automatically for each user, might ultimately be the best solution for achieving a fluid and enjoyable task interaction for all users. For example, for users who prefer greater robot responsibility and input in the SAR-based task, the robot can recommend the “best” choices given the current task conditions and situation, giving the user an informed choice. Alternatively, for users who prefer greater control only once they’ve gained enough experience, the robot can initially make all task-based decisions until the user is ready and confident in making choices. For users who have a clear preference regarding who should make task-based decisions during interaction, the chosen strategy can be implemented continually throughout the sessions.

Clearly, no single fixed user-choice strategy is appropriate for all users; users have varied preferences regarding choice, and those preferences may even change over time.

Therefore, it is important that the strategy employed in SAR systems regarding user choice and autonomy be continually adapted to the specific user engaged in the interaction, thus personalizing the therapeutic intervention.

5.3 Summary

This chapter presented the results of two user studies conducted with older adults and our SAR exercise coach studying intrinsic motivation. The first motivation study showed a strong participant preference for the relational robot over the non-relational robot in terms of enjoyableness of the interaction, companionship, and as an exercise coach, in addition to demonstrating similar evaluations of both robots in terms of usefulness of interaction and social presence. These results illustrate the positive effects of motivational relationship building techniques, namely praise and relational discourse, on participant perceptions of the social agent and interaction in a health-related task scenario, and ultimately on user intrinsic motivation to engage in the task. The results of the second motivation study showed varying participant preferences regarding user choice within the exercise system, suggesting the need for customizable interactions automatically tailored to accommodate the personal preferences of the individual users.

The SAR exercise system was very well received, as demonstrated by both user studies, with high participant evaluations regarding the enjoyableness and usefulness of the interaction, companionship, social presence, intelligence, and helpfulness of the robot coach, and the positive mood and attributed importance of the exercise sessions. The system was also found to be effective in motivating consistent physical exercise throughout the interaction, according to various objective measures, including average gesture completion time, seconds per exercise, and percentage of user gesture attempts requiring robot feedback.

Chapter 6

Embodiment and SAR Evaluation

User Study

This chapter presents a user study that was designed and conducted with older adult participants to evaluate the effectiveness of our SAR approach and system design. Another aim of the study was to investigate the role of physical embodiment in the robot exercise system. Specifically, the study compared the effectiveness and participant evaluation of our physical humanoid robot to those of a computer simulation of the same robot shown on a flat-panel display.

6.1 Study Design

Study participants were divided into two groups, physical robot embodiment vs. virtual robot embodiment, and the study consisted of a total of four 20-minute sessions of exercise interaction with the system, conducted over a two-week period.

The following subsections describe the robot platforms and the between-subjects embodiment comparison method in detail.

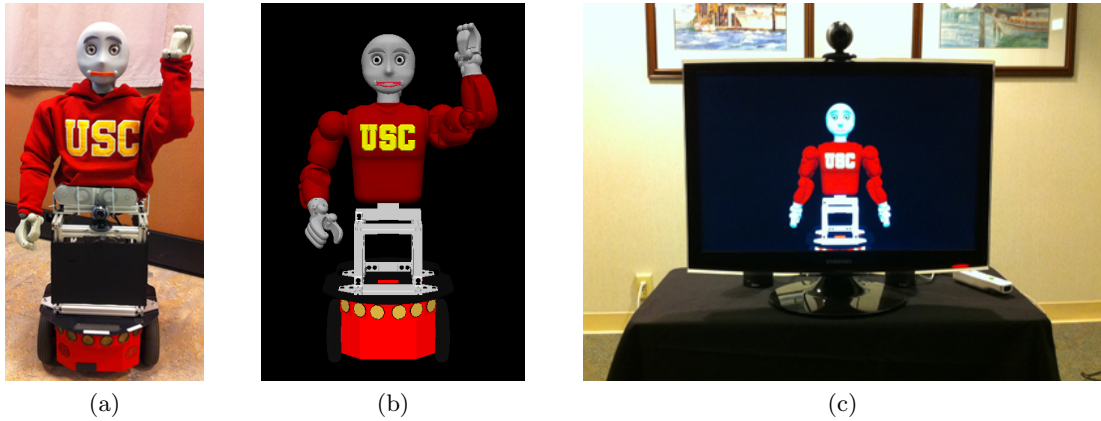


Figure 6.1: (a) Physical robot; (b) virtual robot computer simulation; (c) virtual robot on the screen, with camera.

6.1.1 Robot Platforms

To address the role of the robot's physical embodiment, we used the Bandit humanoid torso robot. A photograph of the physical robot can be seen in Figure 6.1 (a).

The robot's virtual embodiment consisted of a computer simulation of Bandit shown on a 27-inch flat-panel display. The size of the display was chosen to approximate the average size display that would be available in a typical household for use with the robot exercise system, including laptop displays (15 inch), computer monitors (24 inch), and television screens (40 inch). A sample computer simulation image of the virtual robot and a photograph of the virtual embodiment on the flat-panel display are shown in Figure 6.1 (b) and (c), respectively.

For the physical robot embodiment, the USB camera used in the visual user activity recognition procedure was placed at the waist of the torso, whereas for the virtual embodiment, the camera was attached to the top of the television display. The difference in camera location did not affect the accuracy of the visual recognition of the user's movements.

6.1.2 Between-Subjects Design

Survey data were collected at the end of the first and fourth sessions in order to analyze participant evaluations of the robot and the interaction with the exercise system in both conditions over time. Data gathered from the first and fourth session post-session surveys were analyzed using a two-tailed independent two-sample T-test assuming unequal variances to test for significant differences among the participant evaluations of the robot and the interaction across both conditions. Survey results from the fourth session were used to perform the final comparison analysis, as they were less likely to contain scores influenced by the effect of novelty.

6.2 Participant Statistics

We recruited elderly individuals to participate in the study through a partnership with be.group, an organization of senior living communities in Southern California, using flyers and word-of-mouth. We offered a \$50 Target gift card to those willing to participate in all four sessions of the study. Thirty-seven people responded, of whom four were omitted due to inconsistent/incorrect answers to survey questions that were used to identify questionable survey results. Thus, there were a total of 33 participants whose data were analyzed. Half of the participants were placed in the physical robot group ($n = 16$), and the other half were placed in the virtual robot group ($n = 17$). The sample population consisted of 27 female participants (82%) and 6 male participants (18%). Participants' ages ranged from 68-88, and the average age was 76 (S.D. = 6.32).

6.3 Measures

6.3.1 Evaluation of Interaction

There were two dependent measures used to evaluate the interaction with the robot exercise system. The first measure was the *enjoyableness of the interaction*, collected from participant assessments of the interaction according to six adjectives: enjoyable, interesting, fun, satisfying, entertaining, boring, and exciting (Cronbach's $\alpha = .92$). Participants were asked to rate how well each adjective described the interaction on a 10-point scale, anchored by "Describes Very Poorly" (1) and "Describes Very Well" (10). Ratings for the adjective "boring" were inverted to keep consistency with the other adjectives that reflect higher scores as being more positive.

The second measure was the perceived *value or usefulness of the interaction*. Similarly, participants were asked to evaluate how well each of the following four adjectives described the interaction: useful, beneficial, valuable, and helpful (Cronbach's $\alpha = .96$). The same 10-point scale anchored by "Describes Very Poorly" (1) and "Describes Very Well" (10) was used in the evaluation.

6.3.2 Evaluation of Robot

Companionship of the robot was measured from participant responses to nine 10-point semantic differential scales concerning the following robot descriptions: bad/good, not loving/loving, not friendly/friendly, not cuddly/cuddly, cold/warm, unpleasant/pleasant, cruel/kind, bitter/sweet, and distant/close (Cronbach's $\alpha = .89$). These questions were derived from the Companion Animal Bonding Scale of Poresky et al. [1987]. The companionship of the robot was measured to assess potential user acceptance of the robot as an in-home companion, thereby demonstrating the capability of the system to facilitate independent living.

Participants evaluated the *helpfulness of the robot* by rating four robot characteristics on a 10-point scale: useful, valuable, beneficial, and helpful (Cronbach's $\alpha = .96$). The rating scale was anchored by "Not at all" (1) and "Absolutely" (10). Using the same rating scale, the *intelligence of the robot* was measured according to the following four adjectives: competent, clever, intelligent, and smart (Cronbach's $\alpha = .93$).

To help capture the robot's social attributes, we measured both the *social attraction towards the robot* and the *social presence of the robot*. Social attraction was measured by a modified version of the Interpersonal Attraction Scale of McCroskey and McCain [1974]. Participants reported their level of agreement with the following four statements: I think Bandit could be a friend of mine; I think I could spend a good time with Bandit; I could establish a personal relationship with Bandit; I would like to spend more time with Bandit (Cronbach's $\alpha = .88$). The statements were rated on a 7-point scale anchored by "Very Strongly Disagree" (1) and "Very Strongly Agree" (7). Social presence was measured by a 10-point scale anchored by "Not at all" (1) and "Very much" (10) with questions from Jung and Lee [2004] such as the following: While you were exercising with Bandit, how much did you feel as if you were interacting with an intelligent being? (Cronbach's $\alpha = .87$).

To assess perceptions of the exercise capabilities of the system, we measured participant evaluations of the robot as an *exercise partner*. Four items were used, each measured according to a 10-point scale anchored by "Not at all" (1) and "Very much" (10): How much did you enjoy exercising with Bandit?; How likely would you be to recommend Bandit as an exercise partner to your friends?; How much would you like to exercise with Bandit in the future?; How much have you been motivated to exercise while interacting with Bandit? (Cronbach's $\alpha = .93$).

6.3.3 User Performance Measures

To help assess the effectiveness of the SAR exercise system in motivating exercise among the participants, we collected fifteen different objective measures during the exercise sessions regarding user performance and compliance in the exercise task.

Five performance measures were captured during user interaction in the Workout game, including the *average time to gesture completion* (from the moment the robot demonstrates the gesture to successful user completion of the gesture), *number of seconds per exercise completed*, *number of failed exercises*, *number of movement prompts* by the robot to the user due to lack of arm movement, and *feedback percentage*. The feedback percentage measure refers to the fraction of gestures, out of the total given, for which the robot needed to provide verbal feedback to the users regarding their arm positions in order to help guide them to correct gesture completion.

For the Sequence game, we captured four objective measures: the *average time to gesture completion*, *average number of sequences completed*, *average number of gesture pairs completed*, and the *feedback percentage*. In the Memory game, we recorded the *maximum score* over all sessions and all users, *average maximum score* among individual users, and *average time per gesture attempt*. For the Imitation game, the only measure captured was the *number of movement prompts* by the robot due to lack of user arm movement.

The remaining two measures concerned user activity during the entire exercise session: the *average total number of exercises completed*, and *number of breaks* taken.

6.3.4 Relation to Design Principles

The survey evaluation measures and the user performance measures together serve to evaluate our SAR system approach and design principles. Specifically, participant evaluations of each of the five SAR design principles were effectively captured by one or

Table 6.1: Summary of the relations between the five design principles and their related study measures (evaluation and/or performance based)

Design Principle	Related Measures (Survey and/or Performance)
1) <i>Motivating</i>	<p><i>Enjoyableness of interaction</i></p> <p>According to Csikszentmihalyi (1975), intrinsically motivating activities are characterized by user enjoyment.</p>
2) <i>Fluid and Highly Interactive</i>	<p><i>Average time to gesture completion; seconds per exercise; feedback percentage</i></p> <p>The fluidity and real-time interactive nature of the task is most appropriately illustrated by the speed and accuracy of the objective user performance statistics (which are tightly linked to the feedback and responsiveness of the robot).</p>
3) <i>Personable</i>	<p><i>Helpfulness; companionship; social attraction; social presence</i></p> <p>The robot’s personable qualities, such as expressing empathy, reassurance, praise, reference to past user performance, continuity between sessions, humor, politeness, and referring to the user by name, are best characterized by the perceived helpfulness and companionship of the robot. Furthermore, more general assessments of the personable traits of the robot (and of the robot’s embodiment) are captured by the social attraction and social presence measures.</p>
4) <i>Intelligent</i>	<p><i>Intelligence of the robot</i></p> <p>Participant evaluations of the perceived intelligence of the robot provide an appropriate one-to-one assessment of this design principle in practice.</p>
5) <i>Task-Driven</i>	<p><i>Value/usefulness of interaction; robot as exercise partner; all user exercise performance measures</i></p> <p>Participant evaluations of the usefulness of the interaction, and of the robot as an exercise partner, serve to illustrate the perceived effectiveness of the system in encouraging users to exercise (important for establishing user trust in the system). Furthermore, the objective performance measures serve to demonstrate the actual effectiveness of the system in eliciting physical exercise among the users (the goal of the specific healthcare task).</p>

more of the described study measures. Table 6.1 provides a summary of the relations between each of the five design principles outlined in Chapter 3 and their corresponding survey and/or objective performance measures.

6.4 Hypotheses

Based on the related research on embodiment effects, six hypotheses were established for the embodiment comparison portion of this study:

- *Hypothesis 1:* Participants will evaluate the enjoyableness of their interaction with the physical robot more positively than their interaction with the virtual robot.
- *Hypothesis 2:* Participants will evaluate the value/usefulness of their interaction with the physical robot more positively than their interaction with the virtual robot.
- *Hypothesis 3:* Participants will evaluate the helpfulness of the physical robot more positively than that of the virtual robot.
- *Hypothesis 4:* Participants will be more socially attracted to the physical robot than the virtual robot.
- *Hypothesis 5:* Participants will experience a greater sense of social presence when interacting with the physical robot than when interacting with the virtual robot.

6.5 Results

6.5.1 Embodiment Comparison Results

Between-Subjects Comparison Results

A two-tailed independent T-test was performed on the survey data following the fourth exercise session to compare participant evaluations of the robot embodiments and of the overall SAR interaction across the two study groups. Table 6.2 provides the complete set of between-subjects comparison results.

Table 6.2: Results of between-subjects data comparison for all $n = 33$ older adult participants showing means and standard deviations (in parentheses)

Dependent Measure	Physical Robot	Virtual Robot
<i>Interaction Evaluation</i>		
	<i>Between-Subjects Analysis</i>	
Enjoyable	7.51 (1.77)*	6.00 (2.01)
Valuable/Useful	8.14 (1.66)*	6.19 (2.39)
<i>Robot Evaluation</i>		
Helpful	8.11 (1.98)*	6.26 (1.98)
Social Attraction	4.70 (1.40)*	3.61 (1.54)
Social Presence	7.88 (0.94)*	6.47 (2.01)
Companion	7.48 (2.07) [†]	6.23 (1.84)
Intelligence	8.17 (2.02) [†]	6.76 (2.09)
Exercise Partner	7.18 (2.17) [†]	5.76 (2.18)

[†] $p < .10$, * $p < .05$

Consistent with Hypothesis 1, the participants evaluated the interaction with the physical robot embodiment as more enjoyable than the interaction with the virtual robot embodiment ($t[31] = 2.29, p < .03$). Hypothesis 2 was supported by the data as well, as the participants evaluated the interaction with the physical robot as more valuable/useful than the interaction with the virtual robot ($t[29] = 2.72, p = .01$).

Regarding the direct evaluations of both robot embodiments, the participants rated the physical robot as more helpful than the virtual robot ($t[31] = 2.66, p = .01$), consistent with Hypothesis 3, and as more socially attractive ($t[30] = 2.09, p < .05$), supporting Hypothesis 4. Concerning social presence, the data were consistent with Hypothesis 5, as the participants reported feeling a stronger sense of social presence with the physical robot than with the virtual robot ($t[23] = 2.59, p < .02$).

Evaluations of the non-hypothesis-testing system performance measures were also favorable to the physical robot, though not to a significant degree, as the participants

rated the physical robot as somewhat more of a companion ($t[30] = 1.81, p < .08$), more intelligent ($t[31] = 1.96, p < .06$), and a moderately better exercise partner ($t[31] = 1.87, p = .07$) than the virtual robot.

Performance Comparison Results

There were no significant differences found in participant performance between the two study groups, with participants from both the physical robot and virtual robot groups achieving equally high performance. For example, the average gesture completion time in the Workout game during the fourth session for the physical robot group was 2.29 seconds (S.D. = 0.62), compared to 2.08 seconds (S.D. = 0.45) for participants in the virtual robot group ($t[29] = 1.11, n.s.$), and the feedback percentage for the physical robot group was 6.0% (S.D. = 5.7), compared to 4.4% (S.D. = 4.2) for the virtual robot group ($t[29] = 0.96, n.s.$). Further discussion of the exercise performance statistics of the older adult study participants is provided in Chapter 6.5.2.

Discussion of Embodiment Comparison Results

The results of the between-subjects embodiment comparison analysis show a strong participant preference for the physical robot embodiment over the virtual robot embodiment, as hypothesized. More generally, the results illustrate the wide-reaching effect that SAR agent embodiment type has on the interaction and overall perception of the SAR system by the user. In particular, the study results show the influence of agent embodiment on user motivation, enjoyment, and perceptions of the value/usefulness of the interaction, in addition to showing the influence of embodiment on the perceived personable qualities of the robot, including helpfulness, social attraction, and social presence.

As previously mentioned, this study is to the best of our knowledge the first to comprehensively demonstrate the positive effect of physical embodiment in a SAR-guided healthcare scenario, wherein the SAR agent serves as both an instructor and active participant in the healthcare task with target users. Furthermore, our study contrasts with similar studies investigating the role of embodiment in HRI with older adults (e.g., [Heerink et al., 2010]) in that the participant evaluations in our study were based on multiple sessions of interaction with our SAR system. This is an important distinction, because as previously stated, the between-subjects comparison analysis was performed utilizing data collected following the fourth session of interaction instead of the first session in order to minimize the effect of novelty in the evaluations. However, the multiple-session study design was not only useful for attenuating the effect of novelty, but also for tracking user perceptions of the system and performance over time. For instance, it is interesting to note that although the survey results collected after the first session showed a participant preference for the physical robot over the virtual robot, the results were not statistically significant, in contrast to the results following the fourth session.

In addition, the physical robot group's ratings of the interaction and robot coach showed a positive trend across sessions, increasing by 6.2% (S.D. = 6.35) on average from the first session to the fourth session across all dependent measures. Conversely, the virtual robot group's ratings of the interaction and robot coach showed a negative trend across sessions, decreasing by 5.4% (S.D. = 6.64) on average from the first to the fourth session. The difference in the trends regarding participant ratings across sessions for the two study groups was found to be statistically significant ($t[14] = 3.56, p < .01$), showing a clear distinction between the study conditions regarding user perceptions of the system over time, and further demonstrating the positive influence of physical embodiment in SAR-guided interaction.

All of these results illustrate the importance of conducting multiple sessions of interaction for proper system evaluation and comparison, thus validating our study design.

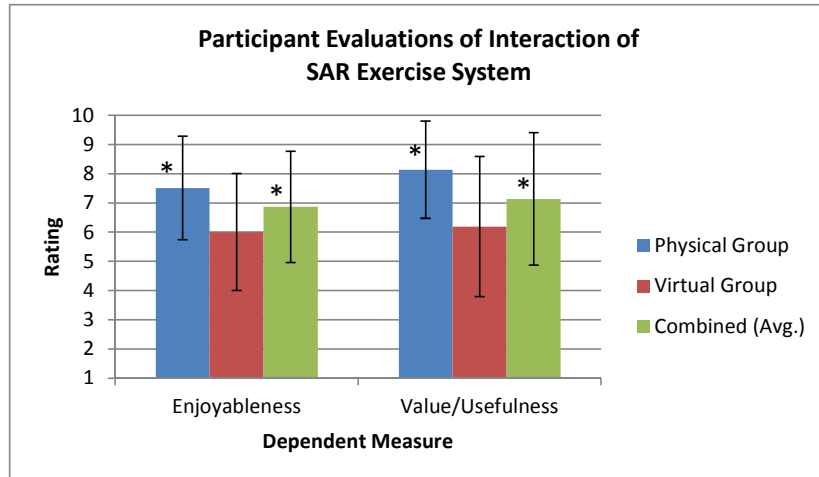
6.5.2 SAR System Evaluation Results

In order to evaluate the effectiveness of the SAR exercise system, we analyzed the combined data of the physical and virtual robot group's fourth session of interaction with the SAR exercise system. Therefore, the SAR system evaluation results regarding user perceptions and user exercise performance were gathered from all 33 older adult participants. The combined data were chosen for the final evaluation of the system for their embodiment-independent nature, and were constructed by averaging the evaluation ratings and user exercise performance of both study groups.

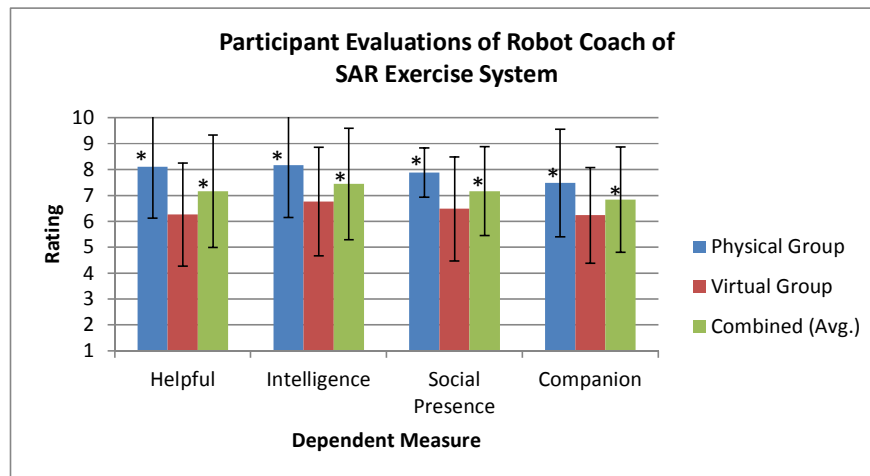
User Evaluations of the SAR System

To analyze the user evaluations of the SAR exercise system, we performed a two-tailed independent T-test assuming unequal variances to test for significant differences between the participant ratings of the subjective measures and a neutral evaluation rating. The neutral evaluation rating distribution was obtained from a uniform sampling of the rating scale (integers from 1 to 10) for the approximate number of participants, and has a mean rating of 5.5 (S.D. = 2.90). This uniform sampling assumes no prior information regarding user perceptions of the system, and thus is deemed neutral.

The combined data of both study groups showed that participants evaluated the interaction with the SAR exercise system as enjoyable ($M = 6.8$, $S.D. = 1.90$) and valuable/useful ($M = 7.1$, $S.D. = 2.26$). The ratings for both measures were found to be significantly more positive than a neutral evaluation (enjoyableness: $t[49] = 2.16$, $p < .05$; usefulness: $t[55] = 2.46$, $p < .02$). These results illustrate the effectiveness of the system



(a)



(b)

Figure 6.2: (a) Plot of participant evaluations of the interaction of the SAR exercise system in terms of enjoyableness and usefulness; (b) plot of participant evaluations of the robot coach of the SAR exercise system in terms of helpfulness, intelligence, social presence, and as a companion. Note: Significant differences ($p < .05$) in comparison to neutral rating distribution are marked by asterisks (*).

in promoting user intrinsic motivation, which is characterized by enjoyment [Csikszentmihalyi, 1975], and in guiding the task-driven interaction toward achieving beneficial health outcomes for the user. Both of these factors show the successful application of

our SAR design principles. A plot of the interaction evaluation results, from each study group and combined, is provided in Figure 6.2 (a).

Regarding user perceptions of the SAR exercise system's robot coach, the combined data showed participants rated the robot highly and significantly more positive than neutral in terms of helpfulness ($M = 7.2$, $S.D. = 2.16$; $t[53] = 2.53$, $p < .02$), intelligence ($M = 7.4$, $S.D. = 2.15$; $t[53] = 2.98$, $p < .01$), social presence ($M = 7.2$, $S.D. = 1.71$; $t[46] = 2.71$, $p < .01$), and as a companion ($M = 6.8$, $S.D. = 2.02$; $t[51] = 2.09$, $p < .05$). The participants also rated the robot coach favorably in terms of social attraction ($M = 4.1$ (on a 7-point scale), $S.D. = 1.55$; $t[63] = 0.35$, $n.s.$) and as an exercise partner ($M = 6.5$, $S.D. = 2.26$; $t[55] = 1.43$, $n.s.$). These results illustrate that the participants perceived the robot coach as having a personable nature and being intelligent, both of which characteristics aid in the development of trust within the human-robot relationship and were design goals of our SAR system approach aimed at providing successful therapeutic interventions. A plot showing participant evaluations of the SAR system's robot coach is shown in Figure 6.2 (b).

The results of the user evaluation of the SAR exercise system were very encouraging, as they showed a notable level of user acceptance of the system, as evidenced by the high ratings across each of the subjective measures, highlighting the effectiveness of our SAR system design principles.

User Exercise Performance Statistics

The collected statistics regarding participant performance in the exercise task were also very encouraging, as they demonstrated a consistently high level of user exercise performance and compliance. As previously stated, there were no significant differences in participant exercise performance between the two study groups; therefore, the

Table 6.3: User exercise performance statistics for all $n = 33$ older adult participants engaging with the SAR exercise system, showing means and standard deviations (in parentheses)

Performance Measure	Mean (std.)
<i>Workout game:</i>	
Time to Gesture Completion (seconds)	2.18 (0.54)
Seconds per Exercise	5.07 (0.58)
Feedback Percentage	5.2% (4.9%)
Number of Failed Gestures	0
Number of Movement Prompts _W	0
<i>Sequence game:</i>	
Time to Gesture Completion (seconds)	5.95 (1.26)
Number of Sequences Completed	5.0 (1.41)
Number of Gesture Pairs Completed	15.1 (4.06)
Feedback Percentage	22.3% (11.8%)
<i>Memory game:</i>	
Maximum Score	6
Average Maximum Score	3.26 (1.31)
Time per Gesture Attempt (seconds)	7.64 (4.03)
<i>Imitation game:</i>	
Number of Movement Prompts _I	0.64 (1.13)
<i>Entire Session:</i>	
Total Number of Exercises Completed	103.59 (27.3)
Number of Breaks Taken	1.08 (1.23)

results presented in this section represent the combined performance statistics of all participants, captured during the fourth session of interaction.

User compliance and performance in the Workout game were high. The average gesture completion time was 2.18 seconds (S.D. = 0.54) and the overall exercise performance averaged 5.07 seconds per exercise (S.D. = 0.58), which included time taken for verbal praise, feedback, and score reporting from the robot. The low percentage of necessary

corrective feedback, averaging 5.2%, combined with zero failures and zero movement prompts during the interaction session, were all very important results, as they suggested that the participants were motivated to do well on the exercises consistently throughout the interaction. Furthermore, the quick pace of the exercise movements accomplished by the participants served to illustrate the fluid, highly interactive nature of the SAR-guided interaction; this is beneficial for increasing user intrinsic motivation and engagement in the task, and is the goal of our second SAR design principle.

The participant statistics for the remaining games also indicated high compliance and performance. A summary of all statistics regarding user performance, including those from the Sequence, Memory, and Imitation games, can be found in Table 6.3. As evidenced by the results, the average time to gesture completion in the Sequence and Memory games was greater than in the Workout game. This result was to be expected, as these two games incorporated a memory component, and during gameplay participants would often forget which gesture poses to complete in sequence. The increased feedback percentage in the Sequence game was also a result of this scenario, with the robot providing additional feedback and re-demonstrating the correct gestures to the users when necessary. It is important to note, however, that even though users had more difficulty in these two games, as expected, the participants were still able to achieve high performance (e.g., in the Sequence game, the participants required no feedback at all to successfully complete nearly 80% of the gestures). Other notable results include the low number of movement prompts during the Imitation game, averaging less than one per user, the high total number of exercises completed during the entire session ($M = 103.59$, $S.D. = 27.3$), and the low number of breaks taken by the users during interaction, all of which demonstrate a high level of participant intrinsic motivation to engage in the exercise games.

Table 6.4: Means, standard deviations, and intercorrelations among dependent measures

Dependent Measure	Mean (Std. Dev.)	Completion Time	Feedback Percentage	Total Exercises	Maximum Score
Gesture completion time in seconds	2.54 (.89)	1.00			
Percentage of gestures where feedback was required	8.2% (7.9)	.91	1.00		
Total exercises performed in session	100 (18.6)	.78	.76	1.00	
Maximum score in Memory game	3.11 (1.35)	.24	.32	.50	1.00

Analysis of User Performance Results

Out of the fifteen user performance measures captured at each session, four were selected for further analysis: average completion time in the Workout game, percent feedback required in the Workout game, total exercises for the entire session, and average maximum score in the Memory game.

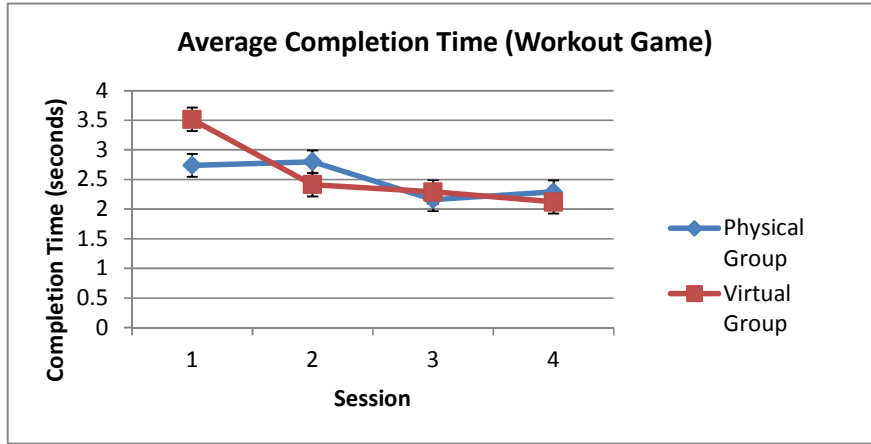
We began our analysis by seeking to determine whether the four measures might all be capturing a single underlying variable – user performance. As the measures use different scales (time, percentage, counts), each was standardized so that the mean was 0 and the standard deviation was 1. In addition, since the first two measures indicate better performance with a lower value and the last two with a higher value, the values of the first two measures were reversed by multiplying the standardized score by -1, so that higher values were uniformly indicative of better performance in all measures. Table 6.4 displays the intercorrelations of the four measures, along with the unstandardized means and standard deviations.

Table 6.5: Analysis of variance testing the fixed effect tests of condition, sessions, and measure on performance

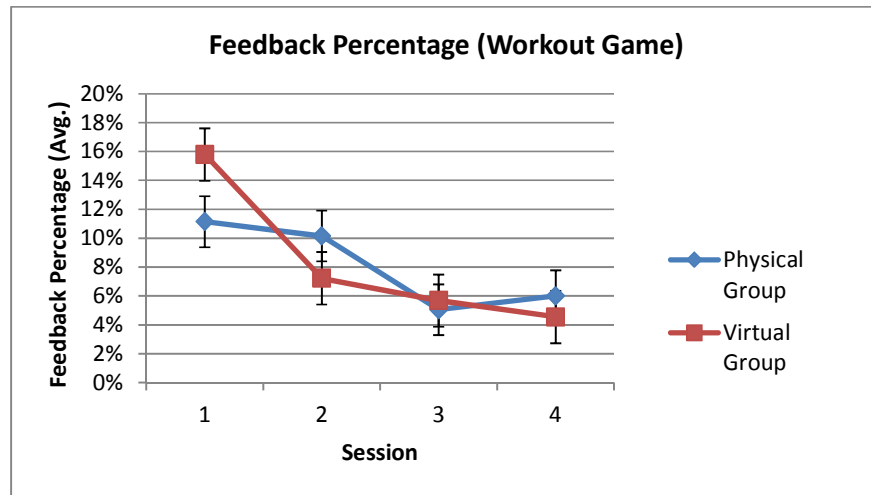
Variable	<i>df</i>	<i>F</i> ratio	<i>p</i> value
Condition (physical vs. virtual)	1, 31		
Session (1 – 4)	3, 465	44.5	< .0001
Condition * Session	3, 465	4.0	< .01
Measure (four)	3, 465		
Condition * Measure	3, 465		
Session * Measure	9, 465	2.3	= .01
Condition * Session * Measure	9, 465		

To evaluate performance across sessions, we performed a mixed models (multi-level) analysis of variance on the standardized scores of the measures within each session. Session is a within-subject factor. Condition, such as whether a physical or virtual robot led the task, is a between-groups measure. The effect of multiple measures on the same participants was controlled for. The analysis showed that, overall, there is no condition main effect, but there is a session main effect (performance improved over time), and there is a significant interaction of condition with session ($F[3, 465] = 4.0, p < .01$); see Table 6.5. The interaction effect arises from the fact that the physical robot prompted better performance in the first session but not in the later sessions. The contrast within the first session between the physical and virtual robot is $F[1, 49] = 3.98, p = 0.5$.

Analyzing the effects of each (unstandardized) performance measure showed that, in comparison with the virtual robot, the physical robot reduced completion time (interaction $F[3, 93] = 5.6, p < .01$) and feedback percentage (interaction $F[3, 93] = 2.9, p < .05$) in the first session but not in the remaining three sessions. Furthermore, according to a contrast, the difference between the physical and virtual robot in the first session is significant for completion time ($F[1, 82] = 7.8, p < .01$) and marginally significant for feedback percentage ($F[1, 81] = 3.3, p = .07$). The other measures do not show significant effects of condition; for example, total exercises increased notably over sessions in



(a)



(b)

Figure 6.3: Participant performance results across all four sessions of interaction for both study groups, showing (a) Average gesture completion time (Workout game); (b) feedback percentage (Workout game). Note: Means are least squares means from the ANOVA; error bars are standard error. The statistical difference between the physical and virtual robot in the first session is significant for completion time ($p < .01$), and marginally significant for feedback percentage ($p = .07$), see text.

both conditions ($F[3, 93] = 31.99, p < .0001$). The plots for average completion time and feedback percentage are shown in Figures 6.3 (a) and (b), respectively.

6.5.3 Study Expansion with Young Adults

To analyze and compare user evaluations and embodiment effects across age groups, we expanded the study to include 33 young adult participants (6 female, 27 male), yielding a combined sample of 66 participants, all of whom engaged in five sessions of interaction with our system for evaluation purposes (330 sessions total). The results of the study with young adults engaging with the SAR exercise system were largely consistent with those observed with the older adult participants. The within-subjects and direct comparison results showed an overwhelming preference for the physical robot embodiment over the virtual robot embodiment, although no significant differences were observed between-subjects.

As expected from the results of both samples, the combined results of older and young adults (66 participants) displayed a strong participant preference for the physical robot over the virtual robot across all comparison methods. Among the combined results, a two-sided exact binomial test showed the physical robot coach received significantly more positive votes than the virtual robot coach upon direct comparison (425 votes vs. 103, $p < .0001$).

6.6 Summary

This chapter presented a multi-session user study conducted with older adults to evaluate the effectiveness of our robot exercise system across a variety of user performance and evaluation measures. The results of the study validate our SAR approach and its effectiveness in motivating physical exercise in older adults; the participants engaged in physical exercise with high performance consistently throughout the interaction sessions, rated the SAR system highly in terms of enjoyableness and usefulness of the

interaction, and rated the robot coach highly in terms of helpfulness, social attraction, social presence, and companionship.

In addition, the role of physical embodiment was investigated in our exercise system in the same multi-session user study by comparing the effectiveness of a physically embodied robot coach to that of a virtually embodied robot coach (a computer simulation of the same robot). Results from interaction with elderly participants ($n = 33$) were presented. The results of the between-subjects embodiment comparison study show a strong user preference for the physical robot embodiment over the virtual robot embodiment in our SAR exercise system. Consistent with our stated hypotheses, participants reported the interaction with the physical robot as being more enjoyable and more valuable/useful than the interaction with the virtual robot. Furthermore, the participants evaluated the physical robot as more helpful, more socially attractive, and as having greater social presence than the virtual robot.

Chapter 7

Spatial Language Understanding for Human-Robot Interaction

This chapter describes our research in the area of user-guided interactions. Specifically, it presents a methodology for autonomous service robots to receive and interpret natural language instructions involving spatial relations from non-expert users. The approach is motivated by related research in linguistics, cognitive science, neuroscience, and computer science, proposes the encoding of spatial language within the robot *a priori* as primitives, and provides a computational framework for human-robot interaction which integrates the proposed model.

7.1 Approach

7.1.1 Semantic Fields

The *semantic field* of a spatial preposition is analogous to a probability density function (pdf), parameterized by schematic figure and reference objects, that assigns weight values to points in the environment depending on how accurately they capture the meaning

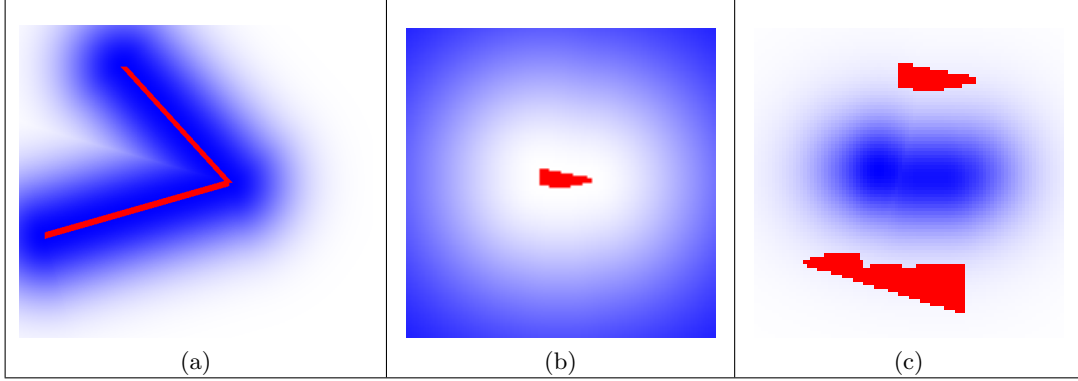


Figure 7.1: Semantic fields for static prepositions (a) near; (b) away from; (c) between.

of the preposition (e.g., points closer to an object have higher weight for the preposition ‘near’). This field representation for the semantics of spatial prepositions, while based on insights gathered from neuroscience research in rats, was shown by O’Keefe [2003] to closely resemble the form and continuous nature of spatial preposition representations demonstrated by humans [Logan and Sadler, 1996]. These types of continuous spatial field functions have also been shown to transfer seamlessly into higher dimensions [Gapp, 1994], thus enabling similar relational comparisons in 2D and 3D space. Example semantic fields are shown in Figure 7.1 for the static prepositions “near”, “away from”, and “between” for illustration purposes. The semantic field for *near* was produced by calculating the weights ($\mathbb{R}[0,1]$) for each point in the environment using the following equation:

$$f_{near}(dist) = exp[-(dist^2)/2\sigma^2] \quad (1)$$

Where *dist* is the minimum distance to the reference object; σ is the width of the field (dropoff parameter) which is context-dependent. The equation in (1) utilizes a Gaussian for the computation of the field, however other exponential or linear functions could instead be applied depending on the domain requirements. For further information regarding static field computation, we refer the reader to [O’Keefe, 2003].

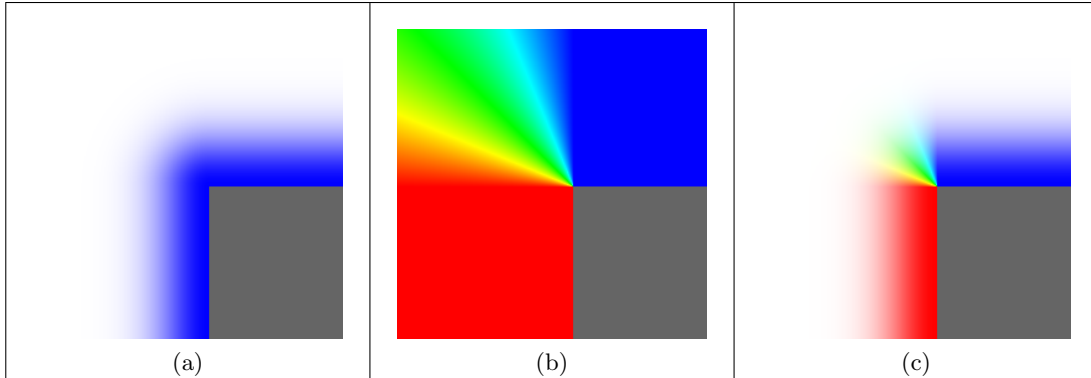


Figure 7.2: Semantic field for “along” (a) near subfield; (b) direction subfield ($90^\circ = \text{red}$, $0^\circ = \text{blue}$); (c) combined field.

7.1.2 Modeling Dynamic Spatial Relations

Modeling DSRs with Local Properties

While appropriate for static relations, the semantic field model, by itself, is not sufficient for representing dynamic spatial relations that involve paths. *Paths* are comprised of a set of points connected by direction vectors that define sequence ordering. Path prepositions include, among others: to, from, along, across, through, toward, past, into, onto, out of, and via. To account for paths in the spatial representation of prepositions, our approach employs multiple methods. The primary method modifies the traditional semantic field model with the addition of a weighted vector field at each point in the environment. As an example, the preposition “along” denotes not only proximity, but also a path parallel to the border of a reference object. Thus, in our proposed model, the semantic field for *along* contains not only weights for each point in the environment to encapsulate proximity, but also weighted direction vectors at each point to encapsulate optimal path direction. Among these direction vectors, those that coincide with the meaning of the relation are favored (in this example, those more parallel to the reference object have higher weight). By multiplying the weights of these two subfields

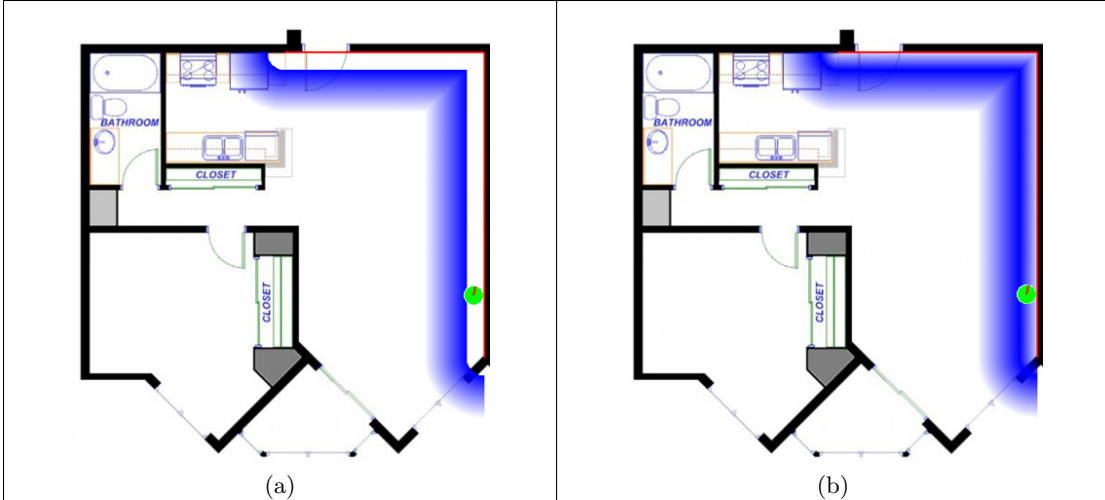


Figure 7.3: Two example alterations to semantic field due to user feedback statement “Move a little away from the wall”. (a) Field moved entirely; (b) Field mean shifted.

together (proximity and path direction) at each point in the environment, we are able to produce the semantic field for the dynamic spatial relation *along* (see Figure 7.2).

The advantage of modeling spatial relations as pdfs, as opposed to using classification-based methods (e.g., [Kollar et al., 2010]), is that generating robot action plans for instruction following is as simple as sampling the pdf, which can be used to find solution paths incrementally (one path segment at a time). In other words, there is no need to search the action space (randomly or exhaustively) to find appropriate solutions by classifying candidate paths as a whole, which may be prohibitive in time-complexity. Furthermore, user teaching, feedback, and refinement of the robot task execution plan can easily be incorporated as an alteration of the pdf. For example, the feedback statement “Move a little away from the wall” could alter the semantic field of the task by attributing higher weight to points further from the wall from the robot’s current location; for example, by shifting the entire field over, or by simply shifting the mean of the field. Figure 7.3 illustrates these two forms of field alterations for the task “Walk along the wall”.

While it is true that some dynamic spatial relations can be modeled by specialized semantic fields that capture optimal path direction at a local level (e.g., along, toward, up, down, etc.), many path prepositions require the existence of certain characteristics achieved at a global level in order to satisfy their meaning.

Modeling DSRs with Global Properties

To represent DSRs with global constraints, our approach identifies four classical AI conditions that each DSR may subscribe to, they are: 1) *pre-condition*, 2) *post-condition*, 3) *continuing-condition*, and 4) *intermediate-condition*. This type of condition-based approach to modeling path prepositions is based on findings in linguistics and cognitive science research on the constraint-based meanings of path prepositions [Bohnenmeyer, 2003; Landau and Jackendoff, 1993]. The methodology is akin to methods developed in the learning by imitation community for task and verb modeling (e.g., [Hewlett et al., 2011; Nicolescu and Matarić, 2005; Pardowitz et al., 2007]). A unique characteristic of our approach, however, is that each condition is represented (typically) by either a semantic field, or by another DSR (which is in turn represented by semantic fields). In the representation, each DSR may have none, one, or multiple of each of the four conditions.

The four conditions enumerated were developed to operate over paths. Formally, a *path* is defined as an ordered set of points (i.e. path $P = \{p_0, p_1, p_2, \dots, p_n\}$), connected by (implicit) direction vectors (from p_i to p_{i+1}). Pre- and post-conditions must be satisfied for the start (p_0) and end (p_n) of the path, respectively. Intermediate-conditions must be satisfied for at least one point in the path, and continuing-conditions must be satisfied for all points in the path. Following this methodology, our condition-based DSR representation may be used both for path classification (e.g., during task learning by demonstration), discussed below, and for path generation (e.g., during robot task execution planning), discussed in Section 7.4.

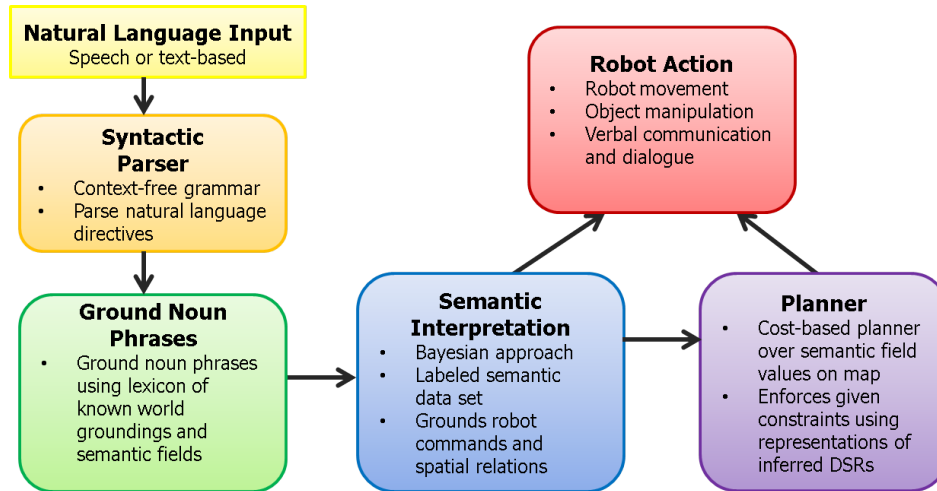


Figure 7.4: Robot software architecture system modules.

7.2 Robot Software Architecture

Our robot software architecture contains five system modules that enable the interpretation of natural language instructions, from speech or text-based input, and translation into agent execution. They include: the syntactic parser, noun phrase (NP) grounding, semantic interpretation, planning, and action modules. Figure 7.4 shows a diagram of the software architecture and system module connections. The following sections discuss the primary modules in detail.

7.2.1 Syntactic Parser

Natural language instructions are received by the syntactic parser as textual input. The text string may be provided by a speech recognizer (e.g., [Nuance, 2013]) or keyboard-based input. While both methods have been implemented with our system, we focus this discussion on well-formed English sentences provided via keyboard input.

The first step of the syntactic parser is to extract the part-of-speech (POS) tags from the natural language text string; these tags identify words in the input as nouns (‘N’),

verbs ('V'), adjectives ('A'), determiners ('Det'), etc. Our system uses the Stanford NLP Parser [Klein and Manning, 2003] for extracting the base POS tags for all words, except for the prepositions ('P'), which are instead identified using a manually created lexicon for single and multi-word prepositions (e.g., "to", "away from", "in line with").

Our system does not attempt to provide a solution for natural language processing in the general case, but instead focuses on directives, and more specifically, on natural language English instructions involving spatial language. To parse these instructions, a phrase structure grammar is utilized. Following are the constituency rules:

$$S \rightarrow V P^* NP$$
$$N' \rightarrow (\text{Det}) A^* N^+$$
$$NP \rightarrow N'$$
$$NP \rightarrow N' P^+ NP$$
$$NP \rightarrow NP \text{ and } NP$$

Here, S defines a valid sentence, NP a noun phrase, and N' a terminal noun phrase. It is important to note that the grammar presented, although limited, is capable of parsing spatial language sentences that do not contain prepositions (e.g., "Enter the room"), those with multiple prepositions (e.g., "Come up on over here"), as well as partial parses of well-formed English sentences (e.g., "PR2, can you please *wait at the counter by the entryway*, thanks").

7.2.2 Grounding Noun Phrases

After the noun phrases in the natural language input are identified by the syntactic parser, the cognitive system attempts to ground the NPs in its representation of the world. Due to the hierarchical nature of NPs, the grounding process is recursive; it first attempts to ground any child NPs before expanding to ground root NPs. To perform this grounding procedure, the nouns in the NP are first checked against the system's

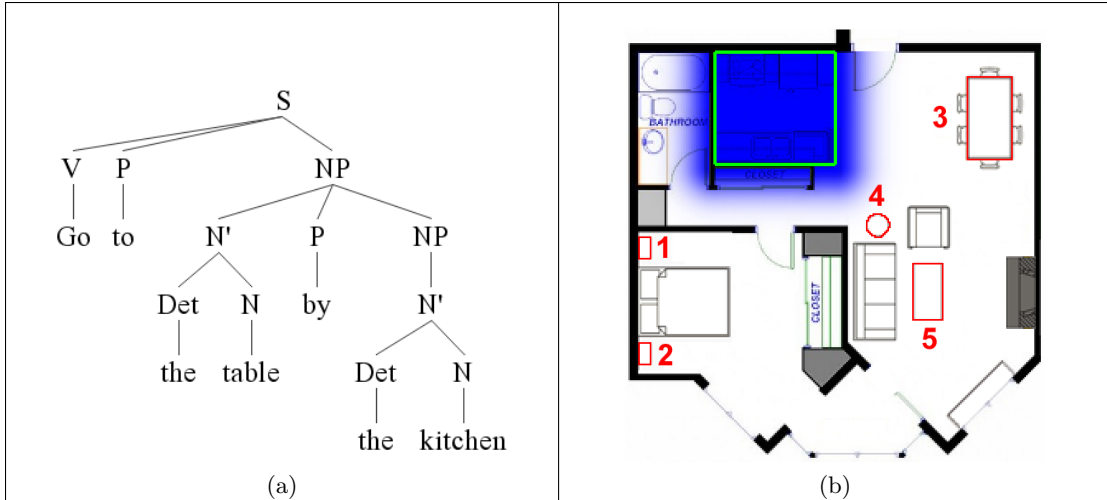


Figure 7.5: (a) Parse tree for “Go to the table by the kitchen”; (b) Semantic field for ‘near kitchen’ with candidate tables.

Table 7.1: Semantic field values of candidate groundings for NP “the table”

Candidate Ground	$\log(\text{Semantic Field Value})$
1	-13.60
2	-46.77
3	-27.67
4	-5.92
5	-28.51

Note: Log semantic field values are reported. Optimal grounding highlighted in bold.

knowledge base of labels for grounds (e.g., objects, rooms, etc.) in the world. These labels are domain-dependent and can either be learned online or, as in our system, loaded from a file along with additional world details, including: a map of the environment, object properties, and the locations of known objects in the map.

If the knowledge base is unable to find a matching label, the grounding process fails, at which point the system may prompt the user for additional information and/or clarification. If a single match is found, the NP is successfully grounded. Lastly, if

multiple matches are found, the system relies on higher-level NPs (for a child NP), or the user (for a root NP), for disambiguation.

In our methodology, disambiguation of multiple matches for a child NP is accomplished in two steps: 1) the semantic field for the prepositional phrase of the child NP's root NP is computed, and 2) each of the candidate grounds are evaluated against the computed semantic field to find the optimal match for the NP.

To illustrate this probabilistic, semantic field-based grounding procedure, consider the instruction “Go to the table by the kitchen”. First, the syntactic parse of the input is obtained (see Figure 7.5 (a)), yielding a single root NP with two children NPs (“the table” and “the kitchen”). In our example world, there is a single ground match for “the kitchen”, but there are five possible groundings for “the table”. To disambiguate among the five candidate groundings, the semantic field for *near* (determined by the use of the preposition “by” in the root NP) is computed for the reference object (i.e., the ground match for the NP “the kitchen”). The field values at each of the candidate ground locations are then evaluated, and the candidate with the highest value is returned as the optimal (most likely) ground match for the NP “the table”. Figure 7.5 (b) shows the semantic field for *near the kitchen* in the example world along with the candidate groundings for “the table”; Table 7.1 lists the field values computed for each of the candidates, for reference.

After all of the NPs in the natural language input have been successfully grounded to known items in the world, the system proceeds to interpret the semantics of the instruction for appropriate robot command execution. If instead grounding errors are encountered in the input, the approach raises the appropriate grounding flags to be handled by higher level reasoning methods and/or human-robot dialogue resolution procedures (see Section 7.5.2).

7.2.3 Semantic Interpreter

Our methodology employs a probabilistic approach to interpreting the semantics of the natural language input. Specifically, the problem statement for the semantic interpretation module is to infer the most likely *command* type, *path* type, and *static spatial relation*, given the observations. The system considers five observations in total, determined by the syntactic parser and grounding modules, including: the verb, the number of NP parameters, the figure type, the reference object type, and the preposition used, if any, in the sentence root.

The command types are domain-dependent, and may include, for example, robot movement, object manipulation, speech production, learned tasks, etc. In evaluating the feasibility of our methodology, our system focuses on two command types: robot movement (translation), and robot orientation. In instructing these types of movement commands, users often utilize spatial relations as opposed to precise quantitative descriptions [Carlson and Hill, 2009]. Therefore, inference of these underlying dynamic and static spatial relations is necessary for correct interpretation of the command. This is especially evident in instructions where path prepositions are not specified (e.g., “Enter the room” vs. “Go into the room”). Static relations are inferred as part of the path specification. For example, the path for to, as described earlier, relies on a static spatial relation to determine the termination condition (e.g., *at* for “to”, *in* for “into”, *out* for “out of”).

The Bayesian inference method utilized by our system is Naïve Bayes, however our methodology allows the use of any probabilistic inference method, leaving the choice up to the system designer. Following is the formula used to perform command inference in our system:

$$\operatorname{argmax}_C P(C|o_1, o_2, \dots, o_n) = \frac{1}{Z} P(C) \prod_i^N P(o_i|C) \tag{2}$$

Where the $N = 5$ observations are the same as those previously listed. The likelihood and prior probabilities are calculated from a database of labeled training data, which could be provided to the system *a priori* and/or gathered incrementally through interaction with the user. The inference of path type and static relations is achieved similarly, with the addition of the inferred command and path type (for static inference) as observations.

7.2.4 Planning

Once the semantic interpreter has inferred the instruction parameters (i.e. the command type, path type, and static spatial relation), the planning module attempts to find a solution for the robot given these command specifications, as well as any other constraints indicated by the user. Constraints are specified to the system the same as instructions, through natural language, and thus their grounding and semantic interpretation are also equivalent.

The A* path planning algorithm is used in our system to find the minimum cost solution for robot action given the command and constraint specifications, which is then passed on to the action module for task execution. In the simplest case, constraints are handled by the planner through modification of the A* cost function. For example, the constraint “Stay away from the TV set” would apply the semantic field of the inferred static relation away from (attached to the reference object) to every point in the environment, and thus points with lower field values would subsequently have higher cost during A* search. More complex constraints would require the planner to segment the search into multiple steps to achieve intermediate goals (see Section 7.4).

7.3 Modeling DSR Representations of “To”, “Through”, and “Around”

7.3.1 “To” Representation

To illustrate our approach for modeling DSRs, consider the path preposition “to”. From linguistics literature, we understand that the path specified by “to” terminates *at* the reference region [Landau and Jackendoff, 1993]. As a result, the DSR representation for “to” in our approach has a single post-condition containing the semantic field for the static spatial relation *at*:

$$to(x) = \begin{cases} \text{pre-condition:} & - \\ \text{cont-condition:} & - \\ \text{int-condition:} & - \\ \text{post-condition:} & at(x) \end{cases} \quad (3)$$

Note that because *at* is represented by a semantic field, it does not return a truth value. Instead, $at(x)$ is a function from points to weight values ($\mathbb{R}[0,1]$). Following is an example semantic field equation for *at*:

$$at(x)(p) = exp[-(dist(x,p)^2)/2\sigma^2] \quad (4)$$

Where $dist(x,p)$ returns the minimum distance between the reference object x and point p ; σ is the width of the field (dropoff parameter) which is context-dependent. By representing conditions as semantic fields, our approach facilitates probabilistic reasoning over paths; an essential quality for path classification, grounding, and generation. As an example, Figure 7.6 shows two sample paths with classification results for the phrase “to the kitchen”, while also displaying the field for $at(the\ kitchen)$. The path

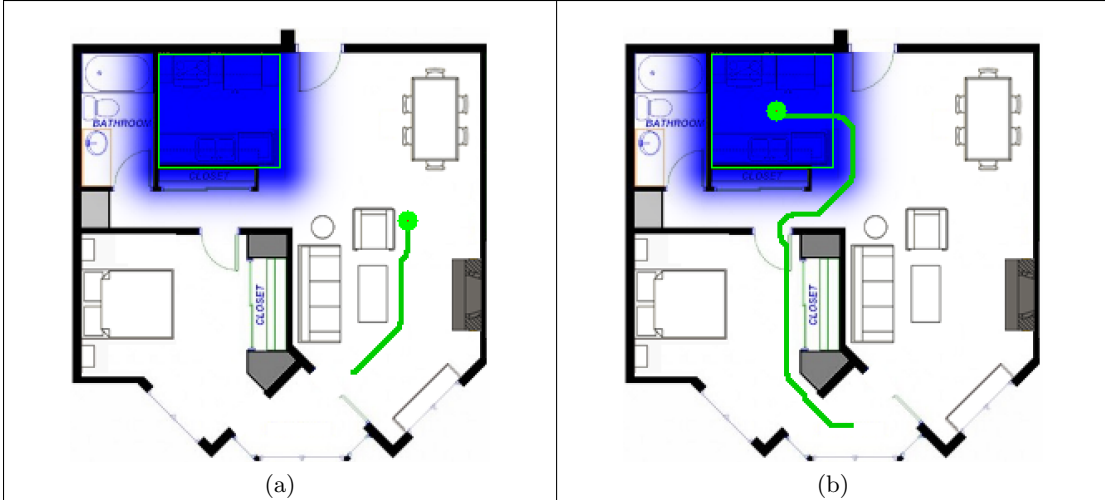


Figure 7.6: Two example paths for “to the kitchen”. (a) Path value = 1.8×10^{-10} ; (b) Path value = 1.0. Note: σ = robot width $\times 2.5$.

values reported correspond to the *at* semantic field values for the path end points (i.e. the post-condition for *to*). As is evident by the results, one path is more acceptable than the other in capturing the meaning of the stated prepositional phrase.

DSRs closely related to that of “to” also have similar representations. For example, the representation for “from” is the reverse of that for “to”, with the *at* field instead being set as a pre-condition. Additionally, the DSR representations for “into”, “onto”, and “out of” are all special cases of *to*, with the *at* field post-condition being replaced by the semantic fields for *in*, *on*, and *out*, respectively.

It is important to note that the DSR representation for “to” is versatile: although *at* is listed as the default post-condition, this determination may change based on context. As an example, consider the phrase “Stand beside the bed”. Here, the (implicit) path relation is *to* and the static relation is *beside*. Hence, the post-condition for *to* would instead be set to the semantic field for *beside*. This substitution is appropriately handled in our methodology by the semantic interpretation module (discussed in Section 7.2.3), which infers path and static relations probabilistically given the natural language input.

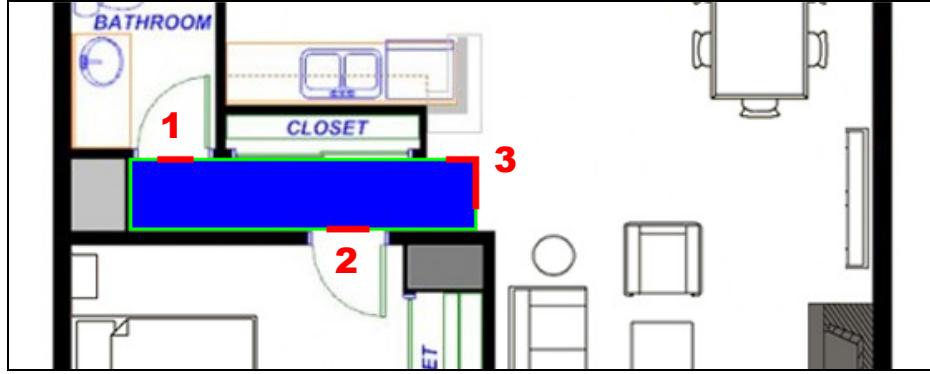


Figure 7.7: Topology of hallway in 2D home environment showing three (configuration space) entrance boundaries.

Table 7.2: DSR Representations for “Through”

Conditions	$through_1(x)$	$through_2(x)$	$through_3(x)$
Pre-	-	$at(B_i(x))$	-
Continuing-	-	$in(x)$	$in(x), along(x)$
Intermediate-	$in(x)$	-	-
Post-	-	$at(B_{j \neq i}(x))$	-

7.3.2 “Through” Representation

The path preposition “through” has a few different semantic interpretations according to the linguistics literature. Therefore, in our approach we developed separate DSR representations for each of them. The most general interpretation asserts that “through” specifies a path with at least one point *in* the reference object [Landau and Jackendoff, 1993]. This definition for “through” can most aptly be characterized by a global DSR representation with a single intermediate-condition containing the semantic field for *in* (see $through_1$ in Table 7.2).

The remaining two interpretations considered by our methodology are both special cases of the first, more general, definition. The second interpretation depends on the topology of the reference object in that it requires that the start and end points, respectively, be coincident with boundaries at separate ends of the reference object [Talmy, 2005]. This definition imposes a path traversing the inside of the reference object, end-to-end. Example uses include “Go through the doorway” and “Walk through the tunnel”. To correctly model this interpretation, the topology (i.e. boundary connectivity) of the reference object must first be determined. While the implementation may vary according to the domain, determining the discrete entrance boundaries for a particular reference object is fairly straightforward in 2D/3D by evaluating edge connectivity. As an example, Figure 7.7 shows the extracted topology of a hallway reference object in a simulated 2D home environment, displaying three separate entrance boundaries. Using the extracted topology, this second definition for “through” can be represented by a DSR with pre- and post-conditions each specifying points at different entrance boundaries, and with a single continuing-condition set to the semantic field for *in*. Table 7.2 presents this representation as *through*₂, with $B_i(x)$ representing the set of all points at entrance boundary i of reference object x .

The third semantic interpretation for “through” is similar to the second except that the path traverses an unbounded segment of the reference object (i.e. does not terminate at object boundaries) [Talmy, 2005]. Thus, this definition simply imposes a path along the inside of the reference object. The DSR representation for this third definition of “through” contains two continuing-conditions with the semantic fields for *in* and *along*, respectively (see *through*₃ in Table 7.2). The *along* semantic field is used in this representation to promote paths that travel parallel to the major axis of the reference object such as to avoid boundedness. Paths that instead travel parallel to the minor axis would be more appropriate for the DSR for “across” [Landau and Jackendoff, 1993;

Algorithm 7.1 Circumcentric Semantic Field Computation

Require: \mathbf{x} is the reference object, $\mathbf{P} = \{p_0, p_1, \dots, p_n\}$ is a path of n points, $diff_{ideal}$ is the ideal path orientation change relative to \mathbf{x} from the start of \mathbf{P} to the end, and let $\theta_{rel}(x, p)$ be a function that returns the orientation of point p relative to x .

circumcentric($\mathbf{x}, \mathbf{P}, diff_{ideal}$)

```
1:  $diff \leftarrow 0$ 
2: for all  $i \in \{1, \dots, n\}$  do
3:    $diff_{new} \leftarrow \theta_{rel}(x, p_i) - \theta_{rel}(x, p_{i-1})$ 
4:    $diff \leftarrow diff + diff_{new}$ 
5: end for
6: return ( $|diff|/diff_{ideal}$ )
```

Talmy, 2005], whose representation in our approach is very similar to that for “through” albeit with the aforementioned distinction.

7.3.3 “Around” Representation

The path preposition “around”, much like “through”, is polysemous. According to Talmy [2005], “around” denotes a circumcentric path (i.e. curved about a center) that can be either revolutionary or rotational. Both path types are similar, except that revolutionary paths refer to curved figure paths about a central reference object (e.g., “The boat sailed around the island”), whereas rotational paths denote a change in orientation of the figure itself (e.g., “John spun around”) [Landau and Jackendoff, 1993]. In the latter case, the figure can also be thought of as the movement path of a point (or points) within the reference object itself during its rotation, thereby illustrating the similarity between the two path types.

In order to represent the DSR for “around”, our approach makes use of a novel semantic field that was developed to quantify the circumcentric nature of a given path. Specifically, the *circumcentric* semantic field maps paths to weight values; where

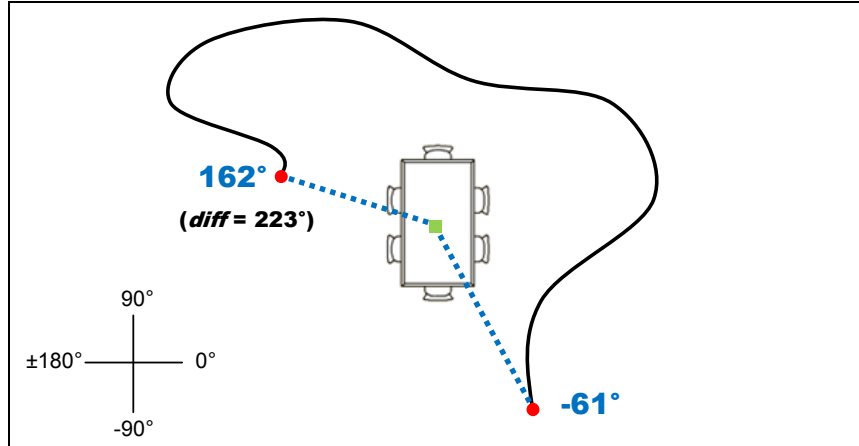


Figure 7.8: Path traveling around a dining table reference object, showing start and end orientations, with resulting *circumcentric* semantic field value = $|223^\circ|/360^\circ = 0.61944$.

Table 7.3: DSR Representations for “Around”

Conditions	$around_{1,2}(x)$	$around_{3,4}(x)$
Pre-	-	-
Continuing-	$\neg in(x)$	$in(x)$
Intermediate-	-	-
Post-	$circumcentric(x, P, [180^\circ_{1,3} 360^\circ_{2,4}])$	

weights are assigned according to the degree to which the orientation of the path changes relative to the center of the specified reference object, from start point to end point. The computation of this field is outlined in Algorithm 7.1, and Figure 7.8 shows an example path with its corresponding *circumcentric* field value for illustration.

In representing the two types of circumcentric paths for “around”, it is important to consider that there are also two termination conditions per path type that are commonly expressed in language: half circle (180°), and full circle (360°). Examples include, as noted by Landau and Jackendoff [1993]: “Go all the way around ...” (360°) vs. “Detour around...” (180°).

Under these specifications, there are a total of four DSR representations for “around”: two revolutionary (half/full circle), and two rotational (half/full circle). Table 7.3 presents all four representations, sequentially labeled as *around*₁₋₄. The continuing-condition for each representation utilizes the semantic field for *in* to express whether or not the figure path is within (i.e. part of) the reference object. Additionally, the post-condition for each representation contains the *circumcentric* semantic field, whose arguments depend on whether the ideal path is a half or full circle.

7.4 Generating Paths for Dynamic Spatial Relations

In searching for robot action solutions for the interpreted command semantics, the planning module must consider not only the inferred command type and spatial relations, but also the *pragmatics* of the natural language instruction. These consist in the unvoiced constraints/specifications that accompany the spoken instructions and further specify the meaning that the speaker intends to convey; which can come from context, prior knowledge, norms, and other factors. Incorporation of specific pragmatic constraints during the planning process is a design decision that depends largely on the domain requirements.

In this section we present the implementations details of the DSR path generation procedures for *to*, *through*, and *around* used in our robot architecture for natural language instruction following. The procedures presented focus on robot movement commands, and illustrate how the representations of DSRs with global properties may be used (combined with the pragmatics of the specific instruction) for the purposes of path generation in robot task planning.

The A* search algorithm is the primary method used for both path planning and topology determination (discussed below) in the planning module; hence, the planner described operates over a discretized representation of the world space.

7.4.1 “To” Path Generation

The DSR representation for “to”, as described in Section 7.3.1, contains a single post-condition with the semantic field for *at*. Therefore, according to this representation, paths that satisfy the relation $to(x)$ are those whose endpoints satisfy the static spatial relation $at(x)$ (defined in (4)).

In the context of our robot architecture, x is the reference object identified during the grounding procedure for the given instruction. In our path generation procedure for *to*, the planner searches for the point p in the free space that maximizes the weight value $at(x)(p)$, which is a real-valued number in the range $[0,1]$, and returns the shortest path to that location from the robot’s current position as a solution.

The pragmatics in our procedure for *to* dictate that if there are multiple points in the free space with maximal weight values (i.e. equal to 1) that are inside the reference object (e.g. in a room), the planner should elect the point furthest from the object edges (i.e. most centrally located) as the end point of the solution path. Figure 7.6 (b) shows an example path for *to* generated under these conditions.

Finally, if the instruction given is part of a command sequence, the pragmatics indicate that expediency in the command solution is favored over optimality. Here, the procedure instead runs A* from the robot’s current location to find the nearest point whose weight value exceeds a certain threshold (e.g., 80% of the maximum weight value in the free space), or whose distance within the reference object exceeds a minimum entry distance (e.g., 1 robot width).

7.4.2 “Through” Path Generation

Paths that satisfy the definition of $through_1$ (see Section 7.3.2) are those with at least one point *in* the reference object x . To generate these types of paths, the planner simply generates paths $into(x)$ using the procedure described above for $to(x)$ and setting the

post-condition to $in(x)$. However, use of the path preposition “through” in directives generally implies the DSR for $through_2$ or $through_3$.

The DSR representation $through_2$ has pre- and post-conditions that require points *at* separate entrance boundaries of x . In planning a solution for $through_2$, the planner first generates a path $into(x)$ to accomplish the pre-condition, and then generates a path $outof(x)$ with the added A* goal constraint that the exit boundary be different than the entrance boundary (determined using the extracted topology of x) to accomplish the post-condition. If at the start of planning the robot is already $in(x)$, the pre-condition is assumed to have been satisfied previously, and the planner subsequently generates a path $outof(x)$ without exit constraints. If there is only one entry boundary for x , the pragmatics dictate the path generation procedure change to that of $through_3$ before accomplishing the post-condition. In addition, if the instruction is part of a command sequence, once a path is generated to within some minimum distance inside x , the pragmatics change the path requirements to $through_1$ and planning for the next command is commenced.

The DSR representation $through_3$ specifies a path $in(x)$ and $along(x)$. To achieve such a path, the planner first generates a path $into(x)$ if the robot is not already $in(x)$, and then runs A* to find the furthest point away from the robot’s location, that is still $in(x)$. The planner then generates a path to this point (staying $in(x)$) to accomplish the $along(x)$ continuing-condition (by default). Alternatively, a path $along(x)$ could be generated using the Voronoi graph of x .

7.4.3 “Around” Path Generation

Considering only the more complex revolutionary cases of “around”, the specific DSR implied depends largely on context, including for example, the topology of the region of space surrounding the reference object. For instance, a lack of 360° connectivity in

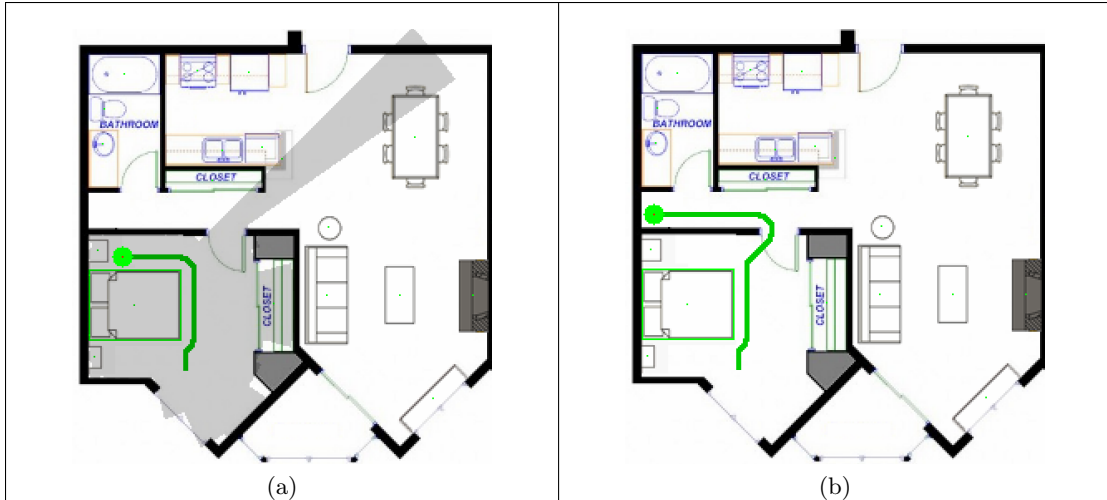


Figure 7.9: Two paths for “Go around the bed” with/without enforcing visibility region. (a) Path value = 0.857; (b) Path value = 0.895.

the region could result in favoring a half circle (180°) interpretation for the DSR. In addition, determination of this topology is required in order to generate appropriate paths for the DSRs of “around”.

To determine if the region of space surrounding reference object x contains 360° connectivity, the planner executes a breadth-first search starting at the robot’s location. As the search progresses, the *circumcentric* semantic field value of each point in the free space is recorded (modified slightly from Algorithm 7.1 with the absolute value of *diff* removed to preserve signed path direction). If meet points are detected from the expanding wavefront of two paths from opposite directions with a combined path orientation difference of 360° , then the topology possesses 360° connectivity.

In the case of *around₂* (full circle), once connectivity is determined the planner first generates a path to the nearest meet point. After which, the planner runs the breath-first search once again to find and plan a path to the nearest meet point on the

opposite side of x , thus completing the 360° loop. If 360° connectivity is not available, the planner instead generates a path to the point with the maximum *circumcentric* field value (to maximize the post-condition weight). Similarly, in the case of *around*₁ (half circle), a path is generated to the point of maximum *circumcentric* field value (stopping at 180°).

Regarding the pragmatics of *around*, consider the instruction “Go around the bed” given by the user to a service robot co-located within the same room. Here, a likely unvoiced constraint is “Stay inside the room”. To incorporate this constraint our planner enforces a global visibility constraint on the free space surrounding the reference object during search. Figure 7.9 highlights the difference between paths generated with and without the visibility constraint, and illustrates its usefulness in practice with end-users. The path values reported correspond to the *circumcentric* semantic field values computed for the paths, with $diff_{ideal} = 180^\circ$ (i.e. the post-condition for *around*₁).

7.5 Parsing Spatial Language Instructions with Figure Objects

7.5.1 Extension for Figure Objects

The syntactic parser of our human-robot interaction framework, as described in Section 7.2.1, is able to interpret a variety of spatial language instructions, including those with hierarchical noun phrases. However, the grammar is only able to capture spatial relationships between an implicit figure object (i.e., the robot) with regards to reference objects specified in a single noun phrase.

To account for the interpretation of explicit figure objects in the instruction semantics, in addition to the previously interpreted reference objects, the phrase structure grammar of our syntactic parser was extended to accept directive sentences with two

Table 7.4: Grammar Constituency Rules for English Directives using Spatial Language

$S \rightarrow V (P_s P_p)^* NP$	$N' \rightarrow (Det) A^* N^+$	$NP \rightarrow N' P_s NP$
$S \rightarrow V NP (P_s P_p)^* NP$	$NP \rightarrow N'$	$NP \rightarrow NP \text{ and } NP$

Part-of-Speech (POS) Tags: V = Verb, P_s = Static Preposition, P_p = Path Preposition, N = Noun, A = Adjective, Det = Determiner

noun phrase parameters. The extended constituency rules, which define valid sentences (S), noun phrases (NP), and terminal noun phrases (N'), are presented in Table 7.4.

The syntactic parser of our framework extracts part-of-speech (POS) tags for all words in the natural language input using the Stanford NLP Parser, with the exception of prepositions, which are instead identified using a manually constructed lexicon. Spatial prepositions in the lexicon are divided into two categories: static (e.g., near, in, on) and path prepositions (e.g., to, from, through), with POS tags of P_s and P_p, respectively. The two categories are not mutually exclusive, as some spatial prepositions are members of both categories (e.g., around). This categorization serves to facilitate the identification of both correct and incorrect preposition usage within noun phrases, as represented by our constituency rules (e.g., “Give me [the ball by the couch]” vs. “Give me [the ball to the couch]”).

7.5.2 Pruning Multiple Parses of a Single Instruction

Static prepositions are often used to express path relations in natural language directive instructions, in substitution of semantically related path prepositions (e.g., the use of “in” instead of “into”; “on” instead of “onto”). This characteristic of natural language often results in the generation of multiple candidate parses for the given directive instruction, each with differing semantics (typically). In these cases, the optimal parse (i.e., the most likely interpretation of the instruction) is determined by evaluating each

Table 7.5: Possible Flags Raised during Parse Pruning

Flag Type	Description
No Ground	No object/label association found
Low Probable Ground	Object/label association found but with low probability for semantic match of prepositional phrase
Multiple Ground	Multiple candidate groundings for NP
Parameter Count	Missing parameter for command
Parameter Type	Figure and/or reference object type mismatch for command and/or prepositional phrase

candidate parse according to both: 1) the resulting parse semantics, and 2) the context of the current environment.

As an example, under the grammar described above, the phrase “Put the cup on the bookcase into the kitchen” has only a single valid parse. In contrast, the phrase “Put the cup on the bookcase in the kitchen” has three possible valid parses, listed below:

- (1) [V Put] [NP the cup] [P_s on] [NP the bookcase in the kitchen]
- (2) [V Put] [NP the cup on the bookcase] [P_s in] [NP the kitchen]
- (3) [V Put] [NP the cup on the bookcase in the kitchen]

In evaluating each candidate parse, our methodology first attempts to ground the NPs of the parse with known objects in the world, and if successful, proceeds to infer the semantics of the instruction as a whole (command, DSR, and static relation), as described in Section 7.2.3. If an error is encountered during the grounding process, or during parameter validation after semantic interpretation of the parse is completed, a flag is raised. Consequently, candidate parses that raise flags are weighted as less likely than parses that do not raise flags. Example flags for grounding and command parameter errors, along with their descriptions, are listed in Table 7.5.

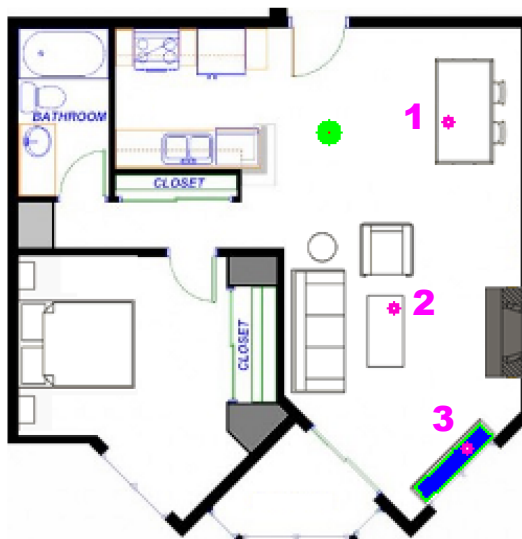


Figure 7.10: Apartment environment with cup locations and $on(\text{bookcase})$ semantic field shown.

If among the candidates a single parse emerges without errors, it is considered to be the optimal parse and is subsequently used for robot task planning and execution. If however, multiple parses are found equally likely, or if all parses raise flags, the user is needed to provide additional clarification before robot task planning can occur.

To illustrate the parse pruning procedure further, consider the three candidate parses listed above in the context of the environment shown in Figure 7.10. Candidates (1) and (3) both fail due to NP grounding errors (low probable ground flag), as the environment does not contain a bookcase in the kitchen. Parse candidate (2) succeeds without errors, as there is a single cup sitting on the bookcase as determined by the probabilistic semantic field grounding procedure (i.e., no multiple ground errors, even with three known cups in the environment), a single ground match for “the kitchen”, and two correctly typed figure (Mobile object) and reference object (Room) parameters for the inferred command of object movement (as determined by the semantic interpreter).

Thus, in this example, parse (2) is chosen as the optimal parse to be used for subsequent robot task planning and execution.

7.6 Object Pick-and-Place Movement Planning

Once the optimal parse is determined (discussed in the previous section) and the results of the semantic interpreter indicate a user instruction of object movement, the robot must first plan to gain possession of the specified figure object, and if applicable, proceed to plan an appropriate placement for the object meeting the requirements of the specified spatial relations of the instruction. This section describes our methodology for accomplishing both of these subtasks using a combination of semantic and pragmatic fields, in addition to taking into account the inherent uncertainty in object grasping in real world scenarios with the use of domain-dependent probabilistic robot-object grasp fields.

7.6.1 Object Pick Up Planning with Grasp Fields

For mobile service robots with onboard manipulators, satisfying an object placement request typically involves first gaining possession of the object in question by picking it up. To address this task, our approach employs a domain-dependent robot-object grasp field centered on the figure object and computed for all points in the environment. This grasp field is analogous to a probability density function, wherein every point in the environment is assigned a weight value in the range $[0, 1]$ that approximates the probability of success of grasping the object with the robot base positioned at that point. The field is domain-dependent as it may incorporate robot characteristics such as arm reach distance, as well as object attributes such as size, weight distribution, handles, etc. These types of fields can either be learned from a corpus of robot-object grasp attempts, or approximated with the use of a general-purpose proximity field based

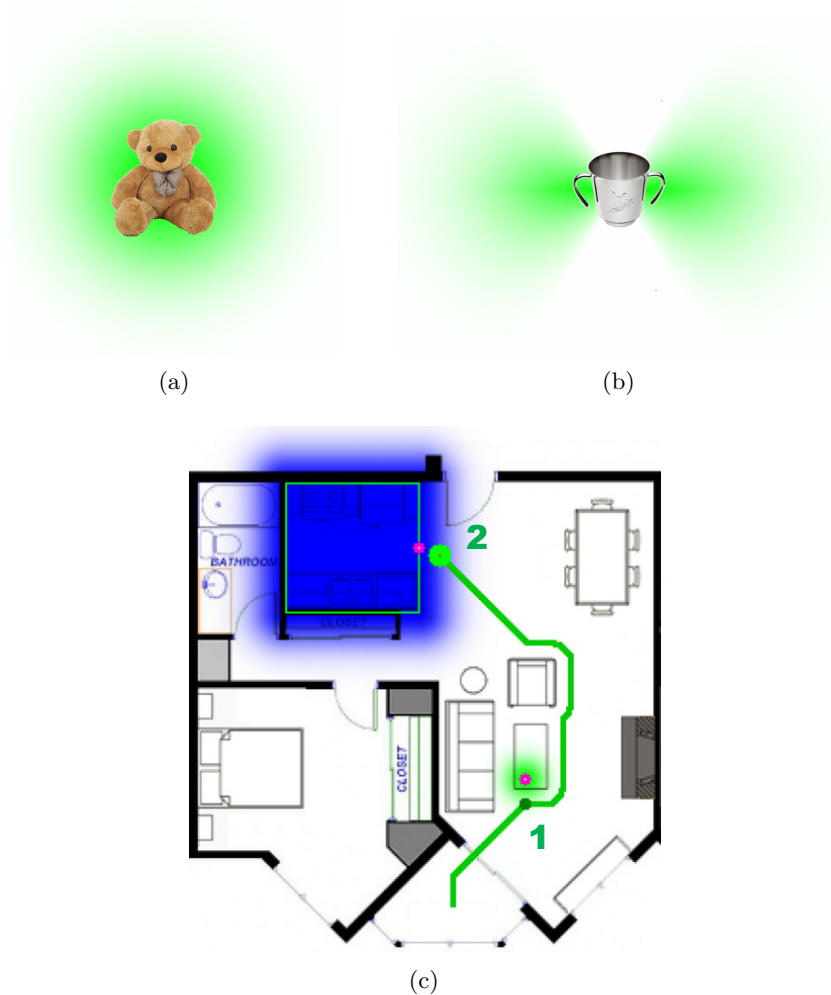


Figure 7.11: (a) Proximity-based grasp field for teddy bear object; (b) Orientation-based grasp field for cup with handles; (c) Task solution for “Take the cup to the kitchen” with grasp field and semantic field shown at the pick (1) and place (2) locations, respectively.

on the robot’s grasp radius. Figure 7.11 shows two example grasp fields for household objects; one proximity-based and the other specific to object orientation.

Once the grasp field of the figure object is computed for all points in the environment, the A* search algorithm is utilized to find a solution path for robot movement to the point of maximum grasp potential in the free space, where the robot can attempt to pick up the object using pre-defined grasp behaviors during task execution.

7.6.2 Object Placement Planning with Semantic Fields

To accomplish the task of object placement, our methodology plans a solution for robot action by first computing the semantic field corresponding to the spatial relation of the instruction (inferred during semantic interpretation) with respect to the reference object, over all points in the robot’s workspace; where the workspace of the robot is defined as all points in the environment reachable by the robot’s end effector.

Once the semantic field is computed, suitable placement points are identified utilizing the maximum field value recorded in the workspace (e.g., all points with weight values 90% of maximum). The point in the environment where the robot base should be positioned as a pre-condition for the object placement action is determined in two steps: 1) candidate robot base target points are identified using the brushfire algorithm, starting at the suitable placement points as initial positions and expanding until reaching the robot’s maximum reach distance, and 2) A* search is run from the robot’s current location until reaching the closest base target point.

The robot movement plan returned by A* is then utilized during task execution for the robot to position itself before attempting to set the object down at the placement point associated with the chosen base target point. Figure 7.11 (c) displays an example solution plan for the object movement instruction “Take the cup to the kitchen”, showing both the grasp field used during pick up planning, and the semantic field for *at*(the kitchen) computed during placement planning.

7.6.3 Pragmatic Fields for Object Placement Planning

In our approach, the use of semantic fields to guide object placement planning successfully enables appropriate placement of figure objects with respect to the spatial relations of the natural language instruction. However, this method only captures the explicit semantics of the instruction, without addressing the *pragmatics* of the task.

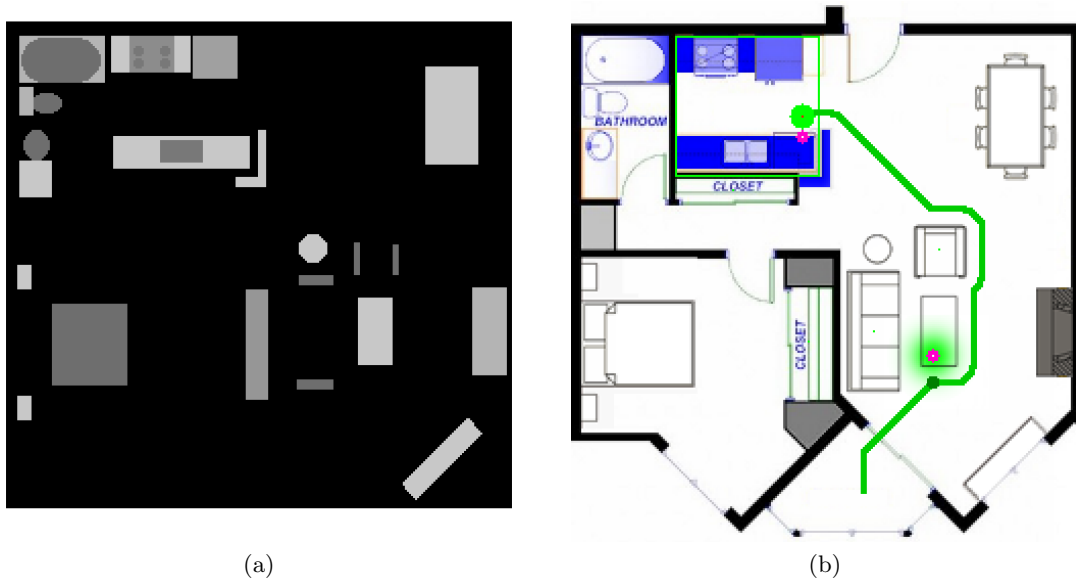


Figure 7.12: (a) Pragmatic field indicating suitable surfaces for object placement (weight values in grayscale); (b) Task solution for “Take the cup to the kitchen” incorporating pragmatic constraints, with combined semantic/pragmatic field shown at object placement location.

Our methodology allows for the incorporation of specific pragmatic constraints with the introduction of spatial *pragmatic fields*. These pragmatic fields are similar to semantic fields in that they assign weight values ($\mathbb{R}[0,1]$) to points in the environment depending on their appropriateness at meeting the goals of the specific spatial pragmatic constraint. As the underlying representations of pragmatic and semantic fields are the same, they can easily be combined for use in robot task planning.

To illustrate the usefulness of incorporating pragmatic fields during planning, consider the example task solution displayed in Figure 7.11 (c) for the instruction “Take the cup to the kitchen”. Here, the robot sets the cup down on the floor at the entryway of the kitchen because, as the cup is located within the kitchen, the explicit semantics of the instruction are satisfied. However, this solution most likely fails to meet the expectations and intentions of the user’s request, which are context-dependent but would

likely indicate placement locations such as: on a counter top, in a cupboard, in the sink, etc. User preferences for such spatial locations, and many others (e.g., on surfaces, away from surface edges, away from obstacles, in drawers), can all be incorporated by multiplying their associated pragmatic fields together with the previously calculated semantic fields for final placement planning. Figure 7.12 (a) shows an example pragmatic field for object placement preference on surfaces (with weighted surfaces), and Figure 7.12 (b) illustrates the use of this field in our described robot task example.

7.7 Summary

This chapter presented our methodology for autonomous service robots to receive and interpret natural language instructions involving spatial relations from non-expert users. Contributions included: the design and implementation details of our robot system modules and software architecture; novel representations for DSRs with local properties and DSRs with global properties that facilitate probabilistic reasoning over paths and that can be applied to both path classification and path generation scenarios; example representations for the DSRs of “to”, “through”, and “around”; implementation details of the path generation procedures utilized by our system for these three DSRs; and discussion of relevant pragmatic constraints along with planning methods developed to address these constraints in multi-step robot execution planning of instruction sequences.

Our approach is capable of addressing both the semantic and pragmatic properties of object movement-oriented natural language instructions, and in particular, proposes a novel computational field representation for the incorporation of spatial pragmatic constraints in mobile manipulation task planning. The design and implementation details of our methodology were presented, including the grammar utilized and our procedure for pruning multiple candidate parses based on context.

Chapter 8

Evaluation of Spatial Language-Based HRI Methodology

This chapter presents an evaluation of our methodology for autonomous service robots to receive and interpret natural language instructions involving spatial relations, with/without natural language constraints and unvoiced pragmatic constraints. Various tests were conducted in both 2D and 3D simulations (using both manually generated and SLAM-based environment maps) to evaluate different components of our methodology, including: the semantic interpretation module, speech recognition, and our approach to following both single instructions and instruction sequences, planning object pick-and-place tasks with figure objects, and planning under unvoiced pragmatic constraints.

8.1 Semantic Inference Accuracy

To evaluate the ability of our approach to follow natural language directives, we first analyzed the effectiveness of the semantic interpretation module to infer the correct command specifications (command type, path type, static relation) given the natural

language input. Our testing domain consisted of a simulated mobile robot operating within a 2D map of a home environment.

A dataset of 128 labeled training examples (each containing a list of observations with correct command specifications), were used in the evaluation of the semantic interpretation module. This dataset included the use of 8 different dynamic spatial relations (path types), 10 separate static spatial relations, 2 commands, and 22 different verbs, each appearing multiple times (and in novel combinations) among the examples. In order to create a training set and a test set for evaluation, the dataset was split into two equal parts, using randomized selection of the examples. A two-fold cross validation was performed on the dataset: the semantic interpreter first utilized the training set to gather probability statistics for the inference process, and was consequently evaluated against the test set. Subsequently the test set and training set were swapped and the inference performance was again evaluated. The results of both evaluations were then averaged to obtain the final inference accuracy results.

The results of the testing show that the semantic interpreter was able to achieve an inference accuracy of 99.2% for commands, 87.8% for paths, and 80.7% for static spatial relations. Table 8.1 contains a summary of these results. Given the relatively small size of the data set, the performance of the semantic interpreter is encouraging. Future work will include performing additional tests to confirm whether or not enhancing the sample size, and/or utilizing a more complex probabilistic model (e.g., Bayesian Network), would result in an increase in inference accuracy.

Table 8.1: Inference Accuracy of Semantic Interpretation Module

Inference Variable	Inference Accuracy
Command	99.2%
Path	87.8%
Static Relation	80.7%

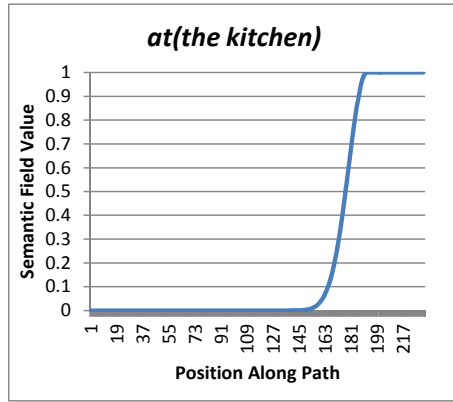
Note. Results of two-fold cross validation of entire semantic dataset with 128 entries.

8.2 Instruction Following Results

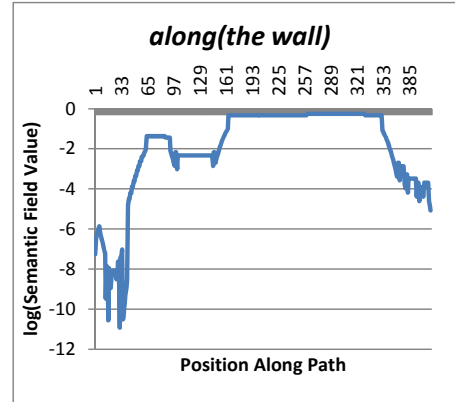
To validate the potential of our methodology towards enabling natural language directive following in service robots, with and without user-specified constraints, we present four example test runs of our system. These examples illustrate the ability of the system to parse natural language input, ground noun phrases, infer command semantics, plan, and execute an appropriate solution while obeying natural language directive constraints.

In the first test run, the command given to the robot was “Go to the room by the entryway”, without constraints. According to the map, the referenced room corresponded to the kitchen, which was correctly grounded by the system using the semantic field for *near*(the entryway). The robot successfully planned and executed the optimal path to the kitchen (see Figure 8.1 (a)). In run #2, the same command was given but with the added constraint “Walk along the wall”, which the robot was also able to account for by utilizing the semantic field values for the dynamic spatial relation *along* in the cost function during the planning process (Figure 8.1 (b)). In runs #3 and #4, the command to the robot was “Stand away from the sink in the bathroom” (differentiating from the kitchen sink), with the addition in run #4 of the constraint “Enter my room” (see Figure 8.1 (c) and (d)).

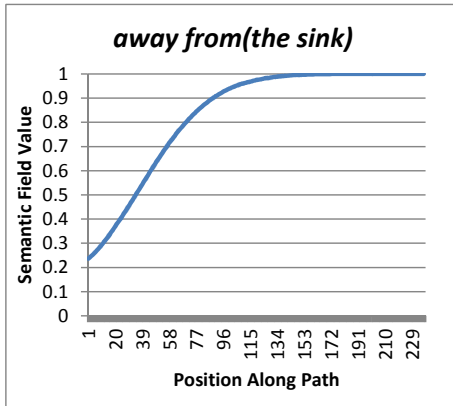
To illustrate the usefulness of the semantic field model towards representing static and dynamic spatial relation primitives for use in path generation and classification, Figure 8.2 shows the progression of the *at*, *along*, *awayfrom*, and *in* semantic field values along the execution paths generated for test runs #1-4, respectively. As demonstrated by the results, the values returned by the semantic fields are highly correlated with the progress made during path execution towards accomplishing the goals of the dynamic spatial relation inferred from the natural language instructions.



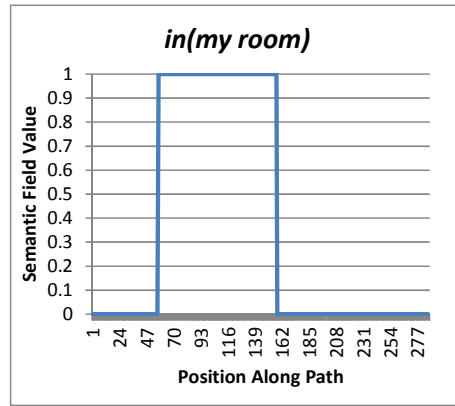
(a)



(b)



(c)



(d)

Figure 8.2: Semantic field values along execution paths in test runs (a) run 1; (b) run 2; (c) run 3; (d) run 4.

Table 8.2: Semantic Inference Results for Instructions and Constraints of Test Runs

Inference Variable	Run 1 instruction	Run 2 constraint	Run 3 instruction	Run 4 constraint
Command	RM	RM	RM	RM
Path	<i>to</i>	<i>along</i>	<i>to</i>	<i>to</i>
Static Relation	<i>at</i>	-	<i>away</i>	<i>in</i>

Note. RM = robot movement command.

As evidenced by the inference results shown in Table 8.2, and all four robot execution paths displayed in Figure 8.1, the system was able to demonstrate its potential by successfully following the natural language directives, with and without constraints, during each of the test runs performed for the purposes of system evaluation.

8.3 Instruction Sequence Following Results

To further evaluate the ability of our robot architecture to follow natural language directives, we conducted two separate test runs of the system testing the robot’s ability to respond to multiple movement commands involving DSRs, provided as a sequence of instructions, both with and without user-specified constraints. The test runs served to evaluate the effectiveness of the semantic interpretation module in inferring the correct command specifications (*command*, *DSR*, *static relation*) given the natural language input, and to demonstrate the DSR path generation capabilities of the system. Our testing domain consisted of a simulated mobile robot operating within a 2D map of a home environment.

The same dataset of 128 labeled training examples from the previous evaluation was utilized for the probabilistic inference procedure of the semantic interpretation module.

The instruction sequence provided to the robot in both test runs, including the natural language constraints that were specified for the individual instructions, is listed in Table 8.3. The sequence of instructions was identical for both runs, with the exception that the constraints listed were specified to the robot for Test Run #2 only. Hence, in Test Run #1 the robot did not operate under any user-specified constraints for the individual instructions. Constraints were provided to the robot in Test Run #2 to illustrate the flexibility of the path generation procedure to operate under user-specified constraints while also accomplishing the goals of the DSR path specification. The planning module accounts for user-specified constraints by introducing modifications to

Table 8.3: Instruction Sequence Given in Test Runs

Type	Natural Language Instruction
Instruction[1]: <i>Constraint:</i>	Go around the bed <i>Stay close to the bed</i>
Instruction[2]:	Travel through the hallway
Instruction[3]: <i>Constraint:</i>	Go around the dinner table <i>Keep away from the kitchen</i>
Instruction[4]: <i>Constraint:</i>	Stand between the tv and the bookcase <i>Travel between the couch and the coffee table</i>
Instruction[5]: <i>Constraint:</i>	Walk through the kitchen <i>Walk along the wall</i>

Table 8.4: Results of Semantic Inference and Pragmatics for Test Run Instructions

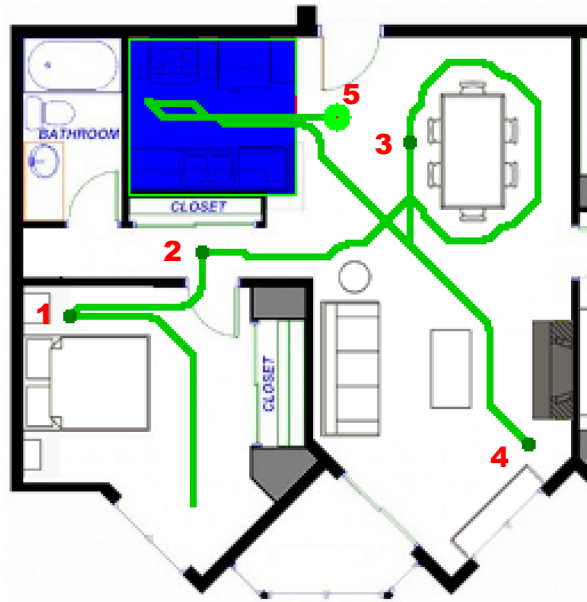
Run #	Semantics	Pragmatics
1	(RM, <i>around</i> , -)	<i>around</i> ₁
2	(RM, <i>through</i> , -)	<i>through</i> ₂ → <i>through</i> ₁
3	(RM, <i>around</i> , -)	<i>around</i> ₂
4	(RM, <i>to</i> , <i>between</i>)	<i>to</i>
5	(RM, <i>through</i> , -)	<i>through</i> ₂ → <i>through</i> ₃ → <i>through</i> ₂

Note. RM = robot movement command.

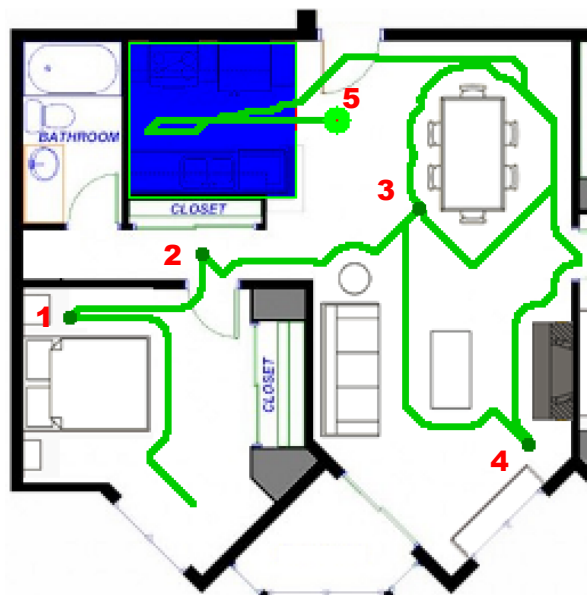
the A* cost function (using the semantic fields of the inferred static relations) during task planning, as detailed in Section 7.2.4.

Results of the inference procedure of the semantic interpretation module, with accompanied pragmatics, for the five instructions given in the instruction sequence for both test runs are provided in Table 8.4. As evidenced by the results, our robot architecture was able to successfully interpret the semantics of the natural language instructions provided by the user during both test runs of the system.

The DSR path generation results for the entire instruction sequence of Test Runs #1 and #2 are provided in Figure 8.3 (a) and (b), respectively. The differences between the paths generated in both test runs highlight the impact of user-specified constraints



(a)



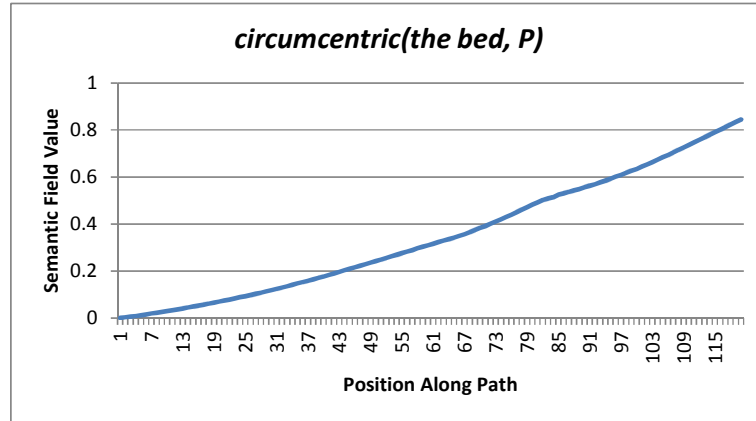
(b)

Figure 8.3: DSR path generation results for entire instruction sequence with and without user-specified constraints. (a) Test Run #1 (no constraints); (b) Test Run #2 (constraints). Note: path endpoints for each instruction are labeled with the instruction number.

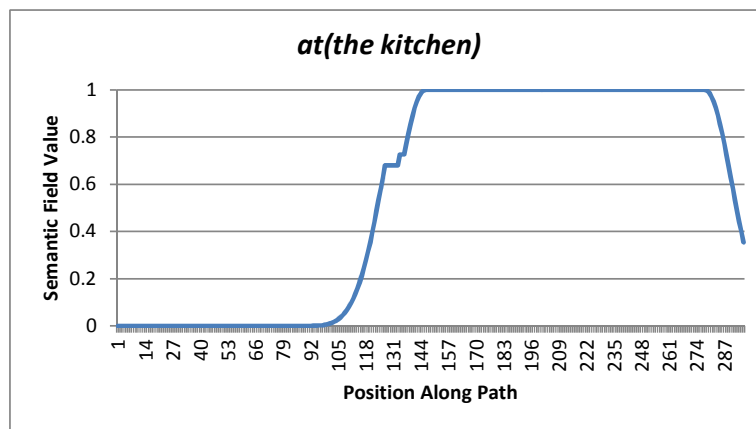
on the resulting robot execution path. For example, in Test Run #1 the robot satisfies the DSR of the first instruction (*around*) by generating and executing the shortest path to the point within the visible region that possesses the maximum *circumcentric* field value among all points considered. In Test Run #2, the robot also generates a path to this point, but due to the user-specified constraint “Stay close to the bed”, the execution path runs along the border of the bed, resulting in a slightly longer path by comparison. This difference in path generation results is also observed for the last instruction in the sequence (“Go through the kitchen”), where in Test Run #2, the robot generates a comparably longer path to the inside of the kitchen by staying close to the edge of the rooms in consideration of the specified constraint “Walk along the wall”.

To illustrate the usefulness of the semantic field model towards representing static and dynamic spatial relation primitives for use in DSR path generation and classification, Figure 8.4 shows the progression of the *circumcentric* and *at* field values along the execution paths generated for instructions 1 and 5, respectively. As demonstrated by the results, the values returned by the semantic fields are highly correlated with the progress made during path execution towards accomplishing the goals of the DSR inferred from the specified natural language instructions.

As evidenced by the semantic inference results shown in Table 8.4, and all robot execution paths for the DSRs of the instruction sequence displayed in Figure 8.3, the robot architecture was able to demonstrate its potential by successfully following the natural language directives, with and without constraints, during each of the test runs performed for the purposes of system evaluation. In addition, the differences observed in the generated DSR paths for both test runs illustrate the capability of our approach to modeling DSRs with global properties in accomplishing natural language instructions in human-robot interaction scenarios under both user-specified constraints as well as unvoiced pragmatic constraints.



(a)



(b)

Figure 8.4: Semantic field values along task execution paths in Test Run #1. (a) *circumcentric* field value along solution path for instruction 1; (b) *at* field value along solution path for instruction 5.

8.4 Speech Recognition Results

As previously mentioned, the syntactic parser module can accept text input from either a speech recognizer or keyboard input. To illustrate the feasibility of our approach to operate with human users in real world environments, we provide the results of our implemented speech recognition module on the 128 natural language instruction training

Table 8.5: Speech Recognition Module Accuracy

Sentence Error Rate (total correct / total sentences)	
Sentence Error Rate	$9/128 = 7.03\%$
Sentence Semantic Error Rate	$6/128 = 4.68\%$
Word Error Rate (substitutions + deletions + insertions / total words)	
Word Error Rate	$(8 + 5 + 2)/686 = 15/686 = 2.18\%$
Word Semantic Error Rate	$(4 + 5 + 2)/686 = 11/686 = 1.6\%$

Semantic error rates exclude errors resulting in semantically equivalent sentences/words.

examples in our test database. Each of the entries was spoken exactly once for analysis, using a headset microphone placed approximately 1 inch from the speaker’s mouth, and with minimal background noise. Table 8.5 presents the accuracy results of the speech recognition module. The speech recognizer used in the module was Nuance’s Dragon NaturallySpeaking [Nuance, 2013].

The low error rate of the speech recognition module observed under our test conditions (low ambient noise and using a user mounted headset microphone), combined with the availability of algorithms to interpret spoken language under various forms of disfluency and repetition (e.g., [Scheutz et al., 2011]), demonstrate the feasibility of obtaining grammatically correct text input from spoken language in real world scenarios for use in our software architecture for service robots.

8.5 Instruction Following Results with SLAM Maps

To demonstrate the generalizability of our approach and its usefulness in practice with real robots in real environments, next we present evaluation results of our robot software architecture using maps of real environments that were generated by real robots

Table 8.6: Instructions Given in Test Runs 5-12

Type[Run #]	Natural Language Instruction
Instruction[5]:	Go to the cafeteria
<i>Constraint:</i>	<i>Walk along the north hallway wall</i>
Instruction[6]:	Go to the cafeteria
<i>Constraint:</i>	<i>Walk along the south hallway wall</i>
Instruction[7]:	Stand away from the desk in my office
<i>Constraint:</i>	<i>Enter the meeting room</i>
Instruction[8]:	Stand between my office and the lab
Instruction[9]:	Relocate to the lounge area next to the lab
Instruction[10]:	Relocate to the lounge area next to the lab
<i>Constraint:</i>	<i>Roll along the central wall</i>
Instruction[11]:	Get to the kitchen
Instruction[12]:	Get to the kitchen
<i>Constraint:</i>	<i>Travel inside the main area</i>

implementing SLAM with onboard laser sensors. We provide the results of spatial language instructions given to a simulated mobile robot within these environments, with and without user-specified natural language constraints, to showcase the ability of our methodology to generate semantic fields, both dynamic and static, to accomplish spatial language tasks in real world task scenarios. The two maps that were used for this additional evaluation were collected from the Radish data set [Howard and Roy, 2003], and consist of a map of a building at the University of Freiburg (FR079), and a map of the interior of the Intel Research Lab in Seattle (intel_lab). The maps were manually annotated to specify landmark locations (e.g., rooms, walls, objects) and were given to the robot *a priori*. In practice, annotation for the robot-generated maps would be accomplished by the user and/or by a qualified technician prior to first use.

Eight additional test runs of our robot architecture were conducted using real world maps generated by robots using onboard laser sensors, as noted above. The natural language instructions, with their associated constraints, given to the robot in test runs #5-12 are provided in Table 8.6. The robot execution paths for each of the test runs,

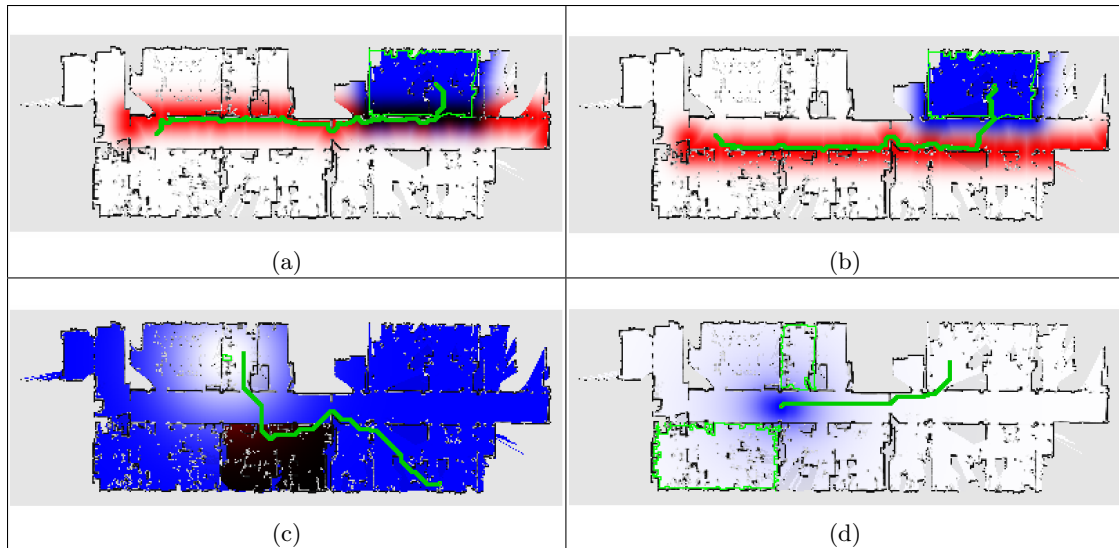


Figure 8.5: Executed paths and semantic fields (command = blue, constraint = red) for test runs (a) run 5; (b) run 6; (c) run 7; (d) run 8.

along with associated semantic fields displayed for reference purposes, are provided in Figures 8.5 and 8.6.

While these results were obtained from simulations in 2D, it is very common for robots operating in real-world environments (such as homes or offices) to utilize a 2D map representation of the environment for localization and spatial task planning. The SLAM maps presented in this section are identical to the maps that would be used by a real robot operating in the actual 3D environments, and the methods employed for spatial language instruction understanding and task following, as presented in this paper, would also be identical. As a last step towards implementing our methodology on a real robot, translation of the discretized plan returned by the planner to continuous robot motor commands (e.g., wheel velocities) is necessary. This translation can be accomplished using a local planner which utilizes the returned A* path to fill a cost map covering the robot’s local environment to determine the wheel velocities that would result in robot movement best following the generated path. This local planner would

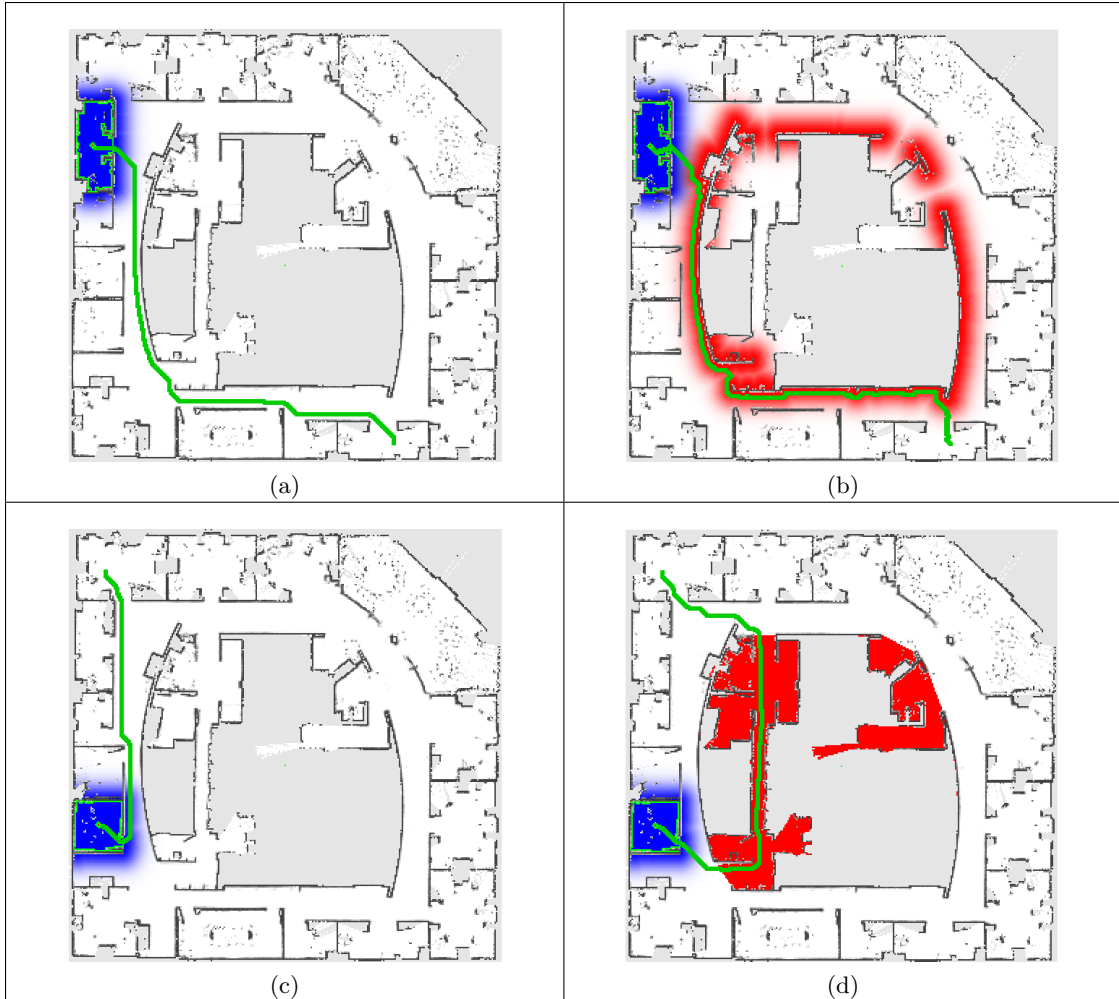


Figure 8.6: Executed paths and semantic fields (command = blue, constraint = red) for test runs (a) run 9; (b) run 10; (c) run 11; (d) run 12.

also be able to respond to dynamic obstacles not represented in the map (e.g., people, objects, etc.) and can be implemented using the dynamic window approach proposed by Fox et al. [1997], for which there is an available ROS package which facilitates its usage in practice [Marder-Eppstein and Perko, 2012].

The evaluation results presented above, regarding speech recognition accuracy and spatial navigation task performance using SLAM maps, demonstrate the feasibility of

our approach for use in practical applications with real robots. The evaluation of the robot software architecture in multiple environments, using both manually-created and robot generated maps, demonstrates the generalizability of the approach and its effectiveness in accomplishing spatial language instruction tasks, with and without user specified constraints, across domains and in novel real world environments.

8.6 Object Pick-and-Place Task Results

To evaluate the ability of our robot architecture to follow natural language directives involving object pick-and-place tasks, we conducted two separate test runs of the system testing the robot’s ability to respond to multiple commands involving object relocation, provided as a sequence of instructions, with and without corresponding user-specified constraints. The test runs served to demonstrate the effectiveness of the semantic interpretation module in inferring the correct command specifications (command, DSR, static relation) given the natural language input, and to demonstrate the path generation capabilities of the system. Our testing domain consisted of a simulated mobile robot operating within a 2D map of a home environment.

A dataset of 189 labeled training examples (each containing a list of observations with correct command specifications), was utilized for the probabilistic inference procedure of the semantic interpretation module. This dataset included the use of 11 different DSRs, 10 separate static spatial relations, 5 commands, and 38 different verbs, each appearing multiple times and in novel combinations among the examples. The instruction sequence provided to the robot in the test runs, including the natural language constraints that were specified for each instruction, is listed in Table 8.7. The two test runs evaluated the same instructions, with the only difference being the addition of the constraints for the second run. The path generation results for the entire instruction sequence of both test runs are provided in Figure 8.7.

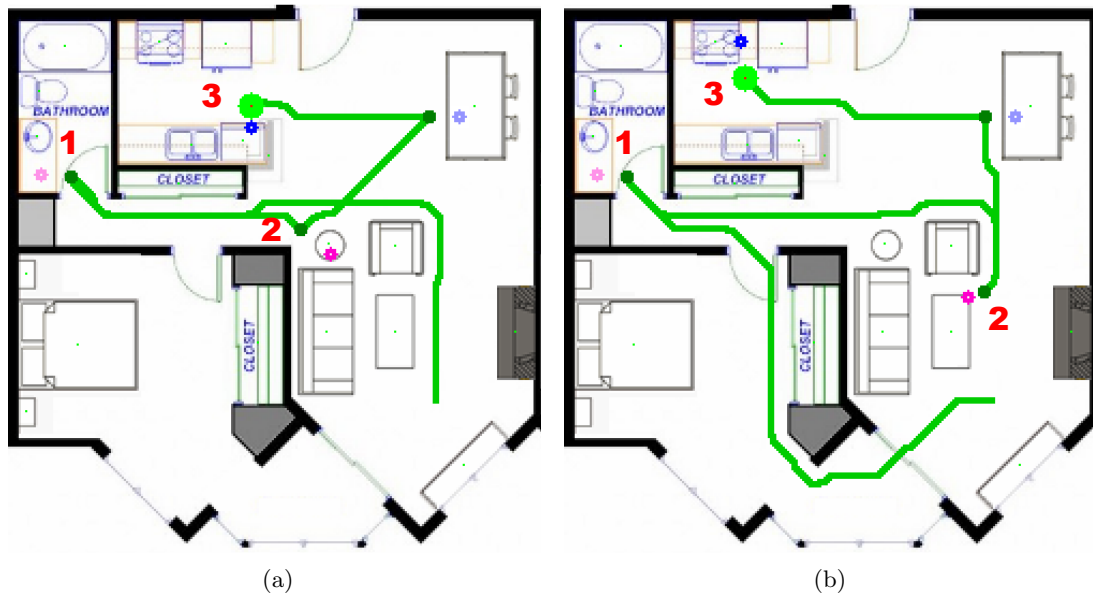


Figure 8.7: Robot execution paths for test runs. (a) Run #1 (no constraints); (b) Run #2 (constraints).

Table 8.7: Instruction Sequence for Test Runs with Corresponding Constraints

Instruction 1	Instruction 2	Instruction 3
Grab the medicine <i>Go inside my room</i>	Drop the medicine in the living room <i>Put the Tylenol between the couch and the tv</i>	Take the olive oil to the kitchen <i>Place the oil near the stove</i>

The path generation procedure for the first test run utilized a combination of spatial semantic fields, and the pragmatic field for surface placement shown in Figure 7.12 (a). For the second test run, the planning procedure remained the same as in the first, except for the addition of pragmatic fields corresponding to the spatial language constraints specified for each instruction. As evidenced by the resulting robot execution paths for the pick-and-place instruction sequence, our spatial language-based HRI framework was able to demonstrate its potential by successfully following the natural language directives, including under user-specified constraints, during each of the test runs performed for the purposes of system evaluation.

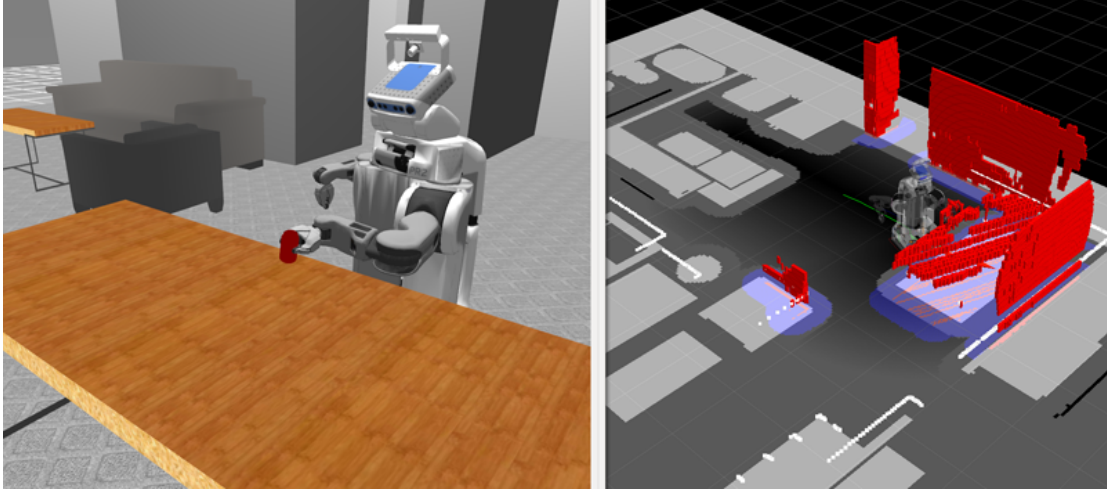


Figure 8.8: PR2 robot executing task for “Put the Coke can on the coffee table” in 3D household environment using Gazebo simulator/ROS framework.

To demonstrate the generalizability of our approach and its usefulness in practice with real robots in real environments, we have implemented our approach in the 3D Gazebo simulator under the ROS framework. Translation of the discretized plan returned by the planner to continuous robot motor commands (e.g., wheel velocities) was accomplished by providing a cost map based on the task solution path to the ROS navigation stack. The local planner employed is able to respond to dynamic obstacles not represented in the map (e.g., people, objects) which facilitates its usage in real world domains. Autonomous pick and place behaviors were also incorporated utilizing software packages available for the PR2 robot. Figure 8.8 shows a snapshot of a successful run of the robot executing a pick-and-place task within the 3D household environment.

8.7 Summary

This chapter presented results from evaluation testing of our methodology for autonomous service robots to receive and interpret natural language instructions involving

spatial relations in a simulated end-to-end system (both in 2D and 3D). Testing included evaluation of: the semantic interpretation module, the speech recognition module, our approach to following both single instructions and instruction sequences with/without natural language constraints and including operation within SLAM generated environment maps, planning object pick-and-place tasks with figure objects, and planning under unvoiced pragmatic constraints. The results presented demonstrate the potential of our methodology for representing dynamic spatial relations, interpreting the semantics of natural language instructions probabilistically, and generating appropriate agent execution plans under user-specified natural language constraints as well as unvoiced pragmatic constraints.

Chapter 9

Spatial Language Discourse with Pragmatics for HRI

This chapter presents a methodology for enabling service robots to interpret spatial language instruction sequences expressed through natural language discourse from non-expert users. In particular, the chapter presents a novel probabilistic algorithm for the automatic extraction of contextually and semantically valid instruction sequences from unconstrained spatial language discourse. Additionally, it presents the design and implementation details of a procedure for reference resolution of anaphoric expressions encountered within the user discourse. Towards application of our human-robot interaction (HRI) methodology on robot platforms in practice with end users, the chapter also discusses a generalized procedure for transfer to physical systems and provides solutions for key pragmatic considerations including the generation of safe robot execution paths for both the robot and people in the environment. The chapter concludes with an evaluation of our spatial language-based HRI framework implemented on a PR2 robot to demonstrate the generalizability and usefulness of our approach in real world applications.

9.1 Interpreting Instructions in Discourse

9.1.1 Probabilistic Extraction of Instruction Sequences

The approach described in Chapter 7 discusses utilization of a phrase structure grammar capable of parsing spatial language directives that instructed a variety of robot tasks, including for example, robot movement commands (e.g., “Go inside the kitchen”), object manipulation/placement commands (e.g., “Put the book on top of the coffee table”), and spatial commands without explicit prepositions (e.g., “Leave the room”). Table 9.1 displays the basic rules of this grammar for illustration purposes; the complete grammar is slightly more complex (see Section 7.5.1). As shown, the non-terminal symbols defined by the constituency rules include those for sentences (S), noun phrases (NP), and terminating noun phrases (N’).

While the grammar presented is capable of capturing many different types of spatial language instructions (including those with hierarchical noun phrases) provided as discrete input, its scope is limited when applied to natural language input taken as a whole. Specifically, it is unable to parse the many non-spatial phrases that users often employ when providing instructions through natural language. In addition, the grammar only allows for a single instruction per sentence, yet in practice, people often sequence multiple instructions together within a single utterance that must therefore be appropriately segmented.

Table 9.1: Grammar for Spatial Language Directives

$S \rightarrow V P^* NP$	$N' \rightarrow (Det) A^* N^+$	$NP \rightarrow N' P NP$
$S \rightarrow V NP P^* NP$	$NP \rightarrow N'$	$NP \rightarrow NP \text{ and } NP$

Part-of-Speech (POS) Tags: V = Verb, P = Preposition, N = Noun, A = Adjective, Det = Determiner

To address these limitations in practice with end users, we have developed a probabilistic parsing procedure capable of extracting a sequence of grammatical instructions (partial parses) from unconstrained natural language input for subsequent robot task planning and execution. Following is an overview of the five steps of the algorithm:

1) *Part-of-Speech (POS) Tag Assignment*: The first step of the algorithm is to determine the POS tags (terminals) for each word of the input text. We use the Stanford NLP Parser [Klein and Manning, 2003] to generate default POS tags; however, because the Stanford parser does not have access to situational context, it occasionally assigns POS tags incorrectly (e.g., “Place” assigned as a noun (N), instead of as a verb (V)). To address this issue, we additionally apply domain-specific POS tags taken from a pre-defined lexicon when there is a disagreement with the default tags, so that the parser may consider both tag options. As a result, there may be multiple assignments generated for a given input text (with exponential growth). In practice, however, there are typically only 1-4 tag assignment arrays for each input, and invalid tag assignments are discarded quickly in the following step due to grammatical incorrectness. For each of the generated tag assignment arrays the algorithm performs steps 2-4; the procedure then concludes with step 5.

2) *Parse Word/POS Tag Array using Grammar*: Given a word/POS tag assignment array, the algorithm proceeds to extract the corresponding high-level tags (non-terminals) for the input as defined by the constituency rules of the grammar (see Table 9.1). The result is that for each word position, there exists a set of non-terminal symbols parsed by the grammar that begin at that index, each with an associated length corresponding to the number of consecutive terminal symbols that serve as constituents for the non-terminal symbol. From this representation, the algorithm only considers non-terminal symbols that denote grammatical sentences (i.e., S).

3) *Find Maximum Probable Sentences*: For each word index of the word/POS tag assignment array, all sentences (symbol S) with maximum length among the available sentences are tested for validity within the context of the environment and the semantics of the inferred instruction. The algorithm only considers sentences of maximum length among those available as a heuristic to avoid evaluating partial sentences unnecessarily. If no sentences exist at the current index, the algorithm moves on to the next index.

Valid sentences are those whose NPs can be grounded uniquely in the world, and whose parameters meet the specifications of the inferred command. An example error would be if the inferred command was [Object Movement] and the grounded NP parameter was [the kitchen], as [the kitchen] is of type [Room] and hence not movable by the robot. In this case a flag would be thrown and the sentence would be deemed invalid. This validation procedure is discussed in detail in Section 7.5.2.

Among the sentences at the current index found to be contextually and semantically valid, the sentence of maximum probability (calculated during the inference process) is chosen as the most likely sentence found at the current index, and the algorithm then skips the word indices covered by the sentence and continues on searching for valid sentences at the next available index.

4) *Form Instruction Sequence Candidate*: All of the valid sentences found within the word/POS tag assignment array (i.e., those with maximum probability at their respective word positions) are then combined to form the optimal instruction sequence candidate for the specific POS tag assignment of the natural language input.

5) *Find Maximum Probable Instruction Sequence*: Once all instruction sequence candidates are gathered, the final instruction sequence returned by the algorithm is that which is of maximum probability among the candidates (determined by multiplying together each of the probabilities of the individual sentences in the sequence). To allow

Table 9.2: Probabilistic Instruction Sequence Extraction Procedure Example

Word Index	0	1	2	3	4	5	6	7	8	9	10	11	12
Word	“PR2	please	go	into	my	room	and	get	me	my	shoes	thank	you”
POS Tag	N	V	V	P	PRP\$	N	CC	V	PRP	PRP\$	N	V	PRP
Non-Terminals	NP ₁	-	S ₄	-	NP ₂	NP ₁	-	S ₄ S ₂	NP ₁	NP ₂	NP ₁	S ₂	NP ₁
Algorithm Steps	<i>No S</i>	<i>No S</i>	<i>Valid S Found len=4</i>	<i>(skip)</i>	<i>(skip)</i>	<i>(skip)</i>	<i>No S</i>	<i>Valid S Found len=4</i>	<i>(skip)</i>	<i>(skip)</i>	<i>(skip)</i>	<i>Invalid S</i>	<i>No S</i>

Note. Input text = “PR2 please go into my room and get me my shoes thank you”. Parser output is shown along with algorithm steps during iteration over the word indices of the POS tag assignment array. Non-terminal symbol lengths are shown in subscripts. Sentences are validated against the context of environment and the inferred command semantics/requirements. Final instruction sequence output by algorithm: {go into my room, get me my shoes}

for fair comparison, candidates are evaluated only against others of equal length (number of sentences), and instruction sequences of greater length are favored.

Table 9.2 illustrates the probabilistic instruction sequence extraction procedure with an example by displaying the word/POS tag assignment array for the input “PR2 please go into my room and get me my shoes thank you”, along with the corresponding parsed non-terminal symbols and resulting algorithm steps. In the example, the procedure finds the following instruction sequence most likely given the natural language input: {Go into my room, Get me my shoes}.

9.1.2 Reference Resolution

In natural language discourse, people often refer to entities, or groundings, which have been previously mentioned or discussed through the use of anaphora. Examples include references to objects (e.g., “it”, “itself”, “this”, “that”) and people (e.g., “he/she”, “him/her”, “him/herself”). In addition, anaphoric expressions typically refer to an

entity introduced by a noun phrase within a recent utterance in the discourse history (usually within one or two past utterances) [Jurafsky and Martin, 2008]. The prevalence of anaphora in natural language discourse necessitates a computational approach for resolving such references for use in real world human-robot interaction scenarios with non-expert users.

In this subsection, we present our approach to resolving anaphoric references to both objects and people in the context of user-guided spatial language discourse. Our reference resolution procedure is similar in nature to those that have been developed previously based on related principles [Carbonell and Brown, 1988; Jurafsky and Martin, 2008], albeit with the distinction of its optimization for use within the framework of our spatial language architecture, and in particular, for its designed integration with our probabilistic instruction sequence extraction procedure (presented in the previous subsection).

At a high level, our procedure for resolving anaphoric references within user discourse can be summarized by the following key concepts: 1) entities represented in the world (e.g., mobile objects, static objects, rooms, people, etc.) are associated internally with numerical identifiers that enable unique identification during the NP grounding process; 2) as these groundings are referenced in the discourse (usually by name) their unique grounding ID numbers are added to a global list of recent references; and 3) upon encountering anaphoric expressions within the discourse, the groundings in the recent references list are used as candidate references in an attempt to uniquely resolve the referential expression to the specific grounding that the user intended to convey.

More specifically, in our approach anaphoric expressions are categorized as either Object References or Human References, depending on whether or not the anaphor encountered refers to a person. In addition, anaphoric references to persons are further categorized by gender (male/female). When adding groundings to the global list of

recent references, an entry pair is made with both the current utterance index and the grounding ID. If a prior entry is found with the same grounding ID, it is removed in favor of the new entry. The utterance index is incremented after every utterance spoken during discourse, and it is included in the global list to enforce the consideration of only the references expressed within the most recent utterances (in our implementation we utilize a history size of three utterances).

In resolving an anaphoric expression, only recent reference groundings with matching type (object vs. human, male vs. female) are allowed as candidates, and the list of candidate groundings is prioritized with the most recent references at the top. During the grounding process, child NPs that can be directly grounded (i.e., not anaphoric) are added to the current list of recent references; alternatively, child NPs that contain anaphora instead merge their candidates list with those of sibling NPs to form one combined candidates list for the parent NP. Once all child NPs are processed and either the figure and/or reference object NP parameters of the spatial language instruction contain anaphora, the command semantics are evaluated for each of the possible candidate groundings until the first successful assignment is found. This greedy approach to resolving the reference is reasonable under the assumption that the list of candidates is ordered with the most likely candidates on top (the most recent are set as a best estimate).

As previously mentioned, our procedure for anaphora resolution was designed to be well integrated with our probabilistic instruction sequence extraction procedure. This integration is actually a crucial necessity, as determining the optimal instruction sequence for a given utterance containing anaphora depends entirely on accurate reference resolution. Furthermore, if multiple POS tag assignments exist for the given utterance, separate reference lists must be concurrently maintained and adjusted according to the evolving context of the different threads of possible discourse under consideration. Yet,

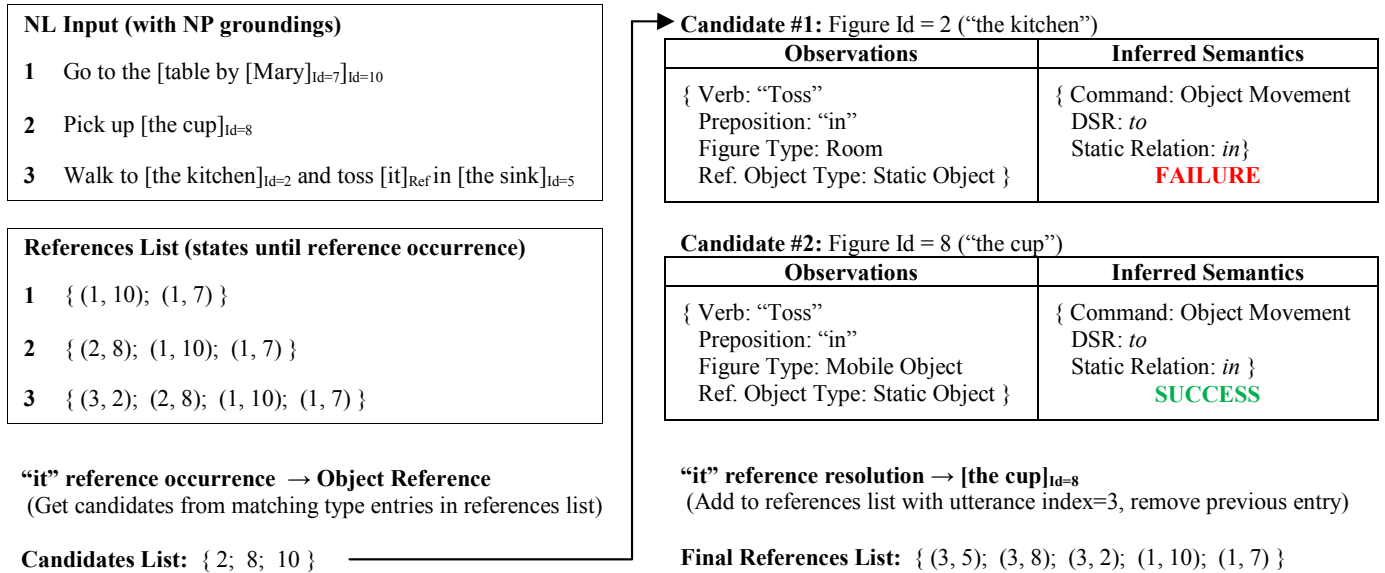


Figure 9.1: Reference resolution example for the instruction “toss it in the sink” expressed by the user during spatial language discourse. Parsed NPs of the natural language input are shown in brackets with their corresponding unique grounding ID numbers as subscripts.

in practice, the integration is seamless: each POS tag assignment is given its own references list that originally is a copy of the most current global references list (most current before utterance processing began). Additionally, in step 3 a temporary list is used for each new sentence grounding and validation check, which is set initially to the most recent references list encapsulating the instructions (of maximum probability) that have already been accepted for the current POS tag assignment’s instruction sequence. Lastly, the resulting references list for each of the candidate instruction sequences is stored until final determination of the maximum probable sequence, whose corresponding references list then becomes the global list.

Figure 9.1 illustrates the reference resolution procedure with an example discourse scenario, displaying the spatial language input and the evolving state of the global references list, among other properties of the algorithm.

9.2 Pragmatics for Physically Embodied Interaction with People

Our spatial language-based HRI methodology, described in Chapter 7, has demonstrated that the A* search algorithm can be used effectively in conjunction with the semantic field model of spatial prepositions to generate robot task solution plans for execution of spatial language instructions provided by the user, including under user-specified natural language constraints. However, the approach was tested only in simulated 2D/3D environments and without modeling direct interaction with people. This section discusses pragmatic considerations in transferring our approach to physical robots for interactions with people in real world environments, and how each was applied in our methodology towards enabling natural human-robot interaction.

Safety is perhaps the most important pragmatic constraint to consider when designing robot systems that are to interact with people. When generating robot task solution plans for given user instructions, it is important that the path/actions taken by the robot be safe for the user, but also for the robot. In our approach described in Chapter 7, the robot task solution plans were generated to achieve optimality in terms of both distance traveled and adherence to user defined constraints, without consideration for the value of generating “safe” solution paths. To address this issue, we incorporated specific pragmatic safety fields into the planning process, one for the robot and another for people within the environment.

Our work presented in Section 7.6 has demonstrated the ease of incorporating pragmatic constraints in our methodology with the use of spatial pragmatic fields. These fields have the same representation as the semantic fields used for computing spatial relations within the environment, and can easily be combined for generating robot task plans that consider both the semantics and the pragmatics of the given instruction. The

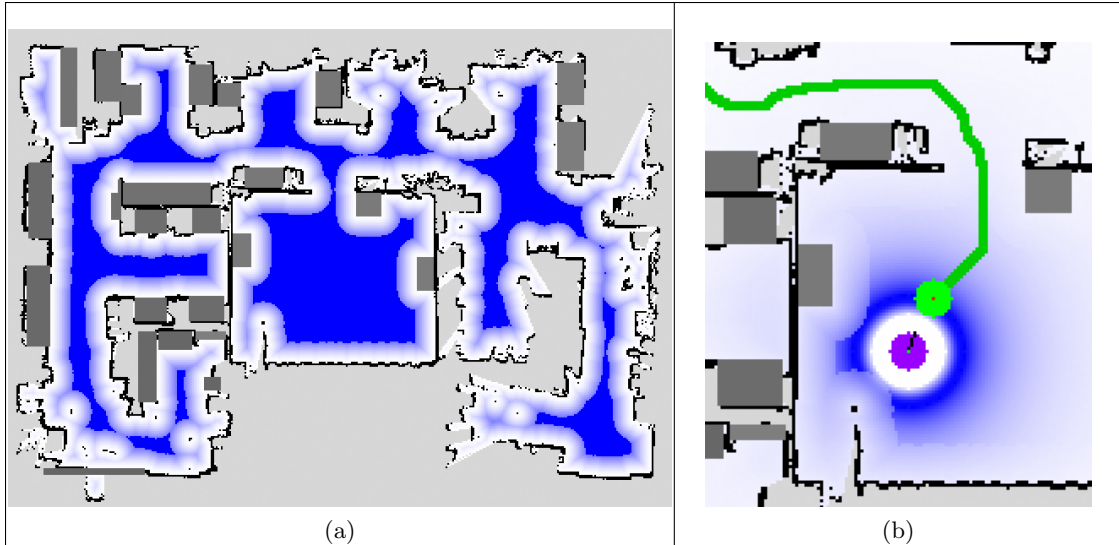


Figure 9.2: (a) SLAM map of laboratory space with pragmatic field for robot safety shown; (b) Example robot approach behavior with combined semantic/pragmatic field shown for at/person safety.

safety field for the robot was generated using a Gaussian function and a safety threshold specifying the minimum desired distance from obstacles and also the function mean (set to $2 \times$ robot radius). Field values for points in the world were designated based on their distance to the nearest obstacle, where distances above the threshold would result in a maximal field value (1.0) and distances below the threshold would be set according to the Gaussian. The safety field for people is generated similarly, with the distance parameter instead referring to the distance from the person. Figure 9.2 (a) shows an example robot safety field computed for a real world laboratory environment using a SLAM map, and Figure 9.2 (b) illustrates use of the person safety field (merged with the at semantic field and pragmatic field for robot safety) in a robot solution for the instruction “Come to me”.

The resulting pragmatic fields have been integrated into the A* cost function of our planning procedure, so as to designate preference for safer solution paths for both the robot and people. Other pragmatic fields can easily be incorporated during planning using our methodology, including for example, those that enforce appropriate approach behaviors (e.g., not from behind) and person-to-person interaction spaces (e.g., do not cross) [Vasquez et al., 2013].

9.3 Generalized Transfer to Robot Systems

Translation from the discretized plan to robot motor commands was accomplished by creating a cost map for the ROS navigation package to use during planning that would strongly favor points along (or close to) the planned solution path. Once created, the cost map is sent to the ROS navigation stack along with the desired goal position in the map. The result is a smooth path that takes into account the motion model of the robot (e.g., omnidirectional vs. differential drive) while following very closely to the path generated by our spatial language-based HRI framework (DWA parameters used: path bias = 30, goal bias = 10). The navigation also takes into account local obstacles encountered during task execution, and is able to quickly re-plan upon encountering an obstruction. Figure 9.3 shows an example of dynamic obstacle avoidance during task execution for obstacles not found in the static map, displaying actual data from a test run with the PR2 robot where a table was introduced into the environment not present in the static map. Figure 9.3 (b) additionally shows the cost map that was generated for task planning (shown in grayscale with values scaling linearly with distance to the original discretized path produced by our spatial planner).

By utilizing the ROS framework to abstract away the generation of robot motor commands from the spatial language task solution, the transfer process is generalized and can easily be replicated for a variety of robot platforms.

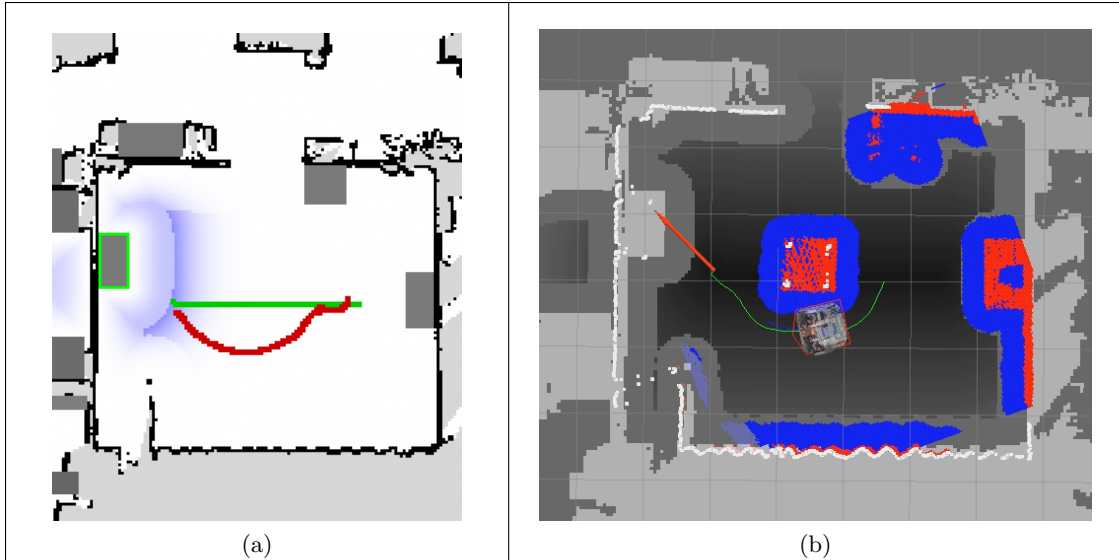


Figure 9.3: Dynamic obstacle avoidance for instruction “Go to the dinner table”. (a) Planned path (green) and actual path (red); (b) Visualization of obstacles detected in robots local map and global plan after robot re-planning.

In order to successfully transfer our spatial language interpretation framework to a physical robot system, there are a few technical challenges that first need to be addressed. The primary of which is a procedure for the translation of the discretized path returned by our planner into appropriate robot motor commands (e.g., wheel velocities) that result in the robot following the desired path. Another major consideration is the autonomous generation of a map of the environment to be used during planning. Last but not least, the robot must be able to localize itself within the generated map of the environment using onboard sensors. Fortunately, previous work conducted by researchers in the field (e.g., [Fox et al., 1997]) have already designed solutions to these challenging research problems, many of which have been packaged within the software framework of the Robot Operating System (ROS) [Willow Garage, 2013] and which are freely available for use. In transferring our approach to a physical robot system (the

PR2 robot platform) we utilized the ROS software packages available for the generation of SLAM maps (gmapping), robot localization (AMCL), and robot navigation planning (global + local planning using DWA).

9.4 Evaluation Results

To evaluate the ability of our robot system to follow natural language directives involving spatial language, we first analyzed the performance of the physical robot platform (PR2 robot) at reaching the desired destination specified in the spatial language instruction. To evaluate the robot’s performance, we conducted a spatial positioning task experiment that consisted of instructing the robot to move to a desired location satisfying a given spatial relation with respect to one or more groundings in the environment, expressed through natural language. The experiment consisted of 14 test instructions given to the robot, two for each static spatial relation analyzed (*near*, *awayfrom*, *between*, *inside*, *outside*, *at*) and two additional test runs for the spatial relation *at*, which is associated with the most common path preposition utilized in spatial instruction tasks (“to”). After each experiment run, the end location of the robot was measured against the goal position generated by our spatial language interpretation framework; specifically the distance between the target end point and actual robot end point was recorded. An example run of the experiment is shown in Figure 9.4 for the instruction “Stand between the printer desk and the whiteboard”, displaying both the planned path and actual path taken by the robot during task execution.

Table 9.3 shows the results of the analysis, which demonstrate the robot’s notable accuracy in estimating its position within the environment, as the distance errors were very small (within 0.2 m). This amount of distance between the final point of the robot and the planned goal point is to be expected, as the ROS navigation module operated with an acceptable goal distance threshold of exactly 0.2 m. The minor differences

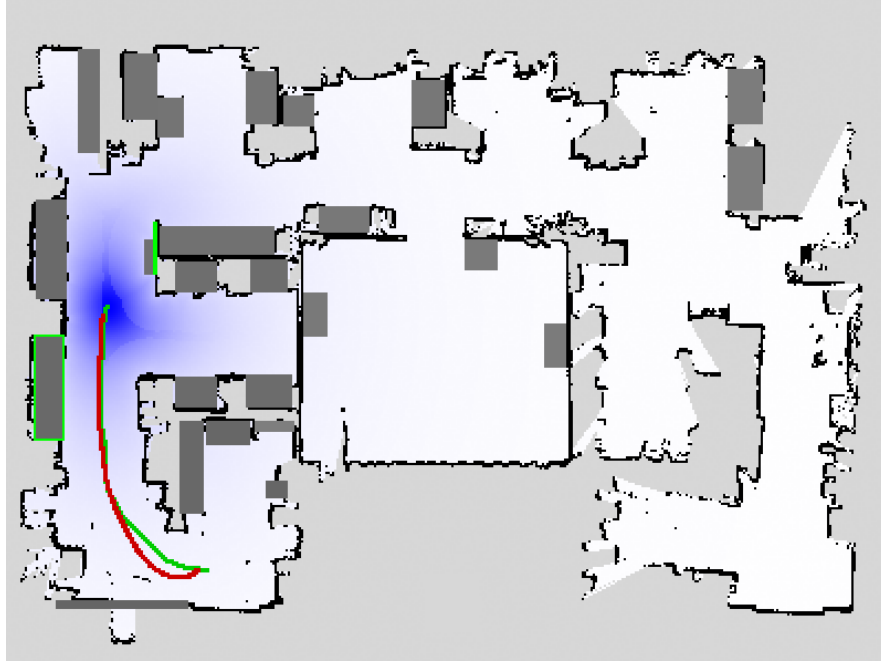


Figure 9.4: Combined semantic/pragmatic field and execution result for the task “Stand between the printer desk and the whiteboard”.

Table 9.3: Accuracy of Final Robot Positions in Spatial Task Experiment

Measure	Mean (Std.)
Distance to Goal Position	0.19 m (0.037 m)
Distance to AMCL Position Estimate	0.07 m (0.038 m)

Note. Distances calculated from the actual robot position measured after task completion.

observed between the robot’s AMCL position estimates and the actual final position highlight the effectiveness of the robot’s onboard localization procedure (implemented in ROS), and demonstrate the ease of which the robot was able to follow the spatial language instructions provided to it by our framework in practice in a real world environment.

Table 9.4: Instructions Given in Test Runs 1-3 with Inference Results for Instruction Sequences

Run	Natural Language Instructions
1	<p>“PR2 can you please head to the dinner table and then pick up the water bottle and take it to my desk so I can have a drink later”</p> <p>Head to the dinner table Pick up the water bottle Take it to my desk</p>
2	<p>“Go ahead and grab the cup if you can and then it would be great if you could go to the kitchen counter and put it on top of it for me”</p> <p>Grab the cup Go to the kitchen counter Put it on top of it</p>
3	<p>“Come into the pen”</p> <p>Come into the pen</p> <p>“Lift up the object close to Juan”</p> <p>Lift up the object close to Juan</p> <p>“Give him it and then step back outside the pen and wait by the entryway”</p> <p>Give him it Step back outside the pen Wait by the entryway</p>

Note. Distinct utterances are listed on separate lines. Instruction sequences inferred by the probabilistic extraction procedure are shown in red.

Next, to demonstrate the capabilities of 1) our probabilistic instruction sequence extraction procedure, 2) our approach to resolving anaphoric expressions, and 3) our integration of pragmatic constraints involving safety fields for interacting with and operating in environments together with humans, we ran three additional test runs of our spatial language discourse interpretation and HRI framework.

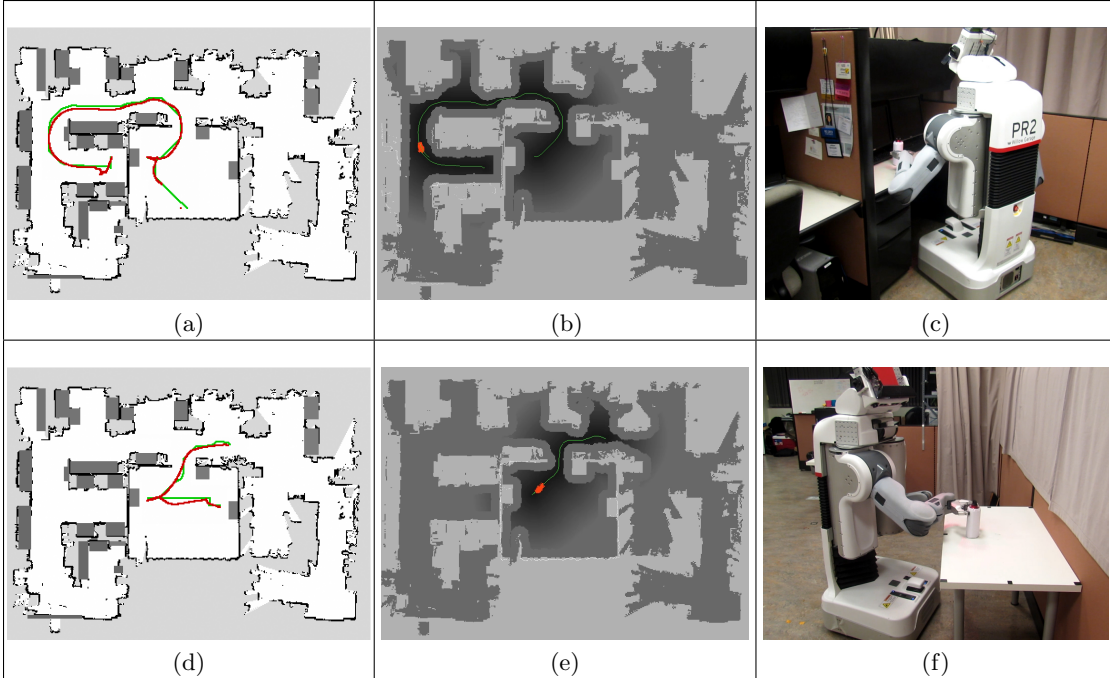


Figure 9.5: (From left to right) Planned (green) and executed paths (red), cost map used for navigation planning with AMCL particles corresponding to robot position estimates, and photograph of PR2 robot just before task termination for test runs 1-2. (a),(b),(c) run 1; (d),(e),(f) run 2.

Each test run involved the user engaging the robot in spatial language discourse, and in particular, providing a series of instructions, with and without the use of anaphora, for the robot to track, resolve references, and execute appropriate task solutions. In total, 11 instructions were evaluated. The spatial language discourse provided to the robot in each of the three test runs are shown in Table 9.4, together with the instruction sequences inferred by our probabilistic instruction sequence extraction procedure. Natural language input given to the robot included those with multiple instructions within a single utterance, which could also contain non-spatial language (e.g., “PR2 can you please head to the dinner table and then pick up the water bottle and take it to my

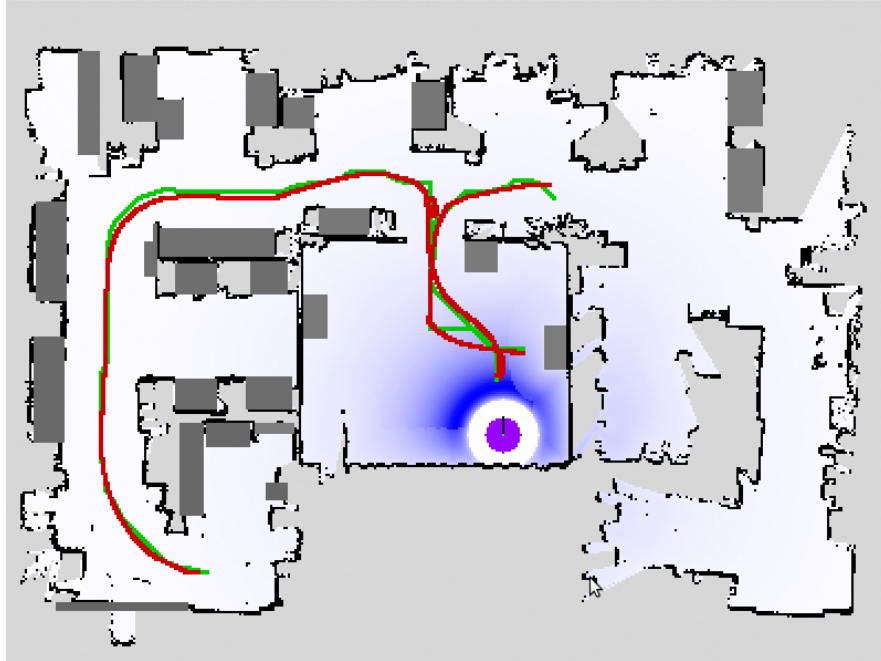
desk so I can have a drink later”), and those with multiple anaphora (“Put it on top of it”) for the robot to attempt to interpret within the context of the discourse.

Figure 9.5 and Figure 9.6 illustrate the performance of the robot during each of the test runs. As evidenced by the results, the robot was able to successfully perform all of the tasks requested by the user in each of the natural language instructions of the test runs. Notable results include the ability of the robot to resolve the multiple anaphoric references in the instruction “Put it on top of it” during the second test run. In this instance, the robot correctly resolved the first reference to the grounding of [the cup], mentioned in the first utterance of the discourse, after disqualifying the initial candidate (most recently grounded NP) of [the kitchen counter] as invalid semantically due to inconsistencies with the parameter requirements of the inferred command of [Object Movement]. Similarly, the instruction “Give him it” expressed by the user during the third test run was correctly resolved by the robot in accordance with the context of the spatial language discourse (“him” → [Juan], “it” → [the object close to Juan]). Figure 9.6 (b) displays the interaction between the robot and the user at the time of object transfer as performed by the robot during execution of this task.

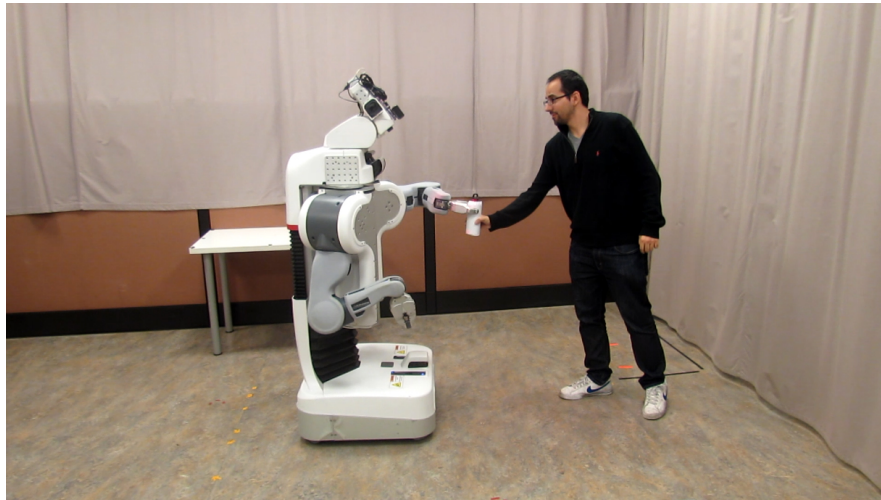
The capability of the robot in successfully interpreting the spatial language discourse expressed during each of the test runs, while also taking into account the pragmatics of the interaction, demonstrates the potential of our approach for use in real world environments with target users.

9.5 Summary

This chapter described the need for enabling autonomous service robots with spatial language understanding and discourse modeling to facilitate natural communication with non-expert users for task instruction and anaphoric reference resolution, and presented a general approach we have developed toward addressing this research challenge. The



(a)



(b)

Figure 9.6: Test run 3 results. (a) Planned path (green) and actual path (red) with semantic/pragmatic fields calculated for hand-off behavior; (b) PR2 robot handing bottle (grounded object referent) to intended person (ground referent for “him”) during task execution.

results obtained from our evaluation testing demonstrate the potential of our methodology for representing dynamic spatial relations, grounding and interpreting the semantics of natural language instructions probabilistically, extracting instruction sequences from unconstrained natural language input, and resolving anaphoric expressions within the context of the current discourse with the user.

Chapter 10

Spatial Language-Based HRI

User Study with Older Adults

This chapter presents a user study we designed and conducted with older adult participants to 1) evaluate the effectiveness and feasibility of our spatial language interpretation framework with end users, and 2) collect data on the types of phrases, responses, and format of natural language instructions given by target users to help inform possible modifications to the spatial language grammar and/or interpretation module of our framework. The study primarily involved manual transcription of users' utterances for accurate online interpretation by the robot; however, the study also incorporated procedures to evaluate the feasibility of fully autonomous human-robot interaction for accomplishing simple household service tasks using a commercial speech recognition engine.

10.1 Study Design

The study consisted of two conditions, Virtual Robot and Physical Robot, both designed to engage the user in human-robot dialogue, and more specifically, to evoke

spatial language instructions from the participant for the robot to interpret and follow according to the context of the discourse and environment. The study design was within-subject, with all participants engaging in both conditions (1 session per condition), and with the second session conducted approximately one week after the first. The order of appearance of the conditions was fixed for each of the participants, with the Virtual Robot condition appearing first as it includes a training session to help the user quickly familiarize themselves with the robot’s capabilities. Each condition (session) lasted 60 minutes, totaling 2 hours of one-on-one interaction, with surveys being administered after both sessions to capture participant perceptions of each study condition independently. The following subsections describe the two conditions in greater detail.

10.1.1 Virtual Robot Condition

In this condition the user interacts with a virtual robot operating within a 2D computer simulated home environment. The overall goal of the scenario is for the robot to execute the tasks expressed to it by the user through natural language. The robot is capable of asking the user clarification queries if it does not understand certain aspects of the given instructions, which typically involve further grounding procedures for the figure and/or reference objects expressed. Thus, the interaction is characterized by human-robot dialogue, with the speech of the robot being generated by the NeoSpeech text-to-speech engine [NeoSpeech, 2009] (the same engine that was used in the SAR exercise coach studies described in Chapters 5 and 6). In this condition, the user is seated in a chair facing a display projecting the simulated home environment, which the user can use to verify the correctness of the robot’s task execution. Commands are issued by the user to the robot using natural language speech, and in all scenarios except one (scenario 2 described below), the spoken instructions are manually transcribed via keyboard input



(a)



(b)

Figure 10.1: (a) Virtual robot condition setup; (b) 2D computer simulated household environment, with example robot task execution path shown for instruction “Pick up the medicine in the guest bathroom”.

by the experimenter in real-time and sent to the robot for interpretation. The interaction setup is shown in Figure 10.1, including the simulated home environment with the robot shown in green, the user in purple, and various movable objects are drawn as smaller colored circles (described below).

At the beginning of the session with the virtual robot, the user is briefed on the four types of tasks/instructions the robot is capable of understanding: 1) *Robot Movement* (e.g., “Go to the kitchen”); 2) *Object Movement* (e.g., “Take the book to my room”); 3) *Object Retrieval* (e.g., “Bring me the bottle from the coffee table”); and 4) *Object Grasp/Release* (e.g., “Pick up/Put down the cup”). The user is encouraged to issue commands using natural speech, in their own words, as if they were commanding a robot in their own home. To help the user communicate effectively with the robot, the participant is given two annotated maps of the simulated home environment which they can refer to at any time during the interaction; one identifies room names that are known/understood by the robot, and the other specifies the names of robot-identifiable appliances/furniture items. The two maps are shown in Figure 10.2. In addition, the user is given a list of objects in the environment that can be moved/transported by the robot, along with their associated color. There are only three types of mobile objects: water bottles (Blue), medicine (Pink), and books (Green).

The Virtual Robot condition examines four different scenarios for data collection and evaluation purposes: 1) *free-form interaction*; 2) *task/instruction training with speech recognition*; 3) *object identification*; and 4) *placement location identification*. The first interaction scenario of the session is a free-form interaction, where the participant is encouraged to give instructions to the robot as they would in their own home, and the choice of task is decided solely by the user (albeit constrained to the four task types discussed above). In the second scenario, the participant undergoes a form of training in the types of instructions the robot is capable of understanding, accomplished by having the user read from a list of written instructions spanning the four task types.

During this scenario, the user wears a headset microphone and provides the instructions directly to the robot via the commercially available Nuance speech recognition engine [Nuance, 2013]. Thus, this scenario serves both as a means for training the user



(a)



(b)

Figure 10.2: Annotated environment maps provided to the user during the session. (a) Room names; (b) appliances/furniture item names.

on what the robot can and cannot understand, and as a method for collecting data on speech recognition accuracy and effectiveness.

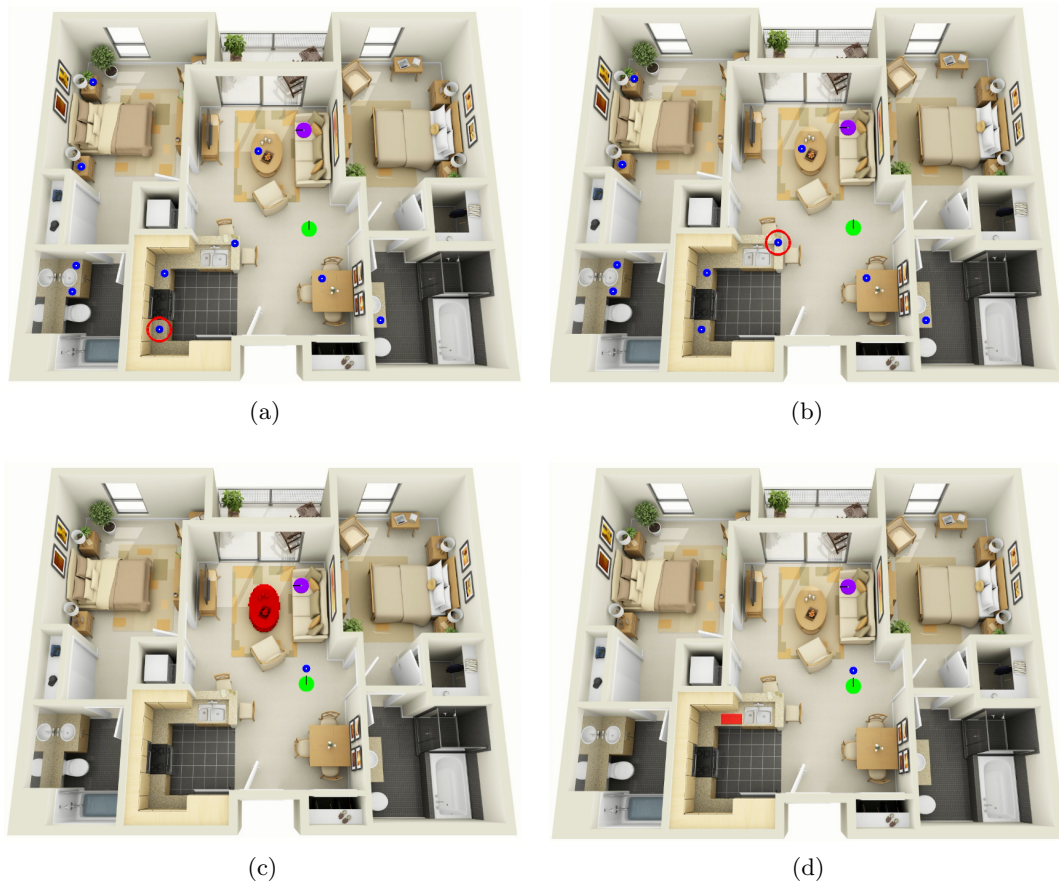


Figure 10.3: Example target objects and placement locations for scenarios 3 and 4 of virtual robot condition. (a) target object left of stove; (b) target object by kitchen sink; (c) target location on coffee table; (d) target location left of kitchen sink.

The third and fourth scenarios are complements of each other and were designed to elicit spatial referencing language from the user, by requiring the user to describe the specific locations of target objects for the robot to pick up, and target locations for the robot to place objects, respectively. In the third scenario, all objects in the household (10 total) are of the same type and color (bottles, medicine, or books). Thus, in this scenario, the user must use spatial language regarding the location of the target object

in order to express the task to the robot. The target object is identified to the user through the simulator by a highlighted red circle surrounding the object. Similarly, in the fourth scenario, the robot is holding the object to be placed, and the target location is expressed to the user through a highlighted red region in the home environment. Both scenarios consist of multiple such target objects/locations expressed in sequence after the user successfully instructs the robot to perform each of the signaled tasks. Figure 10.3 shows screenshots of example task settings for the third and fourth scenario.

10.1.2 Physical Robot Condition

In this condition the user interacts with a physical robot platform situated in the same room as the user. Throughout the session, the user is seated in a chair near the middle of the room. The room is configured with four tables in different locations, each representing a separate area of a typical home environment: the *kitchen*, the *dinner table*, the *bedroom*, and the *coffee table* (in the living room). As with the Virtual Robot condition, the goal of this condition is for the robot to execute the tasks expressed to it by the user through natural language. The primary tasks in this condition require the robot to transport individual household objects to specific locations in the environment as specified by the user. The set of household items used in this condition include: a water bottle, a milk carton, cereal boxes (5 total - all the same brand), medicine (2 total - vitamins and antacid were used as medicine), and one decorative plant. All tables in the environment are appropriately labeled (in bold lettering) so that the user may easily recall the names of the locations represented by each of the tables when giving commands to the robot. Views of the interaction setup for this condition, along with example household items used in the scenario, are shown in Figure 10.4.

The robot platform used in this condition is Bandit, a humanoid torso robot mounted on a MobileRobots Pioneer base (see Section 4.2 for a complete description of the robot).



(a)



(b)



(c)

Figure 10.4: Physical robot condition. (a) View of interaction setup with labeled tables representing typical household areas (from left to right: coffee table, bedroom, dinner table, and kitchen); (b) Bandit, the physical robot platform; (c) Example household items used in the study (from left to right: plant, milk, medicine, bottle, cereal).

Specific adjustments were made to the robot platform to help accomplish the goals of the household service tasks; for example, the robot was modified to include a gripper attachment capable of grasping typical household objects (e.g., bottles, medicine, cereal,

milk, etc.), a Hokuyo laser range finder was added to the base of the robot to aid with navigation and obstacle avoidance, and a PrimeSense Carmine 1.09 RGB-D camera was added to the shoulder of the robot to enable accurate tabletop segmentation and object localization. The physical robot platform used in the study, with all of the modifications described above, is shown in Figure 10.4 (b).

The Physical Robot condition examines three different scenarios for data collection and evaluation purposes: 1) *object identification and placement*; 2) *task-oriented instruction*; and 3) *free-form interaction with speech recognition*. The first interaction scenario can be thought of as a combination of the third and fourth scenarios from the Virtual Robot condition, as the user must command the robot to move a given target object (household item) to a specified target location. Both the object to be moved and the target location are marked by the experimenter (using two green sticky notes) prior to the start of the scenario. Figure 10.5 (a) shows an example setting for the object identification and placement scenario, with a marked target object (bottle) and target location (kitchen).

In the second interaction scenario, the user is given a task that the robot needs to accomplish, and for which they are asked to provide instructions. The task is relayed to the user non-verbally, with the experimenter providing a photograph showing the target state of the environment to be achieved for the task to be completed successfully. Tasks are chosen that require multiple pick-and-place instructions in order to achieve the target goal state from the start state, and typically involve the movement of one or more objects onto a specific table (household area). The exact order of user object placement instructions does not matter, but the user is encouraged to match the relative object positions as closely as possible to the goal state displayed in the image provided. Two example task goal states provided to the user, with specified target object goal locations, are shown in Figures 10.5 (b) and (c).



(a)



(b)



(c)

Figure 10.5: Scenarios 1 and 2 of Physical Robot condition. (a) Object identification and placement scenario, with target object (bottle) and target location (kitchen) marked by green sticky notes; (b), (c) Example task photographs provided to the user displaying task goal states with object target locations (the relative location of objects with respect to one another is important).

The third scenario is a free-form interaction where the user is encouraged to give the robot any command of their choice. In this scenario, the user wears a headset microphone that sends the vocal commands directly to the robot's speech recognition engine. Thus, the interaction between the user and the robot is completely autonomous, without any online transcription of the spoken commands by the experimenter. This scenario was designed to investigate the feasibility of such an interaction with members

of the target user population using commercial speech recognition technology available today [Nuance, 2013].

10.2 Participant Statistics

We recruited elderly individuals to participate in the study through a partnership with be.group, the same senior living organization partner from our previous SAR coaching studies described in Chapters 5 and 6, using flyers and word-of-mouth. We offered a \$20 Target gift card to those willing to participate in the two sessions of the study. In total, 19 older adult participants engaged in both sessions of the study. The sample population consisted of 15 female participants (79%) and 4 male participants (21%). The greater number of female relative to male participants is reflective of the resident statistics of the recruitment facility, and is consistent with our previous studies conducted with older adults from the same partner organization (see Chapters 5 and 6). Participants' ages ranged from 71-97, and the average age was 82 (S.D. = 7.42).

10.3 Measures

10.3.1 Objective Measures

During interaction, there are two possible outcomes after the user gives the robot an instruction: 1) the instruction semantics are inferred without raising flags (semantic errors) and thus the robot plans a task solution that it then executes; or 2) the robot is unable to fully interpret the semantics of the instruction (due to the presence of unknown words/phrases and/or speech recognition errors), and thus a clarification query is posed to the user by the robot in order to resolve any grounding errors and/or confirm the command type (in the case that the command type cannot be inferred with high probability). In this case, the user and robot proceed to engage in a turn-taking

dialogue that terminates if all semantic ambiguities are resolved (at which point the robot executes the instruction), or if the robot deems the progress towards resolving the ambiguities unsatisfactory (e.g., the maximum number of clarification queries is reached, experimentally set to 6 per ground - a threshold set *a priori* by the experimenter). This back-and-forth human-robot dialogue process towards resolving the meaning of a single user instruction is referred to as a *dialogue round*. Many of our objective measures employed the number of dialogue rounds as a normalizing factor.

The objective measures collected (18 total) were chosen to: 1) measure the overall success of the communication between the user and robot; and 2) help characterize the natural language format of spatial tasks and relations expressed by the users to inform possible modifications to our framework.

The performance measures regarding the interaction were: *task success rate* (percentage of tasks that were completed successfully by the robot); *task success rate with repeated attempts* (percentage of tasks that were completed successfully by the robot - after repeated attempts (more than one dialogue round) by the user); *success rate per round* (percentage of rounds that ended in a task execution by the robot); *number of rounds needed to achieve task success*; *number of rounds total during interaction*; and *avg. number of clarification queries per round* (a measure of the fluidity of the interaction and of the comprehension level of the robot).

The remaining data collection measures were: *avg. number of references per round*, *number of total references used*, *maximum number of references used among all participants* (these are all measures of user tendency to use anaphoric references during discourse); *avg. number of yes/no questions posed by the robot*, *avg. number of yes/no responses to yes/no queries* (a measure of user compliance to the robot's questions during clarification procedures). Additionally, in the Physical Robot condition we measured the *total number of instruction sequences of lengths 1-4* (i.e., those containing

one, two, three, and four instructions expressed within a single utterance, respectively), given among all participants. These measures were chosen to evaluate aspects of our approach concerning the interpretation of unconstrained spatial language instructions in user discourse (e.g., instruction sequences, anaphoric references), and dialogue assumptions (e.g., yes/no user responses). Lastly, word count statistics were gathered across both conditions to measure: *verb counts*, *path preposition counts*, and *static preposition counts* for spatial prepositions expressed with noun phrases (e.g., “the table by the kitchen”). All of these data collection measures were gathered to help characterize the format (and meta-format) of natural language instructions and responses expressed by the users, thus helping to inform possible modifications to the grammar and/or semantic interpretation module of our methodology.

10.3.2 Subjective Measures

After each session of the study the participant was asked to fill out three surveys to capture their ratings of the interaction as well as their perceived ease of use of the robot system. The subjective measures included the evaluation of the *enjoyableness of the interaction*, and the perceived *value/usefulness of the interaction*. These measures were the same as those measured in the embodiment comparison user study, described in Chapter 6. In addition, the *intelligence of the robot*, and the *social presence of the robot* were also measured using the same scale as in the embodiment comparison study.

To measure the perceived ease of use of the robot system, primarily targeted towards evaluating the ease of interaction and ability of the user to communicate task goals to the robot, we administered the USE questionnaire [Lund, 2001] as a validated means for measuring usability. This questionnaire records participant responses to various usability related questions posed utilizing a 7-point Likert scale (*Ease of Use*, *Ease of Learning*, and *Satisfaction*).

10.4 Results

10.4.1 Virtual Robot Condition Results

The collected statistics regarding the performance of our spatial language interpretation framework in the Virtual Robot condition were very encouraging. The overall task success rate for the robot averaged 78.6% (S.D.=14.5) among all participants ($n = 19$). This measure refers to the percentage of tasks that were successfully completed by the robot after receiving natural language instructions from the participant in one or more dialogue rounds for each task. Additionally, upon considering only tasks for which the user provided at least one additional dialogue round after initial failure of the first round (i.e., the user employed repeated attempts to achieve task success), the task success rate increased to 82.8% (S.D.=12.8) on average among all of the participants. These results demonstrate the ability of our spatial language framework to correctly interpret and follow natural language instructions provided during user discourse.

The round success rate achieved by our framework was 84.1% (S.D.=11.9). This measure refers to the percentage of dialogue rounds that were successfully interpreted by the robot into a given action sequence (e.g., robot movement, object movement, object retrieval, etc.), and speaks to the ability of our spatial language framework to infer command semantics from natural language input featuring grammatical subcomponents. The high round success rate observed suggests the database of labeled training examples utilized by the semantic interpretation module of our spatial language framework, in addition to the grammar utilized for English directives and accompanying probabilistic extraction procedure, is sufficiently representative of potential inputs as to successfully interpret natural language phrases from target users with high performance. Table 10.1 provides a summary of the collected statistics for all of the objective measures captured during the virtual robot condition of the study.

Table 10.1: Results of Interaction with Participants ($N = 19$) in Virtual Robot Condition

Objective Measure	Mean (Std.)
Task Success Rate	78.6% (14.5)
Task Success Rate (Repeated Attempts)	82.8% (12.8)
Round Success Rate	84.1% (11.9)
Number of Rounds Needed to achieve each Task Success	1.8 (0.45)
Number of Total Dialogue Rounds User Engaged in with Robot	49.2 (13.7)
Number of Clarification Queries Per Round	0.91 (0.37)
Number of Yes/No Queries	9.35 (5.8)
Number of Yes/No Answers	7.4 (4.5)
Number of References Per Round	0.15 (0.23)
Number of References Used	6.9 (9.5)
Maximum Number of References Used in Session	38

The fluidity of the human-robot interaction was also notable, as illustrated by the relatively low number of clarification queries posed by the robot during the dialogue rounds ($M=0.91$, $S.D.=0.37$), and the low number of dialogue rounds needed for the user to achieve success (by instructing the robot) in the tasks presented to them ($M=1.8$, $S.D.=0.45$). The average number of rounds needed by the participants to achieve success in each task (28 total) is shown in Figure 10.6 to illustrate the varying level of difficulty among the tasks.

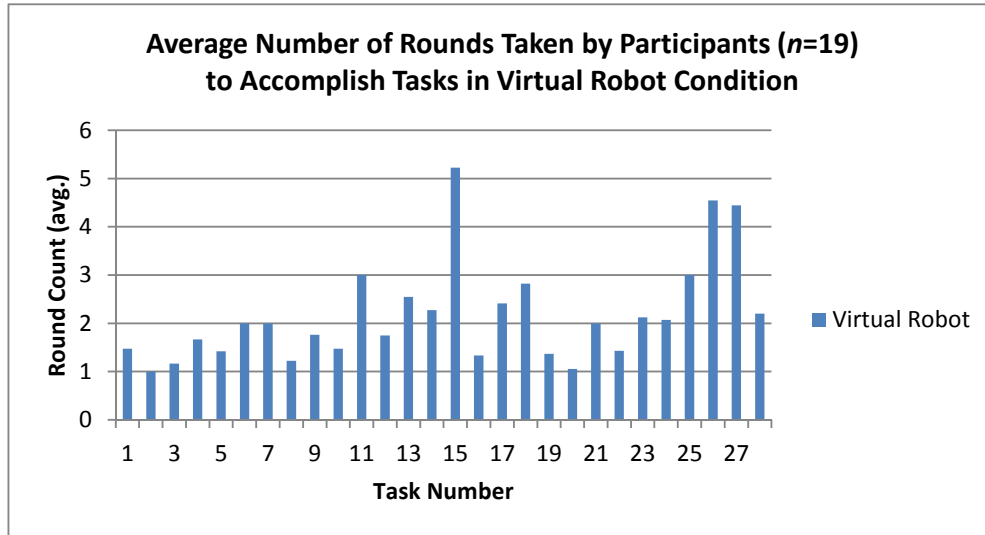


Figure 10.6: Plot of number of dialogue rounds engaged in by the user for each task scenario presented during the Virtual Robot condition as an illustration of the varying level of difficulty among the tasks.

10.4.2 Physical Robot Condition Results

The results of the interaction of the participants with our spatial language framework in the Physical Robot condition were similar to those observed in the Virtual Robot condition, albeit with improved performance overall. Table 10.2 provides a summary of the statistics collected regarding the interaction.

The overall task success rate for the robot averaged 87.4% (S.D.=6.8) among all participants, and increased to 98.0% (S.D.=3.3) among tasks with repeated attempts. The round success rate was observed to be very high at 92.8% (S.D.=6.8), again with a low number of clarification queries posed by the robot during interaction (M=0.47, S.D.=0.22), suggesting an even more fluid interaction than was observed in the Virtual Robot condition. The Physical Robot condition represents the most natural scenario for the user: the user and robot are situated together in the same physical space, and the user can easily identify the spatial locations of objects in the room, both of which

Table 10.2: Results of Interaction with Participants ($N = 19$) in Physical Robot Condition

Objective Measure	Mean (Std.)
Task Success Rate	87.4% (6.8)
Task Success Rate (Repeated Attempts)	98.0% (3.3)
Round Success Rate	92.8% (6.8)
Number of Total Dialogue Rounds User Engaged in	25.8 (7.2)
Number of Clarification Queries Per Round	0.47 (0.22)
Number of Yes/No Queries	3.7 (3.1)
Number of Yes/No Answers	3.2 (2.8)
Number of References Per Round	0.34 (0.26)
Number of References Used	8.5 (6.0)
Maximum Number of References Used in Session	22
Number of Instruction Sequences of Length 1	376 = 72.2%
Number of Instruction Sequences of Length 2	122 = 23.4%
Number of Instruction Sequences of Length 3	22 = 4.2%
Number of Instruction Sequences of Length 4	1 = 0.2%

are helpful for the user to generate natural spatial language commands. It is therefore possible that the increased task success rate, round success rate, and decreased number of clarification queries, were indicative of the increased naturalness of the environment in which the instructions were provided. These results are particularly encouraging as the interaction demonstrated in the Physical Robot condition is most similar to those that would occur in a natural home environment with target users.

Table 10.3: Example Instructions Given by Participants in the Physical Robot Condition with Inference Results for Instruction Sequences

Example	Natural Language Dialogue
1	U: "Go to the kitchen and pick up the milk and take it to the dining room table" Go to the kitchen Pick up the milk Take it to the dining room table - (it referent resolved to milk)
2	U: "I need a bottle between the cereal and the milk" Go to the bottle - (best inference estimate but with low certainty) R: "Would you like me to go to the bottle?" U: "yes" R: "Ok, I'll go to the bottle" U: "Can you pick it up and put it between the cereal and the milk" Pick it up - (it referent resolved to bottle) Put it between the cereal and the milk - (it resolved to bottle)
3	U: "Pick up the medicine" Pick up the (medicine) - (two candidates = 2, 8) R: "Which medicine would you like me to pick up?" U: "the one in front of you" R: "Ok, I'll pick up the medicine" Pick up the [medicine] _{id=8} U: "Take it to the bedroom" Take it to the bedroom - (it referent resolved to [medicine] _{id=8})

Note: Distinct utterances are listed on separate lines. Instruction sequences inferred by the probabilistic extraction procedure are shown in red (with algorithm steps in parentheses).

The number of anaphoric references used per dialogue round during interaction in the Physical Robot condition was also notable, as it represented an increase by a factor of two with respect to the number of references observed in the Virtual Robot condition (M=0.34, S.D.=0.26 vs. M=0.15, S.D.=0.23). This result could be due to the

increased naturalness of the environment, as previously discussed, and also due to the increased complexity of the tasks presented to the user. In the second scenario of the Physical Robot condition, the user was given the task goal state through a photograph, which they then used to help construct instruction sequences to relate to the robot to accomplish the specified goal state. This scenario inherently leads to the user generating multiple instructions that the robot must interpret, and naturally allows for 1) the use of anaphoric references to groundings introduced in prior instructions expressed by the user, and 2) sequencing the instructions within a single utterance. As indicated in Table 10.2, just over one quarter (27.8%) of the instruction sequences provided by the study participants were expressed in utterances containing two or more instructions.

Example human-robot dialogues encountered during interaction in the Physical Robot condition, with extensive use of anaphoric references and multi-instruction utterances, are shown in Table 10.3 for demonstration purposes.

10.4.3 Results of Framework Modification for Prepositional Phrase Attachments

The difficulty of the tasks specified in the Virtual Robot condition, which were designed specifically to elicit spatial referencing language from the user, caused a notable increase in the use of prepositional phrase attachments in the natural language instructions provided by the users. *Prepositional phrase attachments* (PPAs) are prepositional phrases that are attached to the end of sentences to further specify the spatial location of a target region or object. As a demonstration, consider the following instructions provided in a home environment with many candidate water bottles:

- (1) “Pick up the bottle in [NP the kitchen] on the counter by the stove” (2 PPAs)
- (2) “Pick up the bottle in [NP the kitchen] by the stove on the counter” (2 PPAs)

(3) “Pick up the bottle on [NP the counter in the kitchen]” (0 PPAs)

The utterances (1)-(3) all provide directive instructions for the robot to pick up a target object, namely the bottle. Instructions (1) and (2), however, include the use of the same two prepositional phrase attachments to the instruction to further specify the spatial location of the target bottle (“on the counter by the stove”), only expressed in different orders. Instruction (3), on the other hand, describes the spatial location of the bottle without the use of prepositional phrase attachments. The difference between instruction (3) and the others is subtle, but contains an important distinction: the noun phrase expressed (“the counter in the kitchen”) obeys a proper hierarchical ordering of the relative locations in the environment, and thus contains no additional spatial references with the use of prepositional phrase attachments. In short, the NP “the counter in the kitchen” can be correctly resolved to a unique counter in the environment (the one in the kitchen); whereas the NP “the kitchen on the counter by the stove” fails to ground uniquely due to a semantic error (there is no kitchen that is on a counter). Due to the nature of the probabilistic parser of our original spatial language interpretation framework, instructions (1) and (2) were incorrectly discarded as invalid semantically.

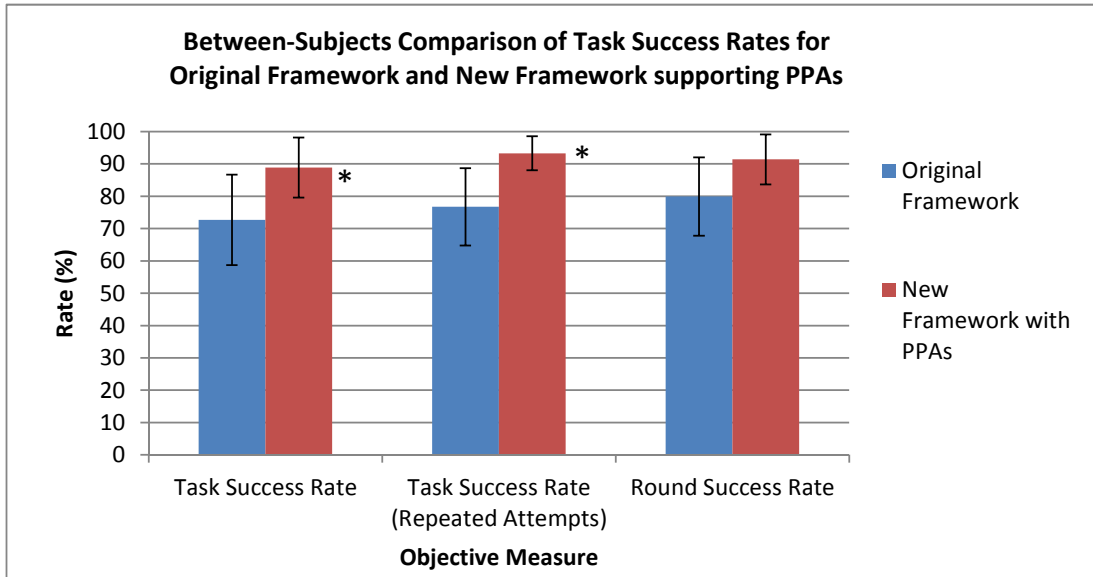
In the context of the Virtual Robot condition, instructions with PPAs caused the robot to pose additional clarification queries to the user to help ground the target object. In some cases, the user would slightly alter the manner in which they provided the instruction to the robot (unbeknown to them) that would allow the framework to correctly interpret the instruction semantics, though in many cases they did not and as a result the robot was unable to interpret the instructions given in the dialogue round.

To address this problem, we extended the spatial language interpretation framework to include the ability to extract PPAs from natural language instructions. The design methodology of the original framework made this extension fairly straightforward, as necessary additions to the probabilistic NP grounding and constraint-based planning

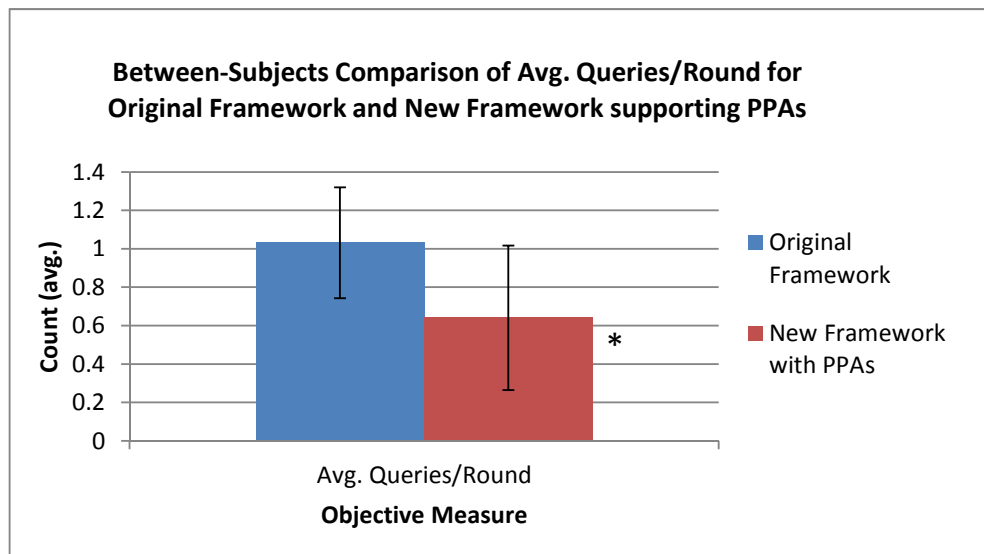
methods were facilitated by the existing spatial semantic field reasoning components of the grounding and planning modules, thereby enabling a seamless extension for the interpretation of PPAs. In particular, the NP grounding module was modified to process individual noun phrases greedily instead of exhaustively: if the module determines that an NP can be grounded uniquely before all child prepositional phrases are processed, the unique grounding is returned and the unprocessed prepositional phrases are returned as PPAs for the NP. These PPAs are then used by the grounding module to help resolve grounding ambiguities within the same expression and/or, if grounding succeeds for the complete instruction, used as semantic field constraints by the planning module when searching for an appropriate task solution (e.g., in (1), “the kitchen” grounds uniquely with two PPAs; *in*(the kitchen) + PPAs are then used to uniquely ground “the bottle”).

Once the extension was completed, the new framework was used for the remainder of the study. In total, 12 participants engaged with the original framework, and the remaining 7 participants engaged with the new framework capable of interpreting PPAs.

To test the effectiveness of the new framework with respect to the original framework, we conducted both a between-subjects analysis and a within-subjects analysis of the performance of both frameworks during interaction with the participants in the Virtual Robot condition. For the between-subjects analysis, the participants were divided into two groups: users who engaged with the original framework ($n = 12$), and users who engaged with the new framework capable of interpreting PPAs ($n = 7$). For the within-subjects analysis, the instruction logs of the 12 participants who engaged with the original framework were relayed to the new framework in an automated procedure performed off-line for data collection purposes. Performance measures were collected to objectively compare the effectiveness of each framework, including the task success rate, round success rate, and the average number of clarification queries posed to the user during the dialogue rounds.

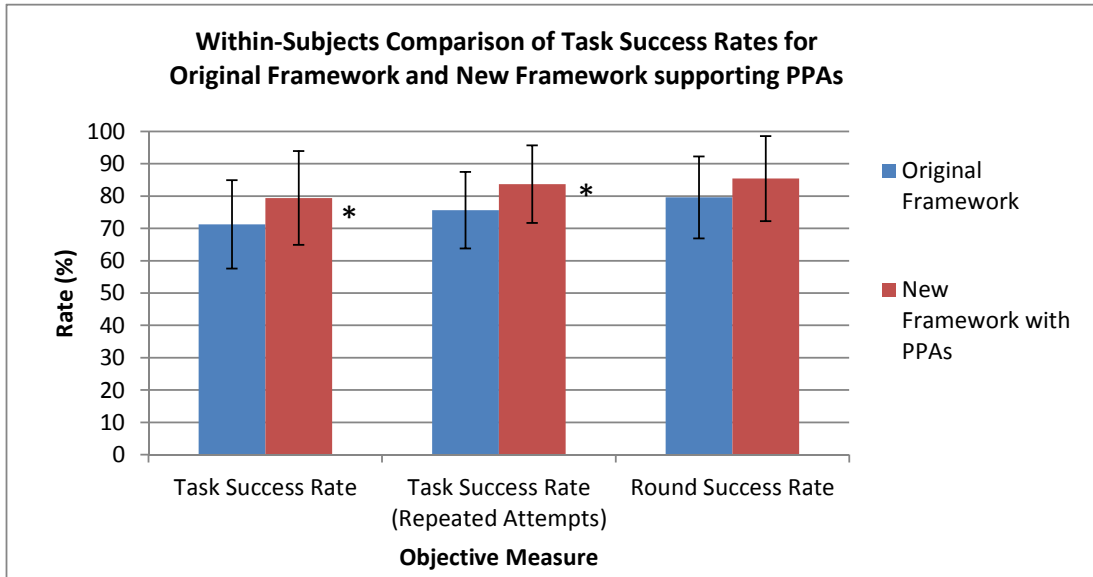


(a)

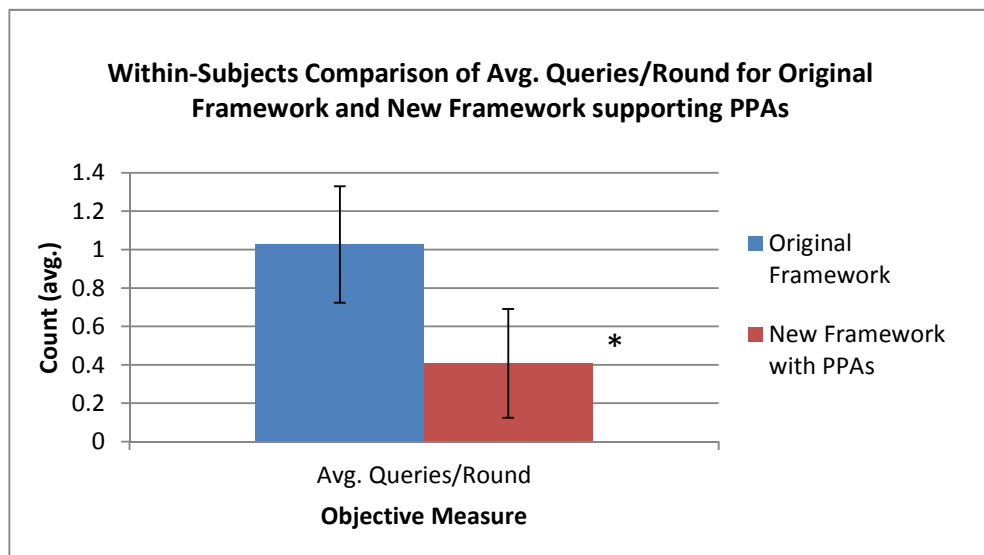


(b)

Figure 10.7: Between-subjects results comparing original framework to new framework capable of interpreting PPAs. (a) Task success rate and round success rate; (b) Average number of clarification queries per dialogue round (Note: Significant results marked by asterisks (*)).



(a)



(b)

Figure 10.8: Within-subjects results comparing original framework to new framework capable of interpreting PPAs. (a) Task success rate and round success rate; (b) Average number of clarification queries per dialogue round (Note: Significant results marked by asterisks (*)).

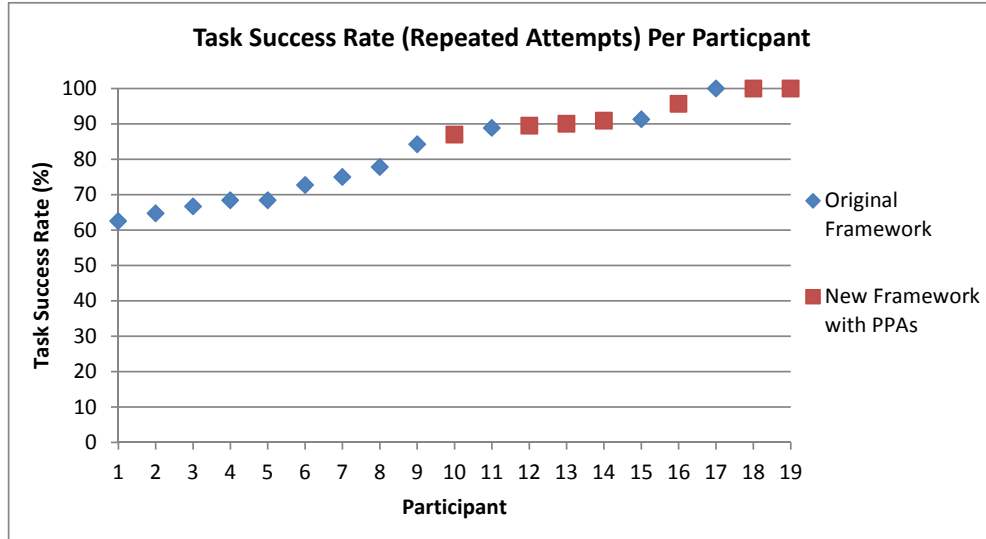


Figure 10.9: Plot of the task success rate (with repeated attempts) achieved in the Virtual Robot condition by each user of the study ($n = 19$), showing both participant groups (original framework vs. new PPA-capable framework) and with the success rates listed in increasing order.

The results of the between-subjects and within-subjects analysis are provided in Figures 10.7 and 10.8, respectively. The results demonstrate the clear performance improvement achieved by the new framework capable of interpreting PPAs, with a task success rate of 89% (S.D.=9.2) for single attempts, and 93% (S.D.=5.2) among repeated attempts, both reaching levels of statistical significance with respect to the original framework ($p < 0.05$). Figure 10.9 further illustrates this result with a plot of the task success rate achieved for each individual user of the study, showing a notable increase in performance among users interacting with the new PPA-capable framework.

10.4.4 Spatial Language Usage Statistics

To analyze the characteristics of the spatial language expressed by the study participants, word count statistics were gathered to measure the number of occurrences of all

Table 10.4: Spatial Language Statistics of Verb and Path Preposition Usage in $N = 1239$ Total Instructions Given by Participants

Verb	Count (%)	Cmd.	% Inf.	Path Prep.	Count (%)	DSR	% Inf.	SSR	% Inf.
pick	414 (33.41%)	AO	100%	up	406 (32.64%)	<i>up</i>	100%	<i>none</i>	100%
put	311 (25.10%)	OM	92%	to	347 (27.89%)	<i>to</i>	100%	<i>at</i>	77%
take	130 (10.49%)	OM	88%	on	220 (17.68%)	<i>to</i>	100%	<i>on</i>	100%
go	126 (10.17%)	RM	100%	(none)	143 (11.50%)	<i>to</i>	75%	<i>at</i>	71%
move	90 (7.26%)	OM	89%	in	50 (4.02%)	<i>to</i>	100%	<i>in</i>	100%
bring	82 (6.62%)	OR	75%	from	16 (1.29%)	<i>to</i>	71%	<i>out</i>	71%
place	29 (2.34%)	OM	92%	between	7 (0.56%)	<i>to</i>	100%	<i>between</i>	100%
get	24 (1.94%)	OR	61%	into	6 (0.48%)	<i>to</i>	100%	<i>in</i>	100%
give	7 (0.56%)	OR	100%	onto	6 (0.48%)	<i>to</i>	100%	<i>on</i>	100%
grab	5 (0.40%)	AO	100%	the left of	6 (0.48%)	<i>to</i>	100%	<i>left-of</i>	100%
turn	4 (0.32%)	RM	100%	front of	5 (0.40%)	<i>to</i>	100%	<i>front-of</i>	100%
remove	3 (0.24%)	OM	67%	near	5 (0.40%)	<i>to</i>	100%	<i>near</i>	100%
wash	3 (0.24%)	RM	67%	next to	5 (0.40%)	<i>to</i>	100%	<i>near</i>	100%
find	2 (0.16%)	RM	100%	the right of	5 (0.40%)	<i>to</i>	100%	<i>right-of</i>	100%
leave	2 (0.16%)	N/A	0%	down	4 (0.32%)	<i>down</i>	100%	<i>none</i>	100%
hold	1 (0.08%)	RM	100%	by	3 (0.24%)	<i>to</i>	100%	<i>near</i>	100%
keep	1 (0.08%)	OM	100%	inside	2 (0.16%)	<i>to</i>	100%	<i>in</i>	100%
milk	1 (0.08%)	OM	100%	on top of	2 (0.16%)	<i>to</i>	100%	<i>on</i>	100%
need	1 (0.08%)	RM	100%	upon	2 (0.16%)	<i>to</i>	100%	<i>on</i>	100%
want	1 (0.08%)	RM	100%	around	1 (0.08%)	<i>around</i>	100%	<i>none</i>	100%
set	1 (0.08%)	N/A	0%	beside	1 (0.08%)	<i>to</i>	100%	<i>in</i>	100%
stand	1 (0.08%)	N/A	0%	close to	1 (0.08%)	<i>to</i>	100%	<i>near</i>	100%
				of	1 (0.08%)	<i>to</i>	100%	<i>behind-of</i>	100%

Note: Word counts (with percentage of total) are shown for all verbs and path prepositions expressed in valid grammatical instructions provided by $n = 19$ participants across both conditions of the study. Initial inference results are shown for the command type, DSR type, and static spatial relation (SSR) as returned by the semantic interpretation module, along with the percentage of inferences (% Inf.) where the indicated inference result was returned for the given verb/preposition. The command types were: RM = Robot Movement, OM = Object Movement, OR = Object Retrieval, and AO = Action on Object (N/A is reported for entries where no inference was made due to grounding errors).

of the different verbs, path prepositions, and static prepositions employed by the participants when issuing instructions to the robot in both study conditions. Table 10.4 shows the counts for each verb and path preposition used in all of the $N = 1239$ valid grammatical instructions issued by the participants that were interpreted by the robot during

interaction. Table 10.4 also provides the most common inference results co-occurring with each verb and path preposition recorded. The inference variables shown were the outputs of the semantic interpretation module, namely: the command type (shown with verbs), and the dynamic spatial relation (DSR) and static spatial relation (SSR) types (shown with path prepositions). The percentage of inference result co-occurrence with the specified verb/preposition is also displayed for each inference variable. As an example, the verb “put” was utilized 311 times during the study by participants when instructing the robot; in 92% of those occurrences, the inferred command type was Object Movement (see Table 10.4). In total, there were four domain-dependent command types available for inference: Robot Movement, Object Movement, Object Retrieval, and Action on Object.

As illustrated by the results, the participants utilized a relatively small set of verbs and path prepositions when issuing instructions to the robot. More specifically, 93% of the 1239 instructions issued employed one of the top six verbs, and 94% of the path prepositions used were among the top five path prepositions encountered (when including (none) as a path preposition option in the case where no path preposition was used in the instruction - e.g., “Bring me the book”). This is an interesting result, as the observed user tendency to reuse the same verbs/prepositions when instructing similar tasks facilitates the probabilistic inference of instruction semantics using relatively small datasets (labeled training examples), especially when employed with our spatial language methodology which separates the inference of command semantics from the grounding of noun phrases for tractability. For reference, the semantic database used during the study consisted of only 372 training examples (labeled with target command, DSR, and SSR types), while the resulting task performance of the robot was quite high (see Tables 10.1 and 10.2). The semantic interpretation module of our approach, by virtue of the Naïve Bayes inference method, is easily capable of performing effective

inference on larger datasets (e.g., with thousands of examples), yet based on the participant language usage statistics and encouraging performance results obtained from the user study, an increase in the number of training examples does not appear to be necessary to achieve high performance.

It must be noted, however, that the inference results shown in Table 10.4 represent only the most common *initial* inference of (command, DSR, and SSR) types, and do not necessarily indicate the final results used by the planner to generate robot task solutions for each verb/path preposition listed. This is because each inference result carries with it a corresponding probability of correctness (or confidence weight), which is used by the dialogue module when deciding whether or not to pose clarification queries to the user (i.e., low confidence values trigger clarification questions), which may alter the final designation of each inference variable. Low confidence inferences are typically caused when the participant utilizes an unknown verb or verb/path preposition combination. However, based on the relatively low number of clarification queries posed during interaction (see Tables 10.1 and 10.2), the high frequency of only a small set of verbs and path prepositions, and the fact that clarification queries often only targeted noun phrase grounding ambiguities, this scenario was rather infrequent.

Table 10.5 shows the word count statistics for static prepositions that were utilized by participants within noun phrases to express static spatial relations (e.g., “Pick up [NP the cup near the TV]”; “Put the bottle on [NP the nightstand to the left of the bed]”). The results were similar to those observed regarding verb and path preposition usage: the participants utilized a fairly small set of prepositions when relaying static spatial relations. However, in this case the results are slightly misleading as one of the most frequently used prepositions, “of”, was often combined with a spatial noun phrase to express the complete spatial relation (e.g., “Put the cup on [NP [NP the left side] of the counter]”; “Put the book down at [NP [NP the front edge] of the coffee table]”).

Table 10.5: Spatial Language Statistics of Static Preposition Usage within Noun Phrases of $N = 1239$ Total Instructions Given by Participants

Static Prep.	Count (%)	Static Prep.	Count (%)
on	240 (29.20%)	between	4 (0.49%)
from	174 (21.17%)	left of	3 (0.36%)
in	156 (18.98%)	on top of	3 (0.36%)
of	143 (17.40%)	the left of	3 (0.36%)
by	18 (2.19%)	off	2 (0.24%)
near	18 (2.19%)	right of	2 (0.24%)
to the right of	13 (1.58%)	behind	1 (0.12%)
next to	11 (1.34%)	beside	1 (0.12%)
to ¹	9 (1.09%)	front of	1 (0.12%)
at	8 (0.97%)	off of	1 (0.12%)
close to	5 (0.61%)	over	1 (0.12%)
to the left of	5 (0.61%)		

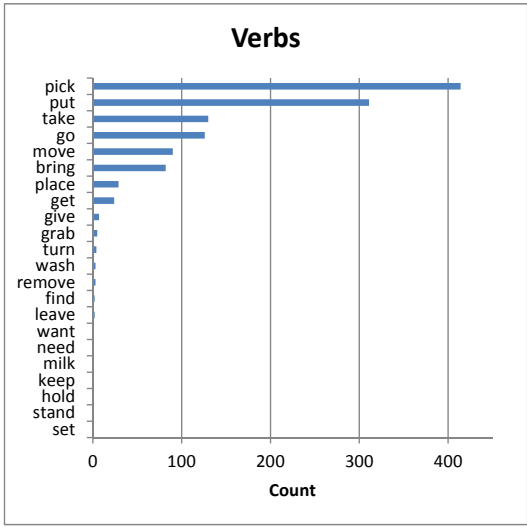
Note: Word counts (with percentage of total) are shown for all static prepositions expressed in valid grammatical instructions provided by $n = 19$ participants across both study conditions.

¹“to” is considered a static spatial preposition in our framework only when paired with a semantic field specifier noun phrase (e.g., “Pick up [_{NP} the bottle to [_{NP} the left]]”).

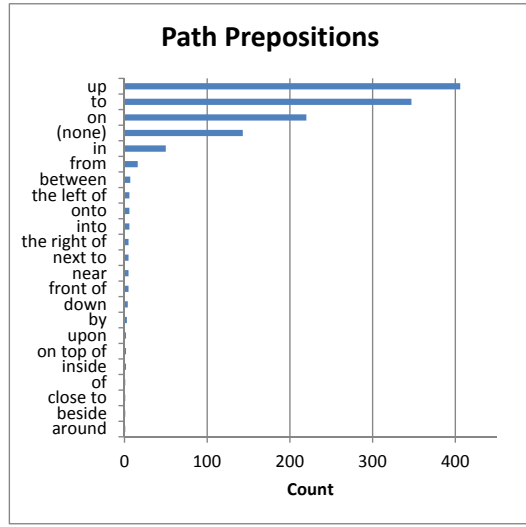
The complete word count statistics for verbs, path prepositions, and static prepositions utilized by the study participants are summarized as histograms in Figure 10.10 for illustration purposes.

10.4.5 Subjective Evaluation Results

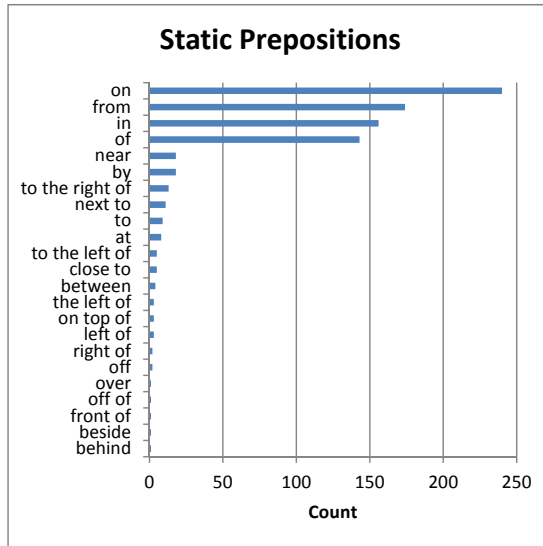
The participant evaluations of the interaction and of the robot, obtained from the surveys administered after the Virtual Robot and Physical Robot conditions, respectively, demonstrated a high rating of our service robotics approach among all of the subjective evaluation items measured. Specifically, the enjoyableness of the interaction (M=8.1, S.D.=1.7), the value/usefulness of the interaction (M=7.9, S.D.=2.3), the intelligence of the robot (M=8.3, S.D.=2.0), and the social presence of the robot (M=7.6, S.D.=1.5), all received high ratings from the participant evaluations, which is very encouraging.



(a)

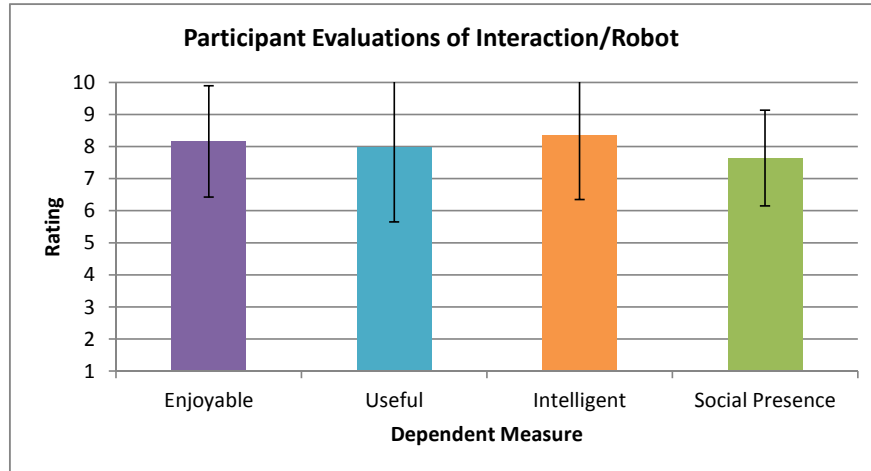


(b)

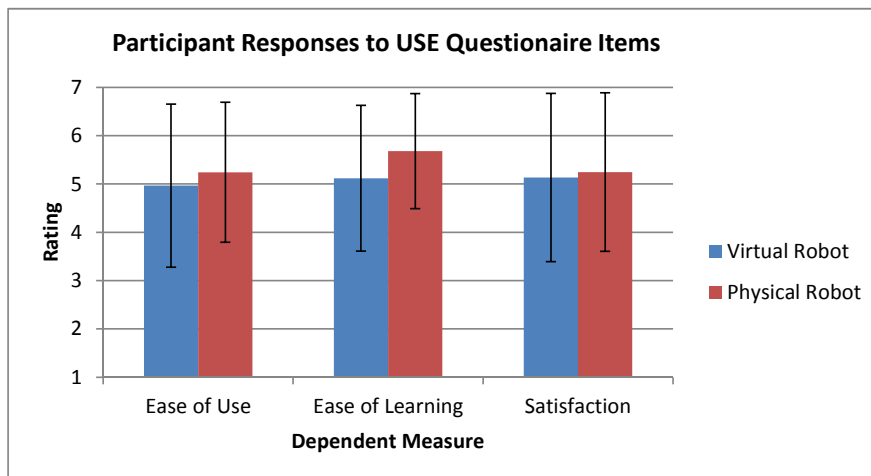


(c)

Figure 10.10: Word count histograms for spatial language encountered in all $N = 1239$ instructions provided by participants during both study conditions. (a) Verb counts; (b) Path preposition counts; (c) Static preposition counts for those expressed within noun phrases as spatial relations.



(a)



(b)

Figure 10.11: Participants evaluation results. (a) Evaluation of interaction and of service robot; (b) Evaluation of interaction with respect to USE questionnaire items.

The results obtained from the participant evaluations of the system with respect to the USE questionnaire items on usability, are also very encouraging as they showed the participants rated the household service robot presented in the study highly in terms of ease-of-use, ease of learning, and satisfaction. Figure 10.11 displays a summary of

the subjective measures captured for the participant evaluation of the interaction and household service robot.

10.5 Summary

This chapter presented a multi-session user study conducted with older adults to evaluate the effectiveness of our spatial language interpretation HRI framework across a variety of objective performance and participant evaluation measures. The results of the study validate our service robotics-based approach and its effectiveness in interpreting and following natural language instructions from target users: the participants commanded the robot to high task success rates in both conditions of the study; our approach incorporating the interpretation of prepositional phrase attachments demonstrated superior success rates; and the participants rated the household service robot highly in terms of enjoyableness and usefulness of the interaction, intelligence, social presence, and usability.

Chapter 11

Summary

The growing population of aging adults is increasing the demand for healthcare services worldwide. Research has shown that regular physical exercise, social interaction, and companionship are critical for maintaining and improving the overall health of elderly individuals [Baum et al., 2003; Dawe and Moore-Orr, 1995; Moak and Agrawal, 2010; Paykel, 1994; Stansfeld et al., 1997], though with the shortfall of nurses and caregivers already becoming an issue [American Association of Colleges of Nursing, 2010; American Health Care Association, 2008; Buerhaus, 2008], access to available care is diminishing. *Socially assistive robotics* (SAR) and *service robotics* have the potential to help address this need.

This dissertation has addressed specific challenges in socially assistive and service robotics-based robot-guided and user-guided interactions. In robot-guided interaction, it presented the approach, design methodology, and implementation details of a novel SAR approach developed to motivate and engage elderly users in simple physical exercise. In the area of user-guided interaction, it presented a novel approach for enabling autonomous service robots to follow natural language commands from non-expert users,

including under user specified constraints, with a particular focus on spatial language understanding.

11.1 SAR Coach for Motivating Therapeutic Behavior

This dissertation presented the design methodology, implementation and evaluation details of a novel SAR approach to motivate and engage elderly users in simple physical exercise. The approach incorporates insights from psychology research into intrinsic motivation and contributes five clear design principles for SAR-based therapeutic interventions. To evaluate the approach and its effectiveness in gaining user acceptance and motivating physical exercise, it was implemented as an integrated system and three user studies were conducted with older adults, to investigate: 1) the effect of praise and relational discourse in the system towards increasing user motivation; 2) the role of user autonomy and choice within the interaction; and 3) the effect of embodiment in the system by comparing user evaluations of similar physically and virtually embodied SAR exercise coaches in addition to evaluating the overall SAR system.

The successful acceptance of the presented SAR approach by elderly users, as evidenced by the high participant evaluations of the system and consistent task performance in all of our user studies, validates the approach, design, algorithms, and effectiveness of our SAR methodology, and illustrates the potential of such technology to help older adults achieve beneficial health outcomes and improve quality of life.

11.2 Spatial Language-Based HRI Framework

This dissertation presented a novel methodology that allows service robots to interpret and follow spatial language instructions, with and without user-specified natural language constraints and/or unvoiced pragmatic constraints. The methodology is generalizable and can be applied across many human-robot interaction domains for a variety of assistive robot behaviors, including user-robot task instruction, teaching, modification, and guidance. In particular, this work contributes a general computational framework for the representation of *dynamic spatial relations* (DSRs), including a novel extension to the *semantic field* model of spatial prepositions, which enables the representation of path prepositions containing both local and global properties. The methodology also contributes a probabilistic approach in the inference of instruction semantics, and in the associated grounding of noun phrases utilizing the proposed computational field representation of spatial relations. The approach allows for robot motion planning and execution of multi-step instruction sequences in real-world continuous domains while providing robustness to sensor noise and environmental uncertainty.

To evaluate the approach and its effectiveness in gaining user acceptance, it was implemented as an integrated system and a user study was conducted with older adults, to 1) evaluate the effectiveness and feasibility of our spatial language interpretation framework with end users, and 2) collect data on the types of phrases, responses, and format of natural language instructions given by target users to help inform possible modifications to the spatial language grammar and/or interpretation module of our framework.

The successful acceptance of the presented service robotics approach by elderly users, as evidenced by the high participant evaluations of the system and high rate of task success in our user study, validates the approach, design, algorithms, and effectiveness

of our service robotics methodology, and illustrates the potential of such technology to help older adults live independently.

11.3 Limitations and Future Work

The high level of task success observed in our spatial language framework user study conducted with older adults is encouraging, however the methodology presented in this dissertation possesses limitations that need to be addressed in future work. For instance, the primary reason for task failure observed in the user study was due to incorrect grounding by the NP grounding module of the framework, attributed to incorrect inference from ambiguous object/task descriptions by the user. To facilitate communication with the user and enhance the fluidity of the interaction, the spatial language framework often employs context-based grounding during interaction that attempts to infer the correct ground referenced by the user in the case of ambiguity. This inference-based grounding procedure works very well when instructions are provided in the current context of the robot’s visual space (e.g. “pick up the medicine” - and there is only one medicine in front of the robot), however when the instructions refer to a future local context (e.g., “go to the bathroom and pick up the medicine”), it has been shown to generate incorrect groundings due to the framework’s current inability to create future local contexts during task planning. Enhancements to the inference grounding procedure represent the most critical and most interesting avenue for future research.

Bibliography

- American Association of Colleges of Nursing. Nursing shortage fact sheet, 2010. Fact sheet.
- American Health Care Association. Summary of 2007 AHCA survey nursing staff vacancy and turnover in nursing facilities, 2008. Report.
- W. A. Bainbridge, J. W. Hart, E. S. Kim, and B. Scassellati. The benefits of interactions with physically present robots over video-displayed agents. *International Journal of Social Robotics*, 3(1):41–52, 2011. doi: 10.1007/s12369-010-0082-7.
- V. Balakrishnan and P. Yeow. Texting satisfaction: Does age and gender make a difference? *International Journal of Computer Science and Security*, 1(1):85–96, 2007.
- C. Bartneck. Interacting with an embodied emotional character. In *Proceedings of the 2003 International Conference on Designing Pleasurable Products and Interfaces*, pages 55–60, New York, 2003. ACM. doi: 10.1145/782896.782911.
- E. E. Baum, D. Jarjoura, A. E. Polen, D. Faur, and G. Rutecki. Effectiveness of a group exercise program in a long-term care facility: A randomized pilot trial. *Journal of the American Medical Directors Association*, 4(2):74–80, 2003. doi: 10.1016/S1525-8610(04)70279-0.
- J. M. Beer, C.-A. 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 *Proceedings of the Seventh Annual ACM/IEEE International Conference on Human-Robot Interaction, HRI '12*, pages 335–342, New York, NY, USA, 2012. ACM. ISBN 978-1-4503-1063-5. doi: 10.1145/2157689.2157806. URL <http://doi.acm.org/10.1145/2157689.2157806>.
- T. W. Bickmore and R. W. Picard. Establishing and maintaining long-term human-computer relationships. *ACM Transactions on Computer-Human Interaction*, 12(2): 293–327, June 2005. doi: 10.1145/1067860.1067867.

- T. W. Bickmore, D. Schulman, and L. Yin. Maintaining engagement in long-term interventions with relational agents. *International Journal of Applied Artificial Intelligence*, 24(6):648–666, 2010.
- J. Bohnemeyer. The unique vector constraint: The impact of direction changes on the linguistic segmentation of motion events. In E. van der Zee and J. Slack, editors, *Representing Direction in Language and Space*, pages 86–110. Oxford University Press, Oxford, 2003.
- P. Buerhaus. Current and future state of the US nursing workforce. *JAMA*, 300(20):2422–2424, 2008. doi: 10.1001/jama.2008.729.
- C. Burgar, P. Lum, P. Shor, and H. Van der Loos. Development of robots for rehabilitation therapy: The palo alto va/standford experience. *Journal of Rehabilitation Research and Development*, 37(6):663–673, 2002.
- R. Cantrell, P. Schermerhorn, and M. Scheutz. Learning actions from human-robot dialogues. In *Proc. IEEE RO-MAN*, pages 125–130. IEEE, 2011.
- J. G. Carbonell and R. D. Brown. Anaphora resolution: A multi-strategy approach. In *Proceedings of the 12th Conference on Computational Linguistics - Volume 1, COLING '88*, pages 96–101, Stroudsburg, PA, USA, 1988. Association for Computational Linguistics. ISBN 963 8431 56 3. doi: 10.3115/991635.991656. URL <http://dx.doi.org/10.3115/991635.991656>.
- L. A. Carlson and P. L. Hill. Formulating spatial descriptions across various dialogue contexts. In K. R. Coventry, T. Tenbrink, and J. Bateman, editors, *Spatial Language and Dialogue*, pages 89–103. Oxford University Press, New York, 2009.
- Centers for Disease Control and Prevention. Morbidity and mortality weekly report, 2003. Department of Health and Human Services.
- C. Chao, M. Cakmak, and A. L. Thomaz. Towards grounding concepts for transfer in goal learning from demonstration. In *Proceedings of the Joint IEEE International Conference on Development and Learning and on Epigenetic Robotics (ICDL-EpiRob)*, volume 2, pages 1–6. IEEE, 2011.
- Y. S. Choi, T. Chen, A. Jain, C. Anderson, J. Glass, and C. Kemp. Hand it over or set it down: A user study of object delivery with an assistive mobile manipulator. In *The 18th IEEE International Symposium on Robot and Human Interactive Communication*, pages 736–743, 2009. doi: 10.1109/ROMAN.2009.5326254.
- S. J. Colcombe and A. F. Kramer. Fitness effects on the cognitive function of older adults: A meta-analytic study. *Psychological Science*, 14(2):125–130, 2003. doi: 10.1111/1467-9280.t01-1-01430.

- S. J. Colcombe, A. F. Kramer, K. I. Erickson, P. Scalf, E. McAuley, N. J. Cohen, A. Webb, G. J. Jerome, D. X. Marquez, and S. Elavsky. Cardiovascular fitness, cortical plasticity, and aging. *Proceedings of the National Academy of Sciences of the United States of America*, 101(9):3316–3321, 2004. doi: 10.1073/pnas.0400266101.
- M. Csikszentmihalyi. *Beyond boredom and anxiety*. Jossey-Bass, San Francisco, 1975.
- M. Csikszentmihalyi. *The evolving self: A psychology for the third millennium*. Harper-Collins, New York, 1993.
- D. Dawe and R. Moore-Orr. Low-intensity, range-of-motion exercise: Invaluable nursing care for elderly patients. *Journal of Advanced Nursing*, 21(4):675–681, 1995. doi: 10.1046/j.1365-2648.1995.21040675.x.
- E. Deci and R. Ryan. *Intrinsic motivation and self-determination in human behavior*. Plenum Press, New York, 1985.
- R. A. Dienstbier and G. K. Leak. Effects of monetary reward on maintenance of weight loss: An extension of the overjustification effect. Paper presented at the American Psychological Association Convention, Washington, D.C., 1976.
- S. Dubowsky, F. Genot, S. Godding, H. Kozono, A. Skwersky, H. Yu, and L. S. Yu. PAMM - a robotic aid to the elderly for mobility assistance and monitoring: A "helping-hand" for the elderly. In *Proceedings of the IEEE International Conference on Robotics and Automation*, volume 1, pages 570–576, 2000. doi: 10.1109/ROBOT.2000.844114.
- J. Fasola and M. J. Matarić. Robot motivator: Increasing user enjoyment and performance on a physical/cognitive task. In *Proceedings of the IEEE International Conference on Development and Learning*, pages 274–279, 2010. doi: 10.1109/DEVLRN.2010.5578830.
- D. J. Feil-Seifer and M. J. Matarić. Defining socially assistive robotics. In *9th International Conference on Rehabilitation Robotics (ICORR)*, pages 465–468, June 2005. doi: 10.1109/ICORR.2005.1501143.
- D. J. Feil-Seifer and M. J. Matarić. Using proxemics to evaluate human-robot interaction. In *Proceedings of the International Conference on Human-Robot Interaction*, pages 143–144, 2010. doi: 10.1109/HRI.2010.5453225.
- D. J. Feil-Seifer and M. J. Matarić. Distance-based computational models for facilitating robot interaction with children. *JHRI*, pages 55–77, Aug 2012. doi: 10.5898/jhri.1.1.feil-seifer. URL <http://dx.doi.org/10.5898/jhri.1.1.feil-seifer>.

- C. D. Fisher. The effects of personal control, competence, and extrinsic reward systems on intrinsic motivation. *Organizational Behavior and Human Performance*, 21(3): 273–288, 1978. doi: 10.1016/0030-5073(78)90054-5.
- D. Fox, W. Burgard, and S. Thrun. The dynamic window approach to collision avoidance. *IEEE Robotics and Automation*, 4(1):23–33, 1997.
- B. French, D. Tyamagundlu, D. Siewiorek, A. Smailagic, and D. Ding. Towards a virtual coach for manual wheelchair users. In *Proceedings of International IEEE Symposium of Wearable Computers*, pages 77–80, 2008. doi: 10.1109/ISWC.2008.4911589.
- H. Fujiyoshi and A. Lipton. Real-time human motion analysis by image skeletonization. In *Proceedings of the Workshop on Applications of Computer Vision*, pages 15–21, October 1998. doi: 10.1109/ACV.1998.732852.
- K.-P. Gapp. Basic meanings of spatial relations: Computation and evaluation in 3d space. In *Proceedings of the 12th National Conference on Artificial Intelligence (AAAI-94)*, volume 2, pages 1393–1398. AAAI Press, 1994.
- L. K. George, D. G. Blazer, D. C. Hughes, and N. Fowler. Social support and the outcome of major depression. *The British Journal of Psychiatry*, 154(4):478–485, 1989. doi: 10.1192/bjp.154.4.478.
- W. Harwin, A. Ginige, and R. Jackson. A robot workstation for use in education of the physically handicapped. *Biomedical Engineering, IEEE Transactions on*, 35(2): 127–131, 1988. ISSN 0018-9294. doi: 10.1109/10.1350.
- N. Hawes, M. Klenk, K. Lockwood, G. S. Horn, and J. D. Kelleher. Towards a cognitive system that can recognize spatial regions based on context. In *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence*, pages 200–206, Palo Alto, CA, 2012. AAAI Press.
- M. Heerink, B. Kröse, V. Evers, and B. Wielinga. Assessing acceptance of assistive social agent technology by older adults: The Almere Model. *International Journal of Social Robotics*, 2(4):361–375, 2010. doi: 10.1007/s12369-010-0068-5.
- D. Hewlett, W. Kerr, T. J. Walsh, and P. Cohen. A framework for recognizing and executing verb phrases. In *2011 Robotics: Science and Systems Workshop: HRI Workshop on Grounding Human-Robot Dialog for Spatial Tasks*, Los Angeles, CA, 2011.
- A. Howard and N. Roy. The robotics data set repository (Radish), 2003. URL <http://radish.sourceforge.net>.

- O. C. Jenkins, C. Chu, and M. J. Matarić. Nonlinear spherical shells for approximate principal curves skeletonization. Technical Report CRES-04-004, University of Southern California Center for Robotics and Embedded Systems, 2004.
- Y. Jung and K. M. Lee. Effects of physical embodiment on social presence of social robots. In *Proceedings of Presence, 2004*, pages 80–87, 2004.
- D. Jurafsky and J. H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Speech Recognition, and Computational Linguistics*. Prentice Hall, 2008.
- L. Kahn, M. Verbuch, Z. Rymer, and D. Reinkensmeyer. Comparison of robot-assisted reaching to free reaching in promoting recovery from chronic stroke. In *International Conference on Rehabilitation Robotics*, pages 39–44, April 2001.
- N. E. Kang and W. C. Yoon. Age- and experience-related user behavior differences in the use of complicated electronic devices. *International Journal of Human-Computer Studies*, 66(6):425–437, 2008. doi: 10.1016/j.ijhcs.2007.12.003.
- J. D. Kelleher and F. J. Costello. Applying computational models of spatial prepositions to visually situated dialog. *Computational Linguistics*, 35(2):271–306, 2009.
- C. Kidd and C. Breazeal. Robots at home: Understanding long-term human-robot interaction. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3230–3235, 2008. doi: 10.1109/IROS.2008.4651113.
- C. Kidd, W. Taggart, and S. Turkle. A sociable robot to encourage social interaction among the elderly. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 3972–3976, 2006. doi: 10.1109/ROBOT.2006.1642311.
- D. Klein and C. D. Manning. Accurate unlexicalized parsing. In *Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics*, pages 423–430, Sapporo, Japan, July 2003. Association for Computational Linguistics. doi: 10.3115/1075096.1075150.
- N. Koenig, L. Takayama, and M. J. Matarić. Learning from demonstration: Communication and policy generation. In *12th International Symposium on Experimental Robotics*, Dec 2010.
- T. Kollar, S. Tellex, D. Roy, and N. Roy. Toward understanding natural language directions. In *Proc. ACM/IEEE Int’l Conf. on Human-Robot Interaction (HRI)*, pages 259–266. IEEE, 2010.
- H. Kress-Gazit, G. E. Fainekos, and G. J. Pappas. Translating structured English to robot controllers. *Advanced Robotics*, 22(12):1343–1359, 2008.

- B. Landau and R. Jackendoff. "What" and "where" in spatial language and spatial cognition. *Behavioral and Brain Sciences*, 16(2):217–265, 1993.
- K. M. Lee. Presence, explicated. *Communication Theory*, 14(1):27–50, 2004. doi: 10.1093/ct/14.1.27.
- G. D. Logan and D. D. Sadler. A computational analysis of the apprehension of spatial relations. In P. E. Bloom, M. A. Peterson, L. E. Nadel, and M. F. Garrett, editors, *Language and Space*, pages 493–529. MIT Press, Cambridge, MA, 1996.
- A. M. Lund. Measuring usability with the use questionnaire. *STC Usability SIG Newsletter*, 8(2), 2001.
- E. Marder-Eppstein and E. Perko. ROS package: base local planner, 2012. URL www.ros.org/wiki/base_local_planner.
- R. B. Margolis and C. R. Mynatt. The effects of external and self-administered reward on high base rate behavior. *Cognitive Therapy and Research*, 10(1):109–122, Feb 1986. doi: 10.1007/bf01173387. URL <http://dx.doi.org/10.1007/bf01173387>.
- M. J. Matarić, J. Eriksson, D. J. Feil-Seifer, and C. J. Winstein. Socially assistive robotics for post-stroke rehabilitation. *Journal of NeuroEngineering and Rehabilitation*, 4:5, 2007. doi: 10.1186/1743-0003-4-5.
- Y. Matsusaka, H. Fujii, T. Okano, and I. Hara. Health exercise demonstration robot TAIZO and effects of using voice command in robot-human collaborative demonstration. In *The 18th IEEE International Symposium on Robot and Human Interactive Communication*, pages 472–477, 2009. doi: 10.1109/ROMAN.2009.5326042.
- C. Matuszek, E. Herbst, L. Zettlemoyer, and D. Fox. Learning to parse natural language commands to a robot control system. In *Proc. of the 13th International Symposium on Experimental Robotics (ISER)*, Québec City, Canada, 2012.
- J. C. McCroskey and T. A. McCain. The measurement of interpersonal attraction. *Speech Monographs*, 41(3):261–266, 1974. doi: 10.1080/03637757409375845.
- M. E. T. McMurdo and L. Rennie. A controlled trial of exercise by residents of old people's homes. *Age and Ageing*, 22(1):11–15, 1993. doi: 10.1093/ageing/22.1.11.
- Z. B. Moak and A. Agrawal. The association between perceived interpersonal social support and physical and mental health: Results from the national epidemiological survey on alcohol and related conditions. *Journal of Public Health*, 32(2):191–201, 2010. doi: 10.1093/pubmed/fdp093.
- S. Mohan, A. Mininger, J. Kirk, and J. E. Laird. Learning grounded language through situated interactive instruction. In *AAAI Fall Symposium on Robots Learning Interactively from Human Teachers (RLIHT)*, pages 30–37, 2012.

- M. Montemerlo, J. Pineau, N. Roy, S. Thrun, and V. Verma. Experiences with a mobile robotic guide for the elderly. In *Proceedings of the AAAI National Conference on Artificial Intelligence*, pages 587–592, 2002.
- NeoSpeech. Text-to-speech engine, 2009. URL www.neospeech.com.
- M. N. Nicolescu and M. J. Matarić. Task learning through imitation and human-robot interaction. In K. Dautenhahn and C. L. Nehaniv, editors, *Models and Mechanisms of Imitation and Social Learning in Robots, Humans and Animals: Behavioural, Social and Communicative Dimensions*. Cambridge University Press, 2005.
- Nuance. Dragon naturally speaking, 2013. URL www.nuance.com.
- J. O’Keefe. Vector grammar, places, and the functional role of the spatial prepositions in English. In E. van der Zee and J. Slack, editors, *Representing Direction in Language and Space*, pages 69–85. Oxford University Press, Oxford, 2003.
- M. Pardowitz, S. Knoop, R. Dillmann, and R. D. Zöllner. Incremental learning of tasks from user demonstrations, past experiences, and vocal comments. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 37(2):322–332, April 2007.
- E. S. Paykel. Life events, social support and depression. *Acta Psychiatrica Scandinavica*, 89:50–58, 1994. doi: 10.1111/j.1600-0447.1994.tb05803.x.
- A. A. Pereira, J. Binney, G. A. Hollinger, and G. S. Sukhatme. Risk-aware path planning for autonomous underwater vehicles using predictive ocean models. *Journal of Field Robotics*, 30(5):741–762, 2013. doi: 10.1002/rob.21472.
- R. H. Poresky, C. Hendrix, J. E. Hosier, and M. L. Samuelson. Companion animal bonding scale: Internal reliability and construct validity. *Psychological Reports*, 60(3):743–746, 1987. doi: 10.2466/pr0.1987.60.3.743.
- A. Powers, S. Kiesler, S. Fussell, and C. Torrey. Comparing a computer agent with a humanoid robot. In *Proceedings of the ACM/IEEE International Conference on Human-Robot Interaction, HRI ’07*, pages 145–152, New York, 2007. ACM. doi: 10.1145/1228716.1228736.
- D. K. Roy. Learning visually grounded words and syntax for a scene description task. *Computer Speech & Language*, 16(3-4):353–385, 2002.
- P. E. Rybski, J. Stolarz, K. Yoon, and M. Veloso. Using dialog and human observations to dictate tasks to a learning robot assistant. *Journal of Intelligent Service Robots*, 1(2):159–167, 2008.

- Y. Sandamirskaya, J. Lipinski, I. Iossifidis, and G. Schöner. Natural human-robot interaction through spatial language: A dynamic neural field approach. In *19th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, pages 600–607. IEEE, 2010.
- M. Scheutz, R. Cantrell, and P. Schermerhorn. Toward humanlike task-based dialogue processing for human robot interaction. *AI Magazine*, 32(4):77–84, 2011.
- M. Scopelliti, M. V. Giuliani, and F. Fornara. Robots in a domestic setting: a psychological approach. *Univ Access Inf Soc*, 4(2):146–155, Dec 2005. doi: 10.1007/s10209-005-0118-1. URL <http://dx.doi.org/10.1007/s10209-005-0118-1>.
- S. Shen, N. Michael, and V. Kumar. Autonomous indoor 3d exploration with a micro-aerial vehicle. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 9–15, 2012. doi: 10.1109/ICRA.2012.6225146.
- M. Skubic, D. Perzanowski, S. Blisard, A. Schultz, W. Adams, M. Bugajska, and D. Brock. Spatial language for human-robot dialogs. *IEEE Transactions on SMC Part C: Special Issue on Human-Robot Interaction*, 34(2):154–167, 2004.
- W. W. Spirduso and P. Clifford. Replication of age and physical activity effects on reaction and movement time. *Journal of Gerontology*, 33(1):26–30, 1978.
- S. A. Stansfeld, G. S. Rael, J. Head, M. Shipley, and M. Marmot. Social support and psychiatric sickness absence: A prospective study of British civil servants. *Psychological Medicine*, 27(1):35–48, 1997. doi: 10.1017/S0033291796004254.
- E. Stice, J. Ragan, and P. Randall. Prospective relations between social support and depression: Differential direction of effects for parent and peer support? *Journal of Abnormal Psychology*, 113(1):155–159, 2004. doi: 10.1037/0021-843X.113.1.155.
- L. Talmy. The fundamental system of spatial schemas in language. In B. Hampe, editor, *From Perception to Meaning: Image Schemas in Cognitive Linguistics*, pages 199–234. Mouton de Gruyter, Berlin, 2005.
- A. Tapus, C. Țăpuș, and M. J. Matarić. User-robot personality matching and assistive robot behavior adaptation for post-stroke rehabilitation therapy. *Intelligent Service Robotics*, 1(2):169–183, 2008. doi: 10.1007/s11370-008-0017-4.
- A. Tapus, C. Tapus, and M. J. Matarić. The use of socially assistive robots in the design of intelligent cognitive therapies for people with dementia. In *International Conference on Rehabilitation Robotics*, pages 924–929, Kyoto, Japan, 2009. IEEE. doi: 10.1109/ICORR.2009.5209501.

- S. Tellex, T. Kollar, S. Dickerson, M. R. Walter, A. G. Banerjee, S. Teller, and N. Roy. Approaching the symbol grounding problem with probabilistic graphical models. *AI Magazine*, 32(4):64–76, 2011.
- V. S. Thomas and P. A. Hageman. Can neuromuscular strength and function in people with dementia be rehabilitated using resistance-exercise training? Results from a preliminary intervention study. *The Journals of Gerontology Series A: Biological Sciences and Medical Sciences*, 58(8):M746–M751, 2003. doi: 10.1093/gerona/58.8.M746.
- D. C. Tway. *A construct of trust*. PhD thesis, The University of Texas at Austin, 1994.
- R. J. Vallerand. The effect of differential amounts of positive verbal feedback on the intrinsic motivation of male hockey players. *Journal of Sport Psychology*, 5(1):100–107, 1983.
- R. J. Vallerand and G. Reid. On the causal effects of perceived competence on intrinsic motivation: A test of cognitive evaluation theory. *Journal of Sport Psychology*, 6(1):94–102, 1984.
- D. Vasquez, P. Stein, J. Rios-Martinez, A. Escobedo, A. Spalanzani, and C. Laugier. Human aware navigation for assistive robotics. In J. P. Desai, G. Dudek, O. Khatib, and V. Kumar, editors, *Experimental Robotics*, volume 88 of *Springer Tracts in Advanced Robotics*, pages 449–462. Springer International Publishing, 2013. ISBN 978-3-319-00064-0. doi: 10.1007/978-3-319-00065-7_31. URL http://dx.doi.org/10.1007/978-3-319-00065-7_31.
- M. Veloso, J. Biswas, B. Coltin, S. Rosenthal, S. Brandao, T. Merili, and R. Ventura. Symbiotic-autonomous service robots for user-requested tasks in a multi-floor building. In *IROS'12 Workshop on Cognitive Assistive Systems*, Algarve, Portugal, 2012.
- K. Wada, T. Shibata, T. Saito, and K. Tanie. Analysis of factors that bring mental effects to elderly people in robot assisted activity. In *Proceedings of the International Conference on Intelligent Robots and Systems*, volume 2, pages 1152–1157, 2002. doi: 10.1109/IRDS.2002.1043887.
- J. Wainer, D. J. Feil-Seifer, D. A. Shell, and M. J. Matarić. The role of physical embodiment in human-robot interaction. In *IEEE Proceedings of the International Workshop on Robot and Human Interactive Communication*, pages 117–122, 2006. doi: 10.1109/ROMAN.2006.314404.
- J. Wainer, D. J. Feil-Seifer, D. A. Shell, and M. J. Matarić. Embodiment and human-robot interaction: A task-based perspective. In *IEEE Proceedings of the International Workshop on Robot and Human Interactive Communication*, pages 872–877, 2007. doi: 10.1109/ROMAN.2007.4415207.

- S. Waldherr, S. Thrun, R. Romero, and D. Margaritis. Template-based recognition of pose and motion gestures on a mobile robot. In *Proceedings of the National Conference on Artificial Intelligence*, pages 977–982, 1998.
- R. S. Weinberg and J. Ragan. Effects of competition, success/failure, and sex on intrinsic motivation. *Research Quarterly*, 50(3):503–510, 1979.
- Willow Garage. Robot Operating System (ROS), 2013. URL www.ros.org.
- C. R. Wren, A. Azarbayejani, T. Darrell, and A. P. Pentland. Pfnder: Real-time tracking of the human body. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7):780–785, 1997. doi: 10.1109/34.598236.
- M. Zuckerman, J. Porac, D. Lathin, R. Smith, and E. L. Deci. On the importance of self-determination for intrinsically motivated behavior. *Personality and Social Psychology Bulletin*, 4:443–446, 1978.