

# Using Spatial Semantic and Pragmatic Fields to Interpret Natural Language Pick-and-Place Instructions for a Mobile Service Robot

Juan Fasola and Maja J Matarić

Interaction Lab, Computer Science Department, Viterbi School of Engineering,  
University of Southern California, Los Angeles, CA, USA  
{fasola, mataric}@usc.edu

**Abstract.** We present a methodology for enabling mobile service robots to follow natural language instructions for object pick-and-place tasks from non-expert users, with and without user-specified constraints, and with a particular focus on spatial language understanding. 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 are also presented, including the grammar utilized and our procedure for pruning multiple candidate parses based on context. The paper concludes with an evaluation of our approach implemented on a simulated mobile robot operating in both 2D and 3D home environments.

## 1 Introduction

Spatially-oriented tasks such as fetching and moving objects are commonly referenced among the tasks most desired by non-expert users for household service robots to perform [3]. In accomplishing 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 learn new tasks and receive guidance and feedback on task execution.

Spatial language plays an important role in instruction-based natural language communication, and is especially relevant for object pick-and-place tasks whose purpose lies in achieving desired spatial relations between specified figure and reference objects [12]. Spatial relations, both dynamic and static, expressed in language are often expressed by prepositions [1]. Therefore, the ability for robots to understand and differentiate between spatial prepositions in spoken language is crucial for their interaction with the user to be successful. However, in performing mobile manipulation tasks, service robots must consider not only the semantics of spatial relations within the natural language instruction, but also the *pragmatics* of the task itself. These consist in the implicit (unvoiced) constraints/specifications that accompany the

spoken instructions and further specify the meaning the speaker intends to convey; which can come from context, prior knowledge, norms, etc.

In this paper, we extend upon our previous work and present a general methodology for enabling mobile service robots to follow natural language instructions for object pick-and-place tasks from non-expert users, with and without user-specified constraints. Our approach is capable of addressing both the semantic and pragmatic properties of natural language instructions, and in particular, proposes a novel computational field representation for the incorporation of spatial pragmatic constraints in mobile manipulation task planning. Furthermore, we describe extensions to our spatial language understanding framework, including a revised grammar for the inclusion of figure objects, and a procedure for pruning multiple candidate parses based on context. We conclude the paper with an evaluation of our approach implemented on a simulated mobile robot in both 2D and 3D home environments to demonstrate the generalizability and usefulness of our approach towards real world applications.

## 2 Related Work

Previous work that has investigated the use and representation of spatial language in human-robot interaction (HRI) includes the work of Skubic et al. [9], who demonstrated a robot capable of understanding and relaying static spatial relations in natural language instruction and production tasks. Computational field models of static relations have also been used in systems for object pick-and-place tasks on a tabletop [8], and for visually situated dialogue [11]. These works implemented pre-defined models of spatial relations, however, researchers have also designed systems capable of learning these types of static spatial relations automatically from training data (e.g., [10]).

Towards enabling HRI-based object pick-and-place instruction understanding and task execution, various approaches have previously been explored. Representative work includes approaches tested on mobile robots with grippers (e.g., [13]), and on humanoid robots with articulated arms (e.g., [14]). Approaches capable of interpreting explicit spatial relations in spoken language instructions include the work of Sisbot et al. [15], who designed a perspective-taking approach to situational awareness using discretized representations of spatial relations, and Brenner et al. [16] who alternatively employed a continuous potential field model. Our approach extends upon this related work by modeling both static and dynamic spatial relations, and by incorporating the implicit pragmatics of the instruction for robot task planning.

Recent work has also investigated approaches to interpreting natural language instructions involving dynamic spatial relations (DSRs). Tellex et al. [6] developed a probabilistic graphical model to infer object pick-and-place tasks for execution by a forklift robot from natural language commands. Kollar et al. [7] employed a Bayesian approach for interpreting route directions on a mobile robot. 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 system designer 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. In addition, as the meanings of spatial relations are learned via task demonstration, the semantics and pragmatics of the instructed tasks are not clearly separable, thus limiting the ability of the robot to generalize the learned task execution to new contexts with different pragmatic constraints. Our methodology employs domain-generalizable spatial relations as primitives, and probabilistic reasoning for the grounding, semantic, and pragmatic interpretation of phrases, thereby allowing for context-based instruction understanding and user modifiable robot execution paths.

### 3 Spatial Language-Based HRI Framework

We have developed a framework for human-robot interaction that incorporates our methodology for representing spatial language, and DSRs in particular, to enable natural language-based interaction with non-expert users. Fundamentally, our approach encodes spatial language within the robot *a priori* as primitives. Static spatial relation primitives are represented using the *semantic field* model proposed by O’Keefe [2], 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).

The software framework 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. In this section we will provide a brief overview of the framework design and module functionality. For a complete description of the software modules, we refer the reader to [4, 5].

The syntactic parser represents the entry point of our software architecture, as it responsible for parsing the user-given natural language instruction into a format that the remaining modules can interpret. 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, grounding hierarchical NPs probabilistically using semantic fields. All groundings are then passed on to the semantic interpreter for final instruction meaning association. The semantic interpretation module utilizes a Bayesian approach to infer the semantics of the given instruction 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). The command type is domain-specific, and may include commands such as: robot movement, object movement, 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 to accommodate the requirements of the specific application domain.

## 4 Parsing Spatial Language Instructions with Figure Objects

### 4.1 Grammar Extension for Figure Objects

The syntactic parser of our human-robot interaction framework has previously been shown to successfully interpret natural language English instructions involving spatial language [4, 5]. However, while our previous parser was able to interpret a variety of spatial language instructions, including those with hierarchical noun phrases, the grammar utilized was 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 noun phrase parameters. The extended constituency rules, which define valid sentences (S), noun phrases (NP), and terminal noun phrases (N'), are presented in Table 1.

**Table 1.** 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<sub>s</sub> = Static Preposition, P<sub>p</sub> = Path Preposition, N = Noun, A = Adjective, Det = Determiner

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<sub>s</sub> and P<sub>p</sub>, 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]”).

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

**Table 2.** 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

**Fig. 1.** Apartment environment with cup locations and *on*(bookcase) semantic field shown

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] [Ps on] [NP the bookcase in the kitchen]
- (2) [v Put] [NP the cup on the bookcase] [Ps 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 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 2.

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 Fig. 1. 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.

## 5 Object Pick-and-Place Movement Planning

Once the optimal parse is determined 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. In this section we will describe 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.

### 5.1 Object Pick Up Planning with Grasp Fields

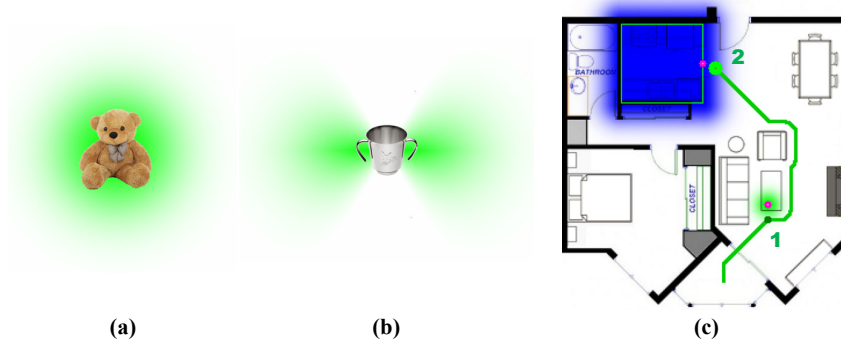
For mobile service robots with onboard manipulators, satisfying an object placement request typically involves first gaining possession of the object in question by picking it up. To address this task, our approach employs a domain-dependent robot-object grasp field centered on the figure object and 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 on the robot's grasp radius. Fig. 2 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.

### 5.2 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  $\geq 90\%$  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.



**Fig. 2.** (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

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. Fig. 2c 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.

### 5.3 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.

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 Fig. 2c 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. Fig. 3a shows an example pragmatic field for object placement preference on surfaces (with weighted surfaces), and Fig. 3b illustrates the use of this field in our described robot task example.



**Fig. 3.** (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

## 6 Evaluation

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 3. 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 Fig. 4.

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 Fig. 3a. 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.





Fig. 4. Robot execution paths for test runs (a) Run #1 (no constraints) (b) Run #2 (constraints)

Table 3. 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>

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. Fig. 5 shows a snapshot of a successful run of the robot executing a pick-and-place task within the 3D household environment. Future work includes testing our approach on the actual PR2 with end users.

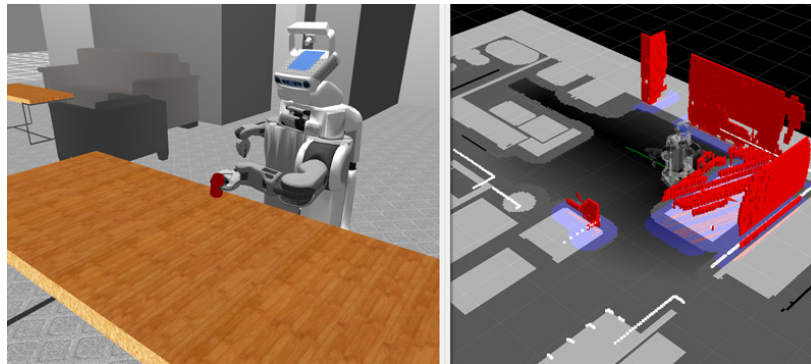


Fig. 5. PR2 robot executing task for “Put the Coke can on the coffee table” in 3D household environment using Gazebo simulator/ROS framework

**Acknowledgements.** This work was supported by the National Science Foundation under Grants IIS-0713697, CNS-0709296, and IIS-1117279.

## References

1. Landau, B., Jackendoff, R.: What and where in spatial language and spatial cognition. *Behavioral and Brain Sciences*, 16(2):217–265 (1993)
2. O’Keefe, J.: Vector grammar, places, and the functional role of the spatial prepositions in English, E. van der Zee and J. Slack, Eds. Oxford: Oxford University Press (2003)
3. Beer, J.M., Smarr, C., Chen, T.L., Prakash, A., Mitzner, T.L., Kemp, C.C., Rogers, W.A.: The domesticated robot: design guidelines for assisting older adults to age in place. In Proc. ACM/IEEE Intl. Conf. on Human-Robot Interaction, 335–342, Boston, MA (2012)
4. Fasola, J., Matarić, M.J.: Using Semantic Fields to Model Dynamic Spatial Relations in a Robot Architecture for Natural Language Instruction of Service Robots. In IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), Tokyo, Japan (2013)
5. Fasola, J., Matarić, M.J.: Modeling Dynamic Spatial Relations with Global Properties for Natural Language-Based Human-Robot Interaction. In Proc. IEEE International Symposium on Robot and Human Interactive Communication (Ro-Man), Aug 26-29, Korea (2013)
6. Tellex, S., Kollar, T., Dickerson, S., Walter, M. R., Banerjee, A. G., Teller, S., Roy, N.: Approaching the Symbol Grounding Problem with Probabilistic Graphical Models. *AI Magazine*. 32(4): 64-76 (2011)
7. Kollar, T., Tellex, S., Roy, D., Roy, N.: Toward Understanding Natural Language Directions. In Proc. ACM/IEEE Int’l Conf. on Human-Robot Interaction (HRI), 259–266 (2010)
8. Sandamirskaya, Y., Lipinski, J., Iossifidis, I., Schöner, G.: Natural human-robot interaction through spatial language: a dynamic neural field approach. 19th IEEE International Symposium on Robot and Human Interactive Communication, pp. 600-607 (2010)
9. Skubic, M., Perzanowski, D., Blisard, S., Schultz, A., Adams, W., Bugajska, M., Brock, D.: Spatial Language for Human-Robot Dialogs. *IEEE Transactions on SMC Part C*, Special Issue on Human-Robot Interaction, 34(2):154-167 (2004)
10. Roy, D.K.: Learning visually grounded words and syntax for a scene description task. In *Computer Speech & Language*. 16(3-4): 353-385 (2002)
11. Kelleher, J.D., Costello, F.J.: Applying Computational Models of Spatial Prepositions to Visually Situated Dialog. *Computational Linguistics*, 35(2):271-306 (2009)
12. Carlson, L.A., Hill, P.L.: Formulating spatial descriptions across various dialogue contexts. In K. Coventry, T. Tenbrink, & J. Bateman (Eds.), *Spatial Language and Dialogue*, 88-103. New York, NY: Oxford University Press Inc. (2009)
13. Nicolescu, M., Matarić, M.J.: Natural Methods for Robot Task Learning: Instructive Demonstration, Generalization and Practice. In International Joint Conference on Autonomous Agents and Multiagent Systems (Agents), pp. 241-248, Melbourne, Australia (2003)
14. Dominey, P.F., Metta, G., Nori, F., Natale, L.: Anticipation and initiative in human-humanoid interaction. In 8th IEEE-RAS International Conference on Humanoid Robots, pp. 693-699 (2008)
15. Sisbot, E.A., Ros, R., Alami, R.: Situation Assessment for Human-Robot Interactive Object Manipulation. In International Symposium in Robot and Human Interactive Communication (RO-MAN), pp. 15-20 (2011)
16. Brenner, M., Hawes, N., Kelleher, J., Wyatt, J.: Mediating between qualitative and quantitative representations for task-orientated human-robot interaction. In Proc. of the 20th International Joint Conference on Artificial Intelligence (IJCAI), pp. 2072-2077 (2007)