

Evaluation of a Spatial Language Interpretation Framework for Natural Human-Robot Interaction with Older Adults*

Juan Fasola and Maja J Matarić, *Fellow, IEEE*

Abstract—We present the design and analysis of a multi-session user study conducted with older adults to evaluate the effectiveness of a human-robot interaction (HRI) framework utilizing a neuroscience-inspired spatial language interpretation approach. The user study was designed to: 1) evaluate the effectiveness and feasibility of the spatial language interpretation framework with target 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 HRI framework. In addition, an expanded study was conducted on Amazon’s Mechanical Turk to gather further natural language instructions from general users for analysis. The results of both studies, evaluated across a variety of objective performance and participant evaluation measures, demonstrate the feasibility of the approach and its effectiveness in interpreting and following natural language instructions from target users, achieving high task success rates and high participant evaluations across multiple measures.

I. INTRODUCTION

Spatially-oriented tasks such as fetching and moving objects are commonly referenced among the tasks most desired by older adults for household service robots to perform [1]. To achieve these types of tasks, autonomous service robots will need to be capable of interacting with and learning from non-expert users in a manner that is both natural and practical for the users. In particular, these robots will need to be capable of understanding natural language instructions in order to accomplish user-defined tasks and receive feedback and guidance on task execution. This capability is especially important in assistive domains, where robots are interacting with people with special needs, as the users may not be able to teach tasks and/or provide feedback by demonstration.

Spatial language plays a key role in instruction-based natural language communication, and is especially relevant for household object pick-and-place tasks whose goals are to satisfy user-desired spatial relations between specified figure and reference objects [2], [3]. Previous work that has investigated the use and representation of spatial language in HRI includes Skubic et al. [4], who developed a robot capable of understanding and relaying static spatial relations in instruction and production tasks. The use of computational models of static relations has also been explored in systems for tabletop pick-and-place tasks [5], and for visually situated dialogue [6]. These works implemented pre-defined

models of spatial relations; however, researchers have also developed systems capable of learning static spatial relations automatically from training data (e.g., [7]).

Recent work has also investigated approaches to interpreting natural language instructions involving dynamic spatial relations (DSRs). Tellex et al. [8] developed a probabilistic graphical model to infer object pick-and-place tasks for execution by a forklift robot. Kollar et al. [9] employed a Bayesian approach for interpreting route directions on a mobile robot. Cantrell et al. [10] demonstrated an approach to learning action verbs through human-robot dialogues.

In this paper, we contribute an evaluation of our spatial language-based HRI framework, proposed in [11], extended to incorporate instruction semantic ambiguity resolution through human-robot dialogue procedures. We implemented the modified framework on a fully autonomous mobile robot and tested the approach with target users. Specifically, we present the design and results of a user study conducted with older adults to: 1) evaluate the effectiveness and feasibility of the 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 further improvements to the spatial language framework. Lastly, we present results of an expanded study conducted on Amazon’s Mechanical Turk (AMT) to increase the evaluation corpus of natural language instructions and to test the modified framework’s effectiveness with a general user population.

II. SPATIAL LANGUAGE FRAMEWORK

The spatial language-based HRI framework evaluated in this paper follows the methodology outlined in our prior work [11] for representing spatial language to enable natural language-based interaction with non-expert users. The principal aspect of the approach is the encoding of spatial language within the robot *a priori* as primitives. Static spatial relation primitives are represented using the semantic field model proposed by O’Keefe [12], where the semantic fields of static prepositions, parameterized by figure and reference objects, assign weight values to points in the environment depending on how accurately they capture the meaning of the preposition (e.g., for the static spatial preposition “near”, points closer to an object have higher weight). In its original form, the framework contains five system modules: the syntactic parser, noun phrase (NP) grounding, semantic interpretation, planning, and action modules. The framework is capable of grounding hierarchical NPs probabilistically using semantic fields, and employs a Bayesian approach

* Research supported by National Science Foundation grants IIS-0713697, CNS-0709296, and IIS-1117279.

J. Fasola is with the University of Southern California, Los Angeles, CA 90089 USA fasola@usc.edu

M. J. Matarić is with the University of Southern California, Los Angeles, CA 90089 USA matari@usc.edu

to infer the semantics of given instructions probabilistically using a database of learned mappings from input observations to instruction meanings. The four observation inputs include: the *verb* and *preposition* utilized, and the associated grounding types observed for the *figure* and *reference objects*. The resulting semantic output of the module includes: the *command* type, the *DSR* type, and the *static spatial relation* (if available). For a complete description of the original framework modules, we refer the reader to [11].

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 the robot plans a task solution and executes it; or 2) the robot is unable to fully interpret the semantics of the instruction (e.g., due to the presence of unknown words/phrases). In the latter case, our approach extends the framework from [11] by introducing a human-robot dialogue procedure that is initiated by the robot with a clarification query that is posed to the user to resolve any grounding errors and/or confirm the command type (in case the command type cannot be inferred with high probability). In this case, the user and robot 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 progress towards resolving the ambiguities is deemed unsatisfactory (e.g., the maximum number of clarification queries is reached, experimentally set to 6 per ground - a threshold set *a priori*). This back-and-forth human-robot dialogue towards resolving the meaning of a single user instruction is referred to as a *dialogue round*, and represents an important addition to the framework.

III. STUDY DESIGN

To evaluate the effectiveness of the spatial language interpretation framework with target users, we conducted a user study with older adults that consisted of two conditions: the Virtual Robot condition and the Physical Robot condition. Both conditions were designed to engage the user in human-robot dialogue, and more specifically, to evoke spatial language instructions from the participant so as to test the effectiveness of the robot in interpreting and following the instructions in accordance with the context of the current discourse and environment. The study design was within-subject, with all participants engaging in both conditions (1 session per condition) approximately one week apart. The order of appearance of the conditions was fixed, with the Virtual Robot condition appearing first as it includes a training session to help the users quickly familiarize themselves with the robot’s capabilities. Each 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. It should be noted that both conditions were designed to be complementary; hence the aim of the study was not to evaluate the effectiveness of one condition over the other, but rather to evaluate the effectiveness of the spatial language framework in responding to user commands under a variety of different



(a)



(b)

Fig. 1. (a) Virtual Robot condition setup; (b) 2D computer simulated home environment, with example robot task execution path shown for instruction “Pick up the medicine in the guest bathroom”.

test scenarios. The following subsections describe the two conditions in greater detail.

A. 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 resolution 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 [13]. In the interaction setup, the user is seated in a chair facing a computer display projecting the simulated home environment, which allows the user to view the robot’s actions and verify the correctness of the robot’s task execution. Commands are issued by the user to the robot using natural language speech, and the spoken instructions are manually transcribed via keyboard input by the experimenter in real-time and sent to the robot for interpretation. The interaction setup is shown in Fig. 1 (a), and the simulated home environment displaying the robot (Light Green), the simulated user position (Purple),



Fig. 2. Example target objects and placement locations for scenarios 1 and 2 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.

and various movable objects drawn as smaller colored circles (described below), are shown in Fig. 1 (b).

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, they are 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. 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 movable objects: water bottles (Blue), medicine (Pink), and books (Green). The object list and both environment naming maps are placed in front of the user, as can be seen in the interaction setup shown in Fig. 1 (a).

The Virtual Robot condition examines two primary scenarios for data collection and evaluation purposes: 1) *object identification*; and 2) *placement location identification*. The two 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 first 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 (pick-up) task to the robot. The target object is identified to the user nonverbally through the simulator by a highlighted red circle surrounding the object. Similarly, in the second scenario, the robot is holding the object to be placed, and the target location is expressed to the user nonverbally through a highlighted red region in the home environment. For both scenarios, multiple such target objects/locations are evaluated, with each provided sequentially after the user successfully instructs the robot to perform each signaled task. In total, 28 task settings were evaluated (17 pick-up, 11 placement). Fig. 2 shows screenshots of example task settings for both scenarios.

B. 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 includes: 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 Fig. 3.

The robot platform used in this condition is Bandit, a humanoid torso robot mounted on a MobileRobots Pioneer base. 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 Fig. 3 (a).

The Physical Robot condition examines two primary scenarios for data collection and evaluation purposes: 1) *object identification and placement*; and 2) *task-oriented instruction*. The first interaction scenario can be thought of as a combination of both scenarios from the Virtual Robot condition, as the user must command the robot to move a

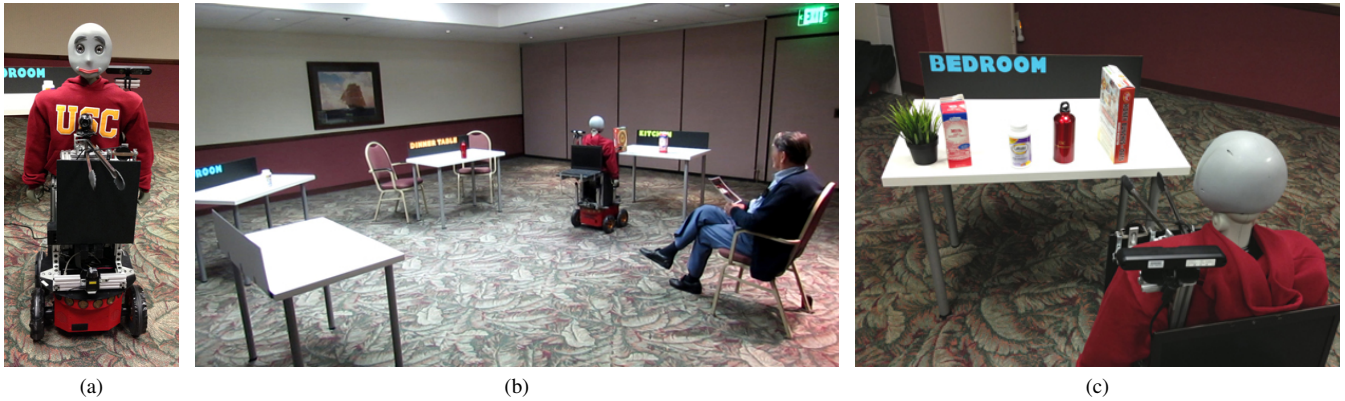


Fig. 3. Physical Robot condition. (a) The physical robot platform; (b) view of interaction setup with labeled tables representing typical household areas (coffee table, bedroom, dinner table, and kitchen); (c) example household items used in the study (from left to right: plant, milk, medicine, bottle, cereal).



Fig. 4. Example task photographs provided to the user displaying task goal states.

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 nonverbally (using two green sticky notes) prior to the start of the scenario.

In the second interaction scenario, the user is given a task that the robot needs to accomplish, and asked to provide instructions to the robot. The task is relayed to the user nonverbally: the user is shown a photograph of the target state of the environment to be achieved for the task to be completed successfully. Tasks were chosen to require multiple pick-and-place instructions to achieve the goal state, and involved the movement of one or more objects onto a specific table (household area). The user is encouraged to match the relative object positions as closely as possible to the goal state displayed in the image. Two example task goal states provided to the user, with specified target object goal locations, are shown in Fig. 4.

IV. PARTICIPANT STATISTICS

Through a partnership with be.group, a senior living organization, we recruited elderly individuals to participate in the study, 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. Participants' ages ranged from 71-97; the average age was 82 (S.D. = 7.42).

V. MEASURES

A. Objective Measures

The objective measures collected (17 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 the framework. Many of our objective measures employed the number of dialogue rounds as a normalizing factor.

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); *round success rate* (percentage of rounds that ended in a task execution by the robot); *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 anaphoric references per round* (e.g., it, this/that, him/her), *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, or four instructions expressed within a single utterance, respectively) given among all participants. These measures were chosen to evaluate aspects of the 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 within 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 the natural language instructions and responses expressed by the users.

B. Subjective Measures

After each session of the study the participants were 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*, the perceived *value/ usefulness of the interaction*, the *intelligence of the robot*, and the *social presence of the robot*. The measures were characterized by a series of related adjectives, and participants were asked to rate how well each adjective described the interaction/robot on a 10-point scale, anchored by “Describes Very Poorly” (1) and “Describes Very Well” (10).

To measure the perceived ease of use of the robot system, we administered the USE questionnaire [14]. 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*).

VI. RESULTS

A. Virtual Robot Condition Results

The collected statistics regarding the performance of the 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 the spatial language framework to correctly interpret and follow natural language instructions provided during user discourse.

The round success rate achieved by the 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 the 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 the 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.

TABLE I
RESULTS OF INTERACTION WITH PARTICIPANTS ($N = 19$) IN
VIRTUAL/PHYSICAL ROBOT CONDITIONS

Objective Measure	Virtual Robot Condition	Physical Robot Condition
Task Success Rate	78.6% (14.5)	87.4% (6.8)
Task Success Rate (Repeated Attempts)	82.8% (12.8)	98.0% (3.3)
Round Success Rate	84.1% (11.9)	92.8% (6.8)
Total Dialogue Rounds	49.2 (13.7)	25.8 (7.2)
Clarification Queries Per Round	0.91 (0.37)	0.47 (0.22)
Yes/No Queries	9.35 (5.8)	3.7 (3.1)
Yes/No Answers	7.4 (4.5)	3.2 (2.8)
Number of References Per Round	0.15 (0.23)	0.34 (0.26)
Total References Used	6.9 (9.5)	8.5 (6.0)
Max References Used in Session	38	22

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 ($M=0.91$, $S.D.=0.37$). Table I provides a summary of the collected statistics for all of the objective measures captured during the Virtual Robot condition of the study.

B. Physical Robot Condition Results

The results of the interaction of the participants with the spatial language framework in the Physical Robot condition were similar to those observed in the Virtual Robot condition, with improved performance overall. Table I 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. By design, the Physical Robot condition represented the more natural scenario for the user: the user and robot were situated in the same physical space, and the user could easily observe/identify the spatial locations of all objects in the room (characteristics that 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.

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. $M=0.15$,

TABLE II

EXAMPLE INSTRUCTIONS BY PARTICIPANTS IN THE PHYSICAL ROBOT CONDITION WITH INFERENCE RESULTS FOR INSTRUCTION SEQUENCES

#	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 [the 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 [the bottle]) Put [it] between [[the cereal] and [the milk]] - ("it" referent resolved to [the 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) and uniquely grounded NPs surrounded by brackets.

S.D.=0.23). This result could be due to the increased naturalness of the environment, as previously discussed, but also due to the increased complexity of the tasks presented to the user. For instance, 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 formulate 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. In total, just over one quarter (27.8%) of the instruction sequences provided by the study participants were expressed in utterances containing two or more instructions. The exact statistics regarding the number of instruction sequences given of length 1-4 were as follows: Length 1) 376 = 72.2%; 2) 122 = 23.4%; 3) 22 = 4.2%; 4) 1 = 0.2%. 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 II.

C. 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 of the different verbs, path prepositions, and static prepositions employed by the participants when issuing instructions to the robot in both study conditions. Table III shows the counts for each verb, path, and static preposition used in all of the $N = 1239$ valid grammatical instructions issued by the participants that were interpreted by the robot during interaction, and 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: the *command* type (shown with verbs), and the *dynamic spatial relation (DSR)* and *static spatial relation (SSR)* types (shown with path prepositions). 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 the spatial language methodology tested, 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 Table I). The semantic interpretation module of the 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 III 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

TABLE III
SPATIAL LANGUAGE STATISTICS OF VERB, PATH, AND STATIC PREPOSITION USAGE IN $N = 1239$ TOTAL INSTRUCTIONS
GIVEN BY PARTICIPANTS DURING STUDY (ACROSS BOTH CONDITIONS)

Verb	Count (%)	Cmd.	% Inf.	Path Prep.	Count (%)	DSR	SSR	Static Prep.	Count (%)
pick	414 (33.41%)	AO	100%	up	406 (32.64%)	<i>up</i>	<i>none</i>	on	240 (29.20%)
put	311 (25.10%)	OM	92%	to	347 (27.89%)	<i>to</i>	<i>at</i>	from	174 (21.17%)
take	130 (10.49%)	OM	88%	on	220 (17.68%)	<i>to</i>	<i>on</i>	in	156 (18.98%)
go	126 (10.17%)	RM	100%	(none)	143 (11.50%)	<i>to</i>	<i>at</i>	of	143 (17.40%)
move	90 (7.26%)	OM	89%	in	50 (4.02%)	<i>to</i>	<i>in</i>	by	18 (2.19%)
bring	82 (6.62%)	OR	75%	from	16 (1.29%)	<i>to</i>	<i>out</i>	near	18 (2.19%)
place	29 (2.34%)	OM	92%	between	7 (0.56%)	<i>to</i>	<i>between</i>	to the right of	13 (1.58%)
get	24 (1.94%)	OR	61%	into	6 (0.48%)	<i>to</i>	<i>in</i>	next to	11 (1.34%)
give	7 (0.56%)	OR	100%	onto	6 (0.48%)	<i>to</i>	<i>on</i>	to ¹	9 (1.09%)
grab	5 (0.40%)	AO	100%	the left of	6 (0.48%)	<i>to</i>	<i>left-of</i>	at	8 (0.97%)
turn	4 (0.32%)	RM	100%	front of	5 (0.40%)	<i>to</i>	<i>front-of</i>	close to	5 (0.61%)
remove	3 (0.24%)	OM	67%	near	5 (0.40%)	<i>to</i>	<i>near</i>	to the left of	5 (0.61%)
wash	3 (0.24%)	RM	67%	next to	5 (0.40%)	<i>to</i>	<i>near</i>	between	4 (0.49%)
find	2 (0.16%)	RM	100%	the right of	5 (0.40%)	<i>to</i>	<i>right-of</i>	left of	3 (0.36%)
leave	2 (0.16%)	N/A	0%	down	4 (0.32%)	<i>down</i>	<i>none</i>	on top of	3 (0.36%)
hold	1 (0.08%)	RM	100%	by	3 (0.24%)	<i>to</i>	<i>near</i>	the left of	3 (0.36%)
keep	1 (0.08%)	OM	100%	inside	2 (0.16%)	<i>to</i>	<i>in</i>	off	2 (0.24%)
milk	1 (0.08%)	OM	100%	on top of	2 (0.16%)	<i>to</i>	<i>on</i>	right of	2 (0.24%)
need	1 (0.08%)	RM	100%	upon	2 (0.16%)	<i>to</i>	<i>on</i>	behind	1 (0.12%)
want	1 (0.08%)	RM	100%	around	1 (0.08%)	<i>around</i>	<i>none</i>	beside	1 (0.12%)
set	1 (0.08%)	N/A	0%	beside	1 (0.08%)	<i>to</i>	<i>in</i>	front of	1 (0.12%)
stand	1 (0.08%)	N/A	0%	close to	1 (0.08%)	<i>to</i>	<i>near</i>	off of	1 (0.12%)
				of	1 (0.08%)	<i>to</i>	<i>behind-of</i>	over	1 (0.12%)

Note: Word counts (with percentage of total) are shown for all verbs, path, and static 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. 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 noun phrase grounding errors). ¹“to” is considered a static spatial preposition in the framework only when paired with a semantic field specifier noun phrase (e.g., “Pick up [_{NP} the bottle to [_{NP} the left]]”)

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 Table I), the high frequency of only a small set of verbs and path prepositions (see Table III), and the fact that clarification queries often only targeted noun phrase grounding ambiguities, this scenario was rare ($\approx 1.5\%$ of total instructions).

Table III also 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 book on [_{NP} [_{NP} the left side] of the counter]”; “Put the cup down at [_{NP} [_{NP} the front edge] of the coffee table]”).

D. 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 the 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 user evaluations, which is very encouraging. No significant differences were found among the participant ratings of the two study conditions.

The results obtained from the participant evaluations of the system with respect to the USE questionnaire items on usability are also encouraging; the participants rated the household service robot presented in the study highly in terms of ease-of-use, ease of learning, and satisfaction. Fig. 5 displays a summary of the subjective measures captured for the participant evaluation of the interaction and household service robot.

VII. AMT STUDY

To expand the evaluation corpus of natural language instructions and to investigate the effectiveness of the spatial language interpretation HRI framework with general users, we conducted a modification of the original study with users on Amazon’s Mechanical Turk (AMT). The 28 tasks given to the participants were exactly the same as those given in the Virtual Robot condition of the original study (see Section III-A). For each task, we collected 100 instructions from

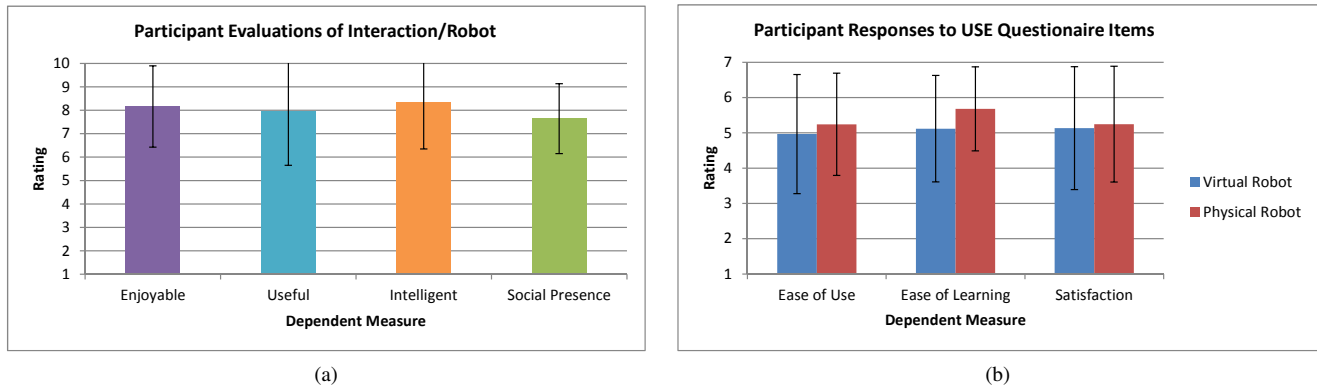


Fig. 5. Participants’ subjective evaluation results. (a) Evaluation of interaction and of service robot; (b) evaluation of interaction with respect to USE questionnaire items.

unique users, resulting in a total of 2800 user-defined natural language instructions. In total, 161 unique users participated in the study, providing instructions for an average of 17.4 tasks (SD = 9.5, Max = 28, Min = 1). The key difference between the two studies was that the AMT interface was not interactive; thus, the robot was unable to engage the user in dialogue to resolve semantic ambiguities. Nevertheless, clarification queries were recorded for analysis.

The results were comparable to those observed in the original user study, demonstrating a task success rate of 80.0% (1775 successes / 2219 task executions), and a round success rate of 79.3% (2219 task executions / 2800 instructions). 581 (20.7%) of the provided instructions did not result in task execution as they contained ambiguities causing the robot to pose a clarification query. Among the 444 failed task executions, 235 were due to incorrect object placement location, 206 were due to incorrect object resolution, and 3 were due to command resolution errors. The failures were largely due to the robot making incorrect context-based grounding inferences from ambiguous instructions (e.g., “Put the bottle on the night table in the guest bedroom.” - ambiguous user instruction as there are two possible night tables, but the robot incorrectly infers unique grounding based on context as only one is visible to the robot).

The high task success rate observed in the expanded study, achieved with a much larger and general user population, confirms the effectiveness and feasibility of our spatial language interpretation approach for interaction with end users.

VIII. CONCLUSIONS

This paper 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 the service robotics-based approach and its effectiveness in interpreting and following natural language instructions from target users, achieving high task success rates and user evaluations in both study conditions. Results of an expanded study on Amazon’s Mechanical Turk with a general user population confirm the effectiveness and feasibility of the approach. Future work

will address current limitations of the framework, including the improvement of context-based grounding and reasoning for ambiguous user instructions.

REFERENCES

- [1] 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 *ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, 2012, pp. 335–342. [Online]. Available: <http://doi.acm.org/10.1145/2157689.2157806>
- [2] L. A. Carlson and P. L. Hill, “Formulating spatial descriptions across various dialogue contexts,” in *Spatial Language and Dialogue*. New York: Oxford University Press, 2009, pp. 89–103.
- [3] B. Landau and R. Jackendoff, ““What” and “where” in spatial language and spatial cognition,” *Behavioral and Brain Sciences*, vol. 16, no. 2, pp. 217–265, 1993.
- [4] 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*, vol. 34, no. 2, pp. 154–167, 2004.
- [5] 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)*. IEEE, 2010, pp. 600–607.
- [6] J. D. Kelleher and F. J. Costello, “Applying computational models of spatial prepositions to visually situated dialog,” *Computational Linguistics*, vol. 35, no. 2, pp. 271–306, 2009.
- [7] D. K. Roy, “Learning visually grounded words and syntax for a scene description task,” *Computer Speech & Language*, vol. 16, no. 3–4, pp. 353–385, 2002.
- [8] 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*, vol. 32, no. 4, pp. 64–76, 2011.
- [9] 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)*. IEEE, 2010, pp. 259–266.
- [10] R. Cantrell, P. Schermerhorn, and M. Scheutz, “Learning actions from human-robot dialogues,” in *Proc. IEEE RO-MAN*. IEEE, 2011, pp. 125–130.
- [11] J. Fasola and M. J. Mataric, “Interpreting instruction sequences in spatial language discourse with pragmatics towards natural human-robot interaction,” in *IEEE International Conference on Robotics and Automation (ICRA)*, June 2014.
- [12] J. O’Keefe, “Vector grammar, places, and the functional role of the spatial prepositions in English,” in *Representing Direction in Language and Space*, E. van der Zee and J. Slack, Eds. Oxford: Oxford University Press, 2003, pp. 69–85.
- [13] NeoSpeech, “Text-to-speech engine,” 2009. [Online]. Available: www.neospeech.com
- [14] A. M. Lund, “Measuring usability with the use questionnaire,” *STC Usability SIG Newsletter*, vol. 8, no. 2, 2001.