



Matteo Baldoni, Jamal Bentahar  
John Lloyd, M. Birna van Riemsdijk (eds.)

## Declarative Agent Languages and Technologies

*Seventh International Workshop, DAL.T 2009  
Budapest, Hungary, May 11th, 2009*

*Workshop Notes*

---

DALT 2009 Home Page:  
<http://www.di.unito.it/~baldoni/DALT-2009/>



## Preface

The workshop on Declarative Agent Languages and Technologies (DALT), in its *seventh edition* this year, is about investigating, studying, and using the declarative paradigm in agent and in multi-agent systems. As one of the well-established workshops in the multi-agent systems area, DALT aims at providing a forum for researchers interested in linking theory to practical applications by combining declarative and formal approaches with engineering and technology aspects of agents and multi-agent systems. The purpose of the DALT workshop is both to advance research on declarative approaches and technologies for agent systems and to foster the scientific exchange among researchers, practitioners, and engineers interested in building complex and practical agent-based systems.

Specification, verification, and validation of agents and multi-agent systems using different formal techniques, such as computational logics, model checking, constraint logic, game theory, and semantic models, along with evaluation of declarative approaches and industrial experiences are the key issues DALT aims to promote. Moreover, DALT has traditionally fostered the development of declarative approaches to engineering agent-based systems. This theme is driven by important application areas, such as the semantic web and service-oriented computing, security, and electronic contracting. We are interested in the foundations, methodologies and techniques for building such systems, in reports on experiences gathered in using them, and in experimental studies analysing their performance.

In the tradition of DALT, the 2009 meeting is being held as a satellite workshop of AAMAS 2009, the 8th International Joint Conference on Autonomous Agents and Multiagent Systems, in Budapest, Hungary. Following the success of DALT 2003 in Melbourne (LNAI 2990), DALT 2004 in New York (LNAI 3476), DALT 2005 in Utrecht (LNAI 3904), DALT 2006 in Hakodate (LNAI 4327), DALT 2007 in Honolulu (LNAI 4897), DALT 2008 in Estoril (LNAI 5397), DALT will aim at providing a discussion forum to both (i) support the transfer of declarative paradigms and techniques to the broader community of agent researchers and practitioners, and (ii) to bring the issue of designing complex agent systems to the attention of researchers working on declarative languages and technologies.

This edition of DALT received *seventeen* high quality submissions, describing works by researchers coming from 12 different countries, *eleven* of which have been selected by the Programme Committee and are included in this volume. Each paper received at least three reviews in order to supply the authors with helpful feedback that could stimulate the research as well as foster discussion. DALT 2009 post-proceedings will appear soon as a volume in the Lecture Notes in Artificial Intelligence by Springer.

We would like to thank all authors for their contributions, the members of the Steering Committee for the valuable suggestions and support, and the mem-

VI

bers of the Programme Committee for their excellent work during the reviewing phase.

March 18th, 2009

Matteo Baldoni  
Jamal Bentahar  
John Lloyd  
M. Birna van Riemsdijk

## Workshop Organisers

Matteo Baldoni	University of Torino, Italy
Jamal Bentahar	Concordia University, Canada
John Lloyd	Australian National University, Australia
M. Birna van Riemsdijk	Delft University of Technology, The Netherlands

## Programme Committee

Thomas Agotnes	Bergen University College, Norway
Marco Alberti	University of Ferrara, Italy
Natasha Alechina	University of Nottingham, UK
Cristina Baroglio	University of Torino, Italy
Rafael Bordini	University of Durham, UK
Jan Broersen	University of Utrecht, The Netherlands
Federico Chesani	University of Bologna, Italy
Amit Chopra	North Carolina State University, USA
Keith Clark	Imperial College London, UK
James Harland	RMIT University, Australia
Andreas Herzig	Paul Sabatier University, France
Koen Hindriks	Delt University of Technology, The Netherlands
Shinichi Honiden	National Institute of Informatics, Japan
Yves Lespérance	York University, Canada
Alessio Lomuscio	Imperial College London, UK
Viviana Mascardi	University of Genova, Italy
Nicolas Maudet	University of Paris-Dauphine, France
John-Jules Ch. Meyer	Utrecht University, The Netherlands
Peter Novak	Clausthal university of Technology, Germany
Enrico Pontelli	New Mexico State University, USA
Azzurra Ragone	Polytechnic of Bari, Italy
Chiaki Sakama	Wakayama University, Japan
Guillermo Simari	Universidad Nacional del Sur, Argentina
Tran Cao Son	New Mexico State University, USA
Eugenia Ternovska	Simon Fraser University, Canada
Wamberto Vasconcelos	University of Aberdeen, UK
Marina De Vos	University of Bath, UK
Michael Winikoff	University of Otago, New Zealand

### **Steering Committee**

Matteo Baldoni	University of Torino, Italy
Andrea Omicini	University of Bologna-Cesena, Italy
M. Birna van Riemsdijk	Delft University of Technology, The Netherlands
Tran Cao Son	New Mexico State University, USA
Paolo Torroni	University of Bologna, Italy
Pinar Yolum	Bogazici University, Turkey
Michael Winikoff	University of Otago, New Zealand

### **Additional Reviewers**

Aaron Hunter  
Tommaso Di Noia

## Table of Contents

Unifying the Intentional and Institutional Semantics of Speech Acts . . . . .	1
<i>Carole Adam, Andreas Herzig, Dominique Longin, Vincent Louis</i>	
Tableaux for Acceptance Logic . . . . .	17
<i>Mathijs de Boer, Andreas Herzig, Tiago de Lima, Emiliano Lorini</i>	
Ontology and Time Evolution of Obligations and Prohibitions using Semantic Web Technology . . . . .	33
<i>Nicoletta Fornara, Marco Colombetti</i>	
Prioritized Goals and Subgoals in a Logical Account of Goal Change – A Preliminary Report . . . . .	49
<i>Shakil M. Khan, Yves Lespérance</i>	
Microscopic Analysis of Edge Creation Process in Trust-based Social Networks . . . . .	66
<i>Babak Khosravifar, Maziar Gomrokchi, Jamal Bentahar</i>	
Computing Utility from Weighted Description Logic Preference Formulas . . . . .	82
<i>Azzurra Ragone, Tommaso Di Noia, Francesco M. Donini, Eugenio Di Sciascio, Michael P. Wellman</i>	
Explaining and Predicting the Behavior of BDI-Based Agents in Role-Playing Games . . . . .	98
<i>Michal P. Sindlar, Mehdi M. Dastani, Frank Dignum, John-Jules Meyer</i>	
Correctness Properties for Multiagent Systems . . . . .	114
<i>Munindar P. Singh, Amit K. Chopra</i>	
Reasoning and Planning with Cooperative Actions for Multiagents Using Answer Set Programming . . . . .	130
<i>Tran Cao Son, Chiaki Sakama</i>	
Social Commitments in Time: Satisfied or Compensated (preliminary report) . . . . .	147
<i>Paolo Torroni, Federico Chesani, Paola Mello, Marco Montali</i>	
Verifying Dribble Agents . . . . .	162
<i>Doan Thu Trang, Brian Logan, Natasha Alechina</i>	
<b>Author Index . . . . .</b>	<b>178</b>





# Unifying the intentional and institutional semantics of speech acts

Carole Adam<sup>1</sup>, Andreas Herzig<sup>2</sup>, Dominique Longin<sup>2</sup>, and Vincent Louis

<sup>1</sup> RMIT University, Melbourne, VIC, Australia  
`carole.adam.rmit@gmail.com`

<sup>2</sup> IRIT, Toulouse, France  
`firstname.surname@irit.fr`

**Abstract.** Research about the semantics of agent communication languages traditionally sees the opposition between the mentalist and social approaches. In this paper we adopt a mixed approach since we propose a logical framework allowing us to express both the intentional and institutional dimensions of a communicative action. We use this framework to give a semantics for some speech acts representing each of Searle's categories except expressives. This semantics relaxes the criticized constraints imposed in FIPA-ACL and also extends this standard with new speech acts and new institutional features to characterise them. It has been implemented in an extension of the Semantic Add-on for the JADE agent development platform, and used in an industrial application in the context of automated B2B exchanges.

## 1 Introduction

Designing efficient Agent Communication Languages is an essential issue in Multi-Agent Systems in order to standardise exchanges between the agents. Research about the semantics of agent communication languages sees the subscribers of the social approach [17, 8, 27] criticize mentalist approaches [24, 15] for only grounding on the agents' private mental attitudes. But one can similarly reproach to social approaches to provide a semantics only based on the agents' public commitments, independently of their mental attitudes. Now these "social attitudes" are mainly descriptive, while mental attitudes allow one to predict the agents' behaviour. Moreover mental attitudes allow agents to reason about social notions. It is thus essential to consider both mental and social attitudes. Some researchers thus propose a mixed approach based both on public and private aspects [18]. But they do not formalise institutional speech acts like declarations. Now such speech acts are essential in new application fields involving communication about norms, roles or powers of agents, for instance in electronic commerce or automated business to business exchanges.

In this paper we thus want to propose an alternative to the well-known standard of agent communication language FIPA-ACL [16] through the following changes: relaxed feasibility preconditions to allow a more flexible utilisation of

the speech acts in various contexts; new institutional speech acts like declarations and promises; and an institutional interpretation of speech acts coupled with their classical intentional interpretation. Therefore we adapt an existing logical framework for the formalisation of institutional notions like roles, powers and norms [10]. We then formalise in this logical framework the institutional interpretation of some specific communicative actions, each one representing one of Searle’s categories of speech acts (except expressive ones) [25]. Our notion of institution is very large (it is a set of rules and facts adopted by a group of agents, like the rules of a game, or the laws of a country) and covers formal, legal institutions as well as informal ones (social rules in a group...).

The paper is structured as follow. Section 2 discusses some other social semantics of speech acts. Section 3 briefly describes the syntax, semantics and axiomatics of our logical framework. The core of the paper (Section 4) is dedicated to the unified semantics of speech acts. We are then able to compare our semantics of ACL with some related ones in more details (Section 5). Finally we conclude about the future prospects opened by this work (Section 6).

## 2 State of the art

The mentalist approach consists in grounding the semantics of speech acts on the agents’ internal mental attitudes. These are represented by belief, desire and intention modalities provided by BDI logics, that are classically used to formalise the reasoning of autonomous agents [23, 28]. This resulted in the design of several standards of agent communication languages like KQML [15] or FIPA [16], this latter one grounding on Sadek’s rational interaction theory [24].

These approaches were criticised a lot for being only based on private concepts (mental attitudes) instead of public verifiable notions (like commitments). Therefore some work exist aiming at enriching BDI logics with deontic operators like obligation [13, 3] or with institutional operators like count as or institutional power [21], in order to formalise the institutional interpretation of speech acts exchanged by the agents. In previous work we used such an extended BDI framework to express the semantics of speech acts with institutional effects [11] but we were limited to declarative speech acts, and the intentional and institutional dimensions were quite blended.

Various other work aims at providing an institutional semantics for speech acts. For example Dignum and Weigand [14] propose a logical framework combining illocutionary and deontic logic to study and model the norms resulting from communication between agents; however, they only consider directive speech acts. Boella *et al.* [1] propose a role-based semantics allowing them to combine social commitments and mental attitudes to express the semantics of speech acts in the context of persuasion dialogues. Actually they rewrite the FIPA feasibility precondition and rational effects of speech acts but replace the private mental attitudes involved by public mental attitudes attributed to the agents’ roles instead of the individual agents. This solves the flaw of mentalist approaches, criticised for grounding on unverifiable mental attitudes, but finally

there is no distinct institutional interpretation of speech acts, that could differ from one institution to another. In the following subsections we give some details about two approaches: Fornara and Colombetti’s approach based on the notion of commitments, and Lorini *et al.*’s approach based on the notion of acceptance.

## 2.1 Fornara and Colombetti: semantics in terms of social commitments

As opposed to the mentalist approach, the social one [26, 27, 8] assumes that private mental attitudes are not verifiable and thus grounds on the concept of public (thus verifiable) commitments [7] to express the semantics of speech acts. All the commitments taken by the agents are stored for possible future reference. The semantics of speech acts is expressed only in terms of such commitments.

For example Fornara and Colombetti [17] ground on Castelfranchi’s notion of commitment [7] to define a library of communicative acts. From the classification of speech acts into four categories (assertives, directives, commissives and declaratives) inspired from Searle’s work [25], they redefine for each category the semantics of its speech acts in terms of social commitments. Thanks to this library, they provide a communication tool based on social commitments, alternative to the FIPA-ACL standard. This tool allows rational agents to reason about the underlying rules of communication and to respect them in order for the system to behave well.

However they are limited to the institutional dimension of speech acts and neglect their relations with the agents’ mental attitudes. Yet agents must be able to reason autonomously about the institution before making their decision to perform a given speech act. Moreover no specific institution is explicit in their commitments, making it impossible to have various commitments in various institutions; therefore it is also impossible for speech acts to have different effects depending on the institution within which they are interpreted. For example the action of nodding one’s head is interpreted in the context of French gestural language as meaning “yes”, while in the context of Bulgarian gestural language it is interpreted as meaning “no”.

## 2.2 Lorini *et al.* : semantics in terms of group acceptance

Lorini *et al.* [22] define a new semantics for speech acts using Gaudou *et al.*’s Acceptance Logic [19].  $\mathcal{AL}$  is a modal logic extended with the notion of acceptance, representing what a group of agents willingly accept to consider as true (even if some (or all) members of the group believe the opposite) in a given institutional context (and that they can refuse in another context). Acceptances influence the agents’ behaviour and utterances in this institutional context. They are represented with the operator  $[C : x]\varphi$  reading “agents in group C accept that  $\varphi$  while functioning as members of this group in the institutional context  $x$ ”.

Institutional notions are not primitive but defined from this notion of acceptance. Thereby institutional facts (facts that are only valid in an institutional context, but not objectively valid) are considered to be equivalent to a group

acceptance in all groups of agents while they function in the considered institutional context. This strong link assumed between acceptance and institutional facts may be a particular rule of the specific “ordinary communication” institution but cannot be generalised to all institutions, particularly legal ones.

The authors then consider the speech act Promise in the institutional context of Ordinary Communication (OC). According to them, if  $i$  informs  $j$  that he is going to perform action  $\alpha$  for him, and  $j$  intends  $i$  to perform this action for him, this counts as a promise at the next instant. The consequence of this promise is that  $i$  is obliged to perform action  $\alpha$  for  $j$ . Moreover the acceptance by these two agents  $i$  and  $j$  while functioning as a group in institution  $OC$  that  $i$  has promised to perform action  $\alpha$  for  $j$  and that  $j$  intends him to do so implies a social commitment of  $i$  towards  $j$  to perform  $\alpha$  for him. This framework is interesting but Lorini *et al.* have only formalised the promise yet. Moreover they do not seem to make a clear distinction between the intentional and institutional preconditions to perform a speech act.

### 3 Our logical framework

We adapt here an existing logical framework for norms, institutional powers and roles defined in [10]. It is a multi-modal logic with modal operators of belief, intention, obligation, institutional facts and consequences, and action.

#### 3.1 Syntax

Let  $AGT = \{i, j, \dots\}$  be a finite set of agents. Let  $ACT = \{\alpha, \beta, \dots\}$  be the set of actions. We suppose that some actions in  $ACT$  are of the form  $i:\alpha$ , where  $i$  is the author of action  $\alpha$  (the agent who performs it). Let  $ATM = \{p, q, \dots\}$  be the set of atomic formulas. Let  $INST = \{s, t, \dots\}$  be the set of institutions. Complex formulas are denoted by  $\varphi, \psi, \dots$ . The language of our logic is defined by the following BNF grammar:

$$\varphi ::= p | \neg\varphi | \varphi \vee \varphi | B_i\varphi | Ch_i\varphi | I_i\varphi | D_s\varphi | \varphi \Rightarrow_s \varphi | O\varphi | before_\alpha\varphi | after_\alpha\varphi$$

where  $p$  ranges over  $ATM$ ,  $\alpha$  over  $ACT$ ,  $i$  over  $AGT$ , and  $s$  over  $INST$ . The classical boolean connectives  $\wedge, \rightarrow, \leftrightarrow, \top$  (tautology) and  $\perp$  (contradiction) are defined from  $\vee$  and  $\neg$  in the usual manner. The operators  $done_\alpha\varphi, happens_\alpha\varphi, P\varphi, F\varphi$  and  $power(i, s, \varphi, \alpha, \psi)$  will be defined as abbreviations.

#### 3.2 Semantics and axiomatics

We only give here the informal meanings of our operators. It is sufficient to know that they have a Kripke semantics in terms of possible worlds. We also give some useful axioms. This framework is adapted from Demolombe and Louis’ logic of norms, roles and institutional powers [10]. But please notice that actually, the details of the semantics of operators is not important, and any other institutional logic would work.

**Belief, intention and action**  $B_i p$  means that agent  $i$  believes that  $p$ .  $Ch_i p$  means that agent  $i$  prefers  $p$  to be true. These two normal operators have a standard  $KD45$  axiomatics.  $I_i p$  means that agent  $i$  intends that  $p$ . Its axiomatics is that defined for FIPA by Sadek [24]. In particular intention is linked with belief by the following mix axioms:

- introspection:  $I_i p \leftrightarrow B_i I_i p$
- automatic dropping of achieved intentions:  $I_i p \rightarrow \neg B_i p$

$before_\alpha$  and  $after_\alpha$  are normal modal operators defined in standard tense logic in linear time version [6].  $done_\alpha \varphi = \neg before_\alpha \neg \varphi$  means that action  $\alpha$  has just been performed, and  $\varphi$  was true before.  $happens_\alpha \varphi = \neg after_\alpha \neg \varphi$  means that action  $\alpha$  is about to be performed and  $\varphi$  will be true just after.

**Institutional modalities** Finally this framework also provides some specific operators to formalise institutional concepts. These operators have a parameter  $s$  specifying the institution within which they are valid. Here we consider an institution as a set of institutional facts and rules that a group of agents (the “members” of this institution) adopt. This is a general view that can account for various institutional contexts, be they formal institutions or informal ones: the law of a country, a contract between two parties in a business relationship, a social structure, the rules of a game...

An **institutional fact** is a fact that is recognised to be valid in the context of a given institution, but that can make no sense in itself; *i.e.* it is not a physically observable fact (what Searle calls a “brute fact”) but something written in the registry of this institution. For example the fact that two people are married, or that one is authorised to drive a truck, is only valid *w.r.t.* the law of a country; all deontic fact should also be encapsulated in an institutional fact to make the institution in which they hold explicit. We represent these institutional facts with the operator  $D_s \varphi$  meaning that for institution  $s$ , it is officially established that  $\varphi$  holds. In particular if  $\varphi$  is an agent’s mental attitude, then  $D_s \varphi$  can be understood as this agent’s commitment (either a propositional commitment if  $\varphi$  is a belief, or a commitment in action if  $\varphi$  is an intention).

Institutional facts can be deduced from other facts thanks to the rules of the institution. For example the presentation of an invoice by a provider to his client counts as an obligation for the client to pay it. The existence of the invoice is physically observable, while the obligation is only valid in an institutional context. We represent these normative consequences with the primitive operator  $p \Rightarrow_s q$ , meaning that according to the norms holding in institution  $s$ ,  $p$  entails  $q$ . This operator is known in the literature as *count as*, and has been first formalised by Sergot and Jones [21]. The following mix axioms explicit the link between institutional facts and normative consequences:

$$\begin{aligned} (\varphi \Rightarrow_s \psi) &\rightarrow D_s(\varphi \rightarrow \psi) && \text{(SD)} \\ (\varphi \Rightarrow_s \psi) &\rightarrow (\varphi \rightarrow D_s \varphi) && \text{(SC)} \end{aligned}$$

From these axioms and the properties of  $D_s$  (see [10, p.8] for details) we can deduce:

$$(\varphi \Rightarrow_s \psi) \rightarrow (\varphi \rightarrow D_s \psi) \quad (\text{SP})$$

A particular case of normative consequence concerns the consequences of the performance of an official procedure. Actually some agents can have the power when performing a given procedure under some conditions to create new institutional facts. We represent these institutional powers as an abbreviation  $\text{power}(i, s, \text{cond}, \alpha, \varphi) = ((\text{done}_{i:\alpha} \top \wedge \text{cond}) \Rightarrow_s \varphi)$ . Intuitively this means that  $i$  has the power in institution  $s$ , by performing action  $\alpha$  and if condition  $\text{cond}$  holds, to see to it that  $\varphi$  becomes officially true in institution  $s$ . For example a mayor has the power in the law of the French Republic, by performing a declaration, and on condition that the two people agree, to marry them. Obviously these powers result from the agent’s role in the institution, but this is not the focus of this paper so we will not remind how roles are formalised in the original framework (the interested reader can refer to [10] for details on this point).

**Deontic modalities** We have a modality for impersonal obligation to be:  $O\varphi$  reads “it is obligatory that  $\varphi$ ”, and its axiomatic is that of the Standard Deontic Logic [20], that is  $KD$ . *Obligations to do* can be expressed as obligations to be in a state where the obliged action has been performed. Obligations are impersonal since no agent is explicitly responsible for their fulfilment, but such an agent can implicitly appear in their content. For instance  $O\text{done}_{i:\alpha} \top$  means that it is obligatory (for no one in particular) to be in a state where  $i$  has just performed action  $\alpha$ ; this can be understood as “ $i$  has the obligation to perform action  $\alpha$ ”.

Permissions and interdiction are defined from obligations in a standard way:  $P\varphi = \neg O\neg\varphi$  means that it is permitted that  $\varphi$ , and  $F\varphi = O\neg\varphi$  means that it is forbidden that  $\varphi$ .

Please notice that no institution is explicit as a parameter of this obligation modality. But such obligations will be encapsulated in institutional facts to express the institution in which they are valid. For example  $D_s O\varphi$  means that “in institution  $s$ , it is obligatory that  $\varphi$ ”.

## 4 Semantics of speech acts

### 4.1 Preliminary remarks

**Intentional and institutional dimensions** The FIPA-ACL standard [16] defines features allowing one to give an **intentional dimension** to the observation and interpretation of a communicative action: the feasibility precondition (the appropriate mental attitudes to perform the speech act) and the rational effect (this is a formula  $\varphi$  representing the content of the speaker  $i$ ’s intention that he intends the receiver  $j$  to know; so the performance of the speech act allows any observer  $k$  to deduce this corresponding intentional effect:  $B_k I_i B_j I_i \varphi$ ). Please

notice that the performance of the speech act does **not** automatically allow one to deduce its rational effect, but only its intentional effect, meaning that any agent  $k$  believes that the speaker  $i$  intends the hearer  $j$  to recognize its ( $i$ 's) intention to achieve the rational effect  $\varphi$ . However, nothing ensures that  $i$  indeed achieves  $\varphi$ , his speech act may fail, for example the hearer may not obey an order, or may not believe an assertion. Thus the rational effect can only be deduced under some constraining hypotheses such as the sincerity and competence hypotheses used in FIPA.

In a similar way, we want to provide here the **institutional dimension** of the observation and interpretation of a communicative action relative to one or several institutions. This institutional interpretation is composed of the following features:

- a permission condition that is necessary and sufficient for the speaker to be allowed to perform this speech act;
- a power condition that also needs to be true for the speech act to have an institutional effect;
- an explicit institutional effect that is obtained when the speech act is performed while permission and power conditions were true.

We will thus be able to combine the intentional and institutional dimensions of communicative actions (formalised as speech acts [25]), both essential to fully characterise their interpretation. Lorini *et al.* have also investigated such a unified approach but they have only formalised the interpretation of a promise in the context of ordinary communication; we aim at being much more generic. In particular we formalise one speech act from each of Searle's categories of illocutionary forces, except the expressive one.

Actually we have relaxed some of the (widely criticised) strong constraints imposed by FIPA-ACL semantics on the appropriate context of performance of speech acts. Instead of imposing these conditions as strong constraints, we have moved them into the permission preconditions of the speech act. The agents are thus physically able to disobey these constraints, but it is forbidden by the interaction norms, and they may incur sanctions for such violations.

**Notations** In the sequel we use the following abbreviations:

- FP = feasibility preconditions
- RE = rational effect
- PermC = (institutional) permission condition
- PowC = power condition
- EE = institutional explicit effect

Speech acts are actions of the form  $Force(sp, ad, inst, content)$  where  $sp \in AGT$  is the speaker,  $ad \in AGT$  is the addressee,  $inst \in INST$  is the institutional context,  $content$  is the propositional content and can be any formula of our language, and  $Force \in \{inform, promise, command, declare\}$  is the illocutionary force.

**Action laws** We now explain how the intentional and institutional dimensions of actions interact by providing the action laws governing the performance of speech acts.

We notice that  $FP$  is a factual executability precondition, while  $PermC$  is an ideal one. But even ideal worlds are submitted to physical world laws, *i.e.*  $PermC$  is not sufficient for the action to be executable,  $FP$  also has to be true. For example a mayor has the permission to marry people by making a declaration, but the declaration must be executable; thus if he is voiceless one day, he will be unable to marry anyone.

We thus have the following executability laws. The factual executability law ( $FEL_\alpha$ ) means that an action happens iff its feasibility precondition is true and the agent chooses to perform it. The ideal executability law ( $IEL_\alpha$ ) means that ideally, an action should happen iff it is permitted and feasible.

$$\begin{aligned} happens_\alpha \top &\leftrightarrow (FP(\alpha) \wedge Ch_i happens_{i:\alpha} \top) && (FEL_\alpha) \\ O(happens_\alpha \top &\rightarrow PermC(\alpha)) && (IEL_\alpha) \end{aligned}$$

We also have the following effect laws. The rational effect law ( $REL_\alpha$ ) means that if the power precondition of an action is false, then only its rational effect can be deduced after its performance. The power effect law ( $PEL_\alpha$ ) means that if the power condition of an action is true, then both its rational and institutional effects can be deduced after its performance.

$$\begin{aligned} \neg PowC(\alpha) &\rightarrow after_\alpha RE(\alpha) && (REL_\alpha) \\ PowC(\alpha) &\rightarrow after_\alpha RE(\alpha) \wedge EE(\alpha) && (PEL_\alpha) \end{aligned}$$

From these laws we can deduce the following theorems clarifying the factual executability and effects of  $\alpha$  depending on the different combinations of its feasibility and power preconditions. If  $FP(\alpha)$  is false then  $\alpha$  is not executable.

$$\neg FP(\alpha) \rightarrow after_\alpha \perp$$

If  $FP(\alpha)$  is true but  $PowC(\alpha)$  is false,  $\alpha$  is about to happen after which its rational effect will be true.

$$(FP(\alpha) \wedge \neg PowC(\alpha)) \rightarrow (happens_\alpha \top \wedge after_\alpha RE(\alpha))$$

Finally if both  $FP(\alpha)$  and  $PowC(\alpha)$  are true,  $\alpha$  is about to happen after which both its rational and institutional effects will be true.

$$(FP(\alpha) \wedge PowC(\alpha)) \rightarrow (happens_\alpha \top \wedge after_\alpha (RE(\alpha) \wedge EE(\alpha)))$$

## 4.2 Assertives: inform

The assertive speech act *Inform* commits the speaker to the truth of a proposition. The notation  $inform(i, j, s, \varphi)$  reads “agent  $i$  informs  $j$  in institution  $s$  that  $\varphi$  is true”.



**Intentional interpretation** As we said before we have relaxed FIPA constraints on the executability of speech acts. We thus impose no feasibility precondition here.

$$FP(\text{inform}(i, j, s, \varphi)) = \top$$

The rational effect (the content of the speaker  $i$ 's intention that he intends the receiver  $j$  to know) is that  $j$  believes the promised proposition  $\varphi$  to be true:

$$RE(\text{inform}(i, j, s, \varphi)) = B_j \varphi$$

**Institutional interpretation** The permission precondition to inform  $j$  that  $\varphi$  in institution  $s$  includes the constraints removed from the factual feasibility preconditions: the speaker should not believe that the hearer already knows if  $\varphi$ , and he should not be already committed on  $\neg\varphi$  in the same institution.

$$PermC(\text{inform}(i, j, s, \varphi)) = \neg D_s B_i B_i f_j \varphi \wedge \neg D_s B_i \neg \varphi$$

Now the institutional effect of *Inform* is to retract possible opposite commitments contracted before and to assert a new commitment on  $\varphi$ . Indeed, even if agent  $i$  was previously committed on  $\neg\varphi$  (and therefore was not permitted to inform anyone that  $\varphi$ ), he may violate that obligation. But these two commitments are inconsistent so the previous one must be retracted while asserting the new contradictory one. Though one can still detect that the opposite commitment was true when  $i$  performed the action and that he has thus violated the rules of the institution. Actually due to the seriality of  $D_s$  we have that  $D_s B_i \varphi \rightarrow \neg D_s B_i \neg \varphi$ . So the explicit institutional effect of inform is the new commitment:

$$EE(\text{inform}(i, j, s, \varphi)) = D_s B_i \varphi$$

This effect is always obtained and does not depend on particular powers of  $i$ , so the power condition is trivial.

$$PowC(\text{inform}(i, j, s, \varphi)) = \top$$

**Example** For example in the context of B2B exchanges, if a provider sends his catalogue to a client, this counts as information about the prices given in this catalogue. As an effect of this action, the provider is thus committed to these prices during the validity of his catalogue. In the specific institution  $s$  constituted by the contract between the provider and the client, we assume that we have a specific rule forbidding to contradict one's commitments, which takes the form  $D_s O(D_s B_i p \rightarrow \text{after}_\alpha D_s B_i p)$ , for every speech act  $a$ , where  $p$  is the proposition denoting that the price is 100. This means that in institution  $s$ , it is obligatory that if an agent  $i$  is committed to believe that the price of an item is 100, then after any speech act he is still committed to this (in other words it is forbidden to retract this commitment by any speech act). From this we can deduce that the provider is obliged to respect the given prices, *i.e.*  $D_s O(D_s B_i p \rightarrow \text{after}_{\text{Inform}(i, j, s, \neg p)} \perp)$  (it is obligatory that if  $i$  is committed to  $p$ , then the action of informing any agent  $j$  that  $\neg p$  is not feasible).

### 4.3 Commissives: promise to

This commissive speech act commits the speaker on a course of action. The notation  $promise\text{-}to(i, j, s, \alpha)$  reads “ $i$  promises to  $j$  in institution  $s$  to perform action  $\alpha$ ”.

**Intentional interpretation** We begin with specifying the intentional dimension of this speech act, that is not given in FIPA-ACL. A promise-to is feasible if the speaker believes that the hearer intends him to perform the concerned action<sup>3</sup>. For example it makes no sense that a child promises to his father to play, while it makes sense to promise him to make his schoolwork. So:

$$FP(promise\text{-}to(i, j, s, \alpha)) = B_i I_j done_\alpha \top$$

The rational effect pursued by the speaker is that the hearer be aware of his intention to perform the promised action:

$$RE(promise\text{-}to(i, j, s, \alpha)) = B_j I_i done_{i:\alpha} \top$$

**Institutional interpretation** In an institutional context  $s$ , this promise to perform an action  $\alpha$  is permitted on condition that the action  $i:\alpha$  is not explicitly forbidden itself. So the permission precondition is the following:

$$PermC(promise\text{-}to(i, j, s, \alpha)) = \neg D_s O \neg happens_{i:\alpha} \top$$

The institutional effect consists in ratifying in institution  $s$  the speaker’s intention to perform action  $\alpha$ ; so after  $promise\text{-}to(i, j, s, \alpha)$  the speaker has stored in the registry of  $s$  its intention to perform  $\alpha$ , which is similar to him being committed in  $s$  to this course of action.

$$EE(promise\text{-}to(i, j, s, \alpha)) = D_s I_i done_{i:\alpha} \top$$

This is thus similar to the  $inform(i, j, s, \varphi)$  speech act except that a promise stores a commitment in action while an inform stores a propositional commitment.

There is no power precondition, so the institutional effect of a promise is always reached.

$$PowC(promise\text{-}to(i, j, s, \alpha)) = \top$$

**Example** A client  $c$  promises to pay his provider  $p$  once the ordered goods have been delivered. The action to pay is denoted by  $\alpha_{pay}$ . This promise is valid in the context of a B2B exchange contract, that is a particular institution denoted  $b2b$  here. So this promise is formalised as:  $promise(c, p, b2b, \alpha_{pay})$ . This promise is permitted since obviously the promised action to pay is not forbidden:  $\neg D_{b2b} O \neg happens_{\alpha_{pay}} \top$ . So when the client receives the delivery, his promise allows to deduce his commitment (ratified intention) to pay:  $D_{b2b} I_c done_{c:\alpha_{pay}} \top$ , that is the institutional effect of this speech act.

<sup>3</sup> Please notice that threats such as “I promise that I will kill you” cannot be considered as promises in the sense of Searle.

#### 4.4 Directives: command

This directive speech act is an attempt from the speaker to make the hearer perform some action. The notation  $command(i, j, s, \alpha, cond)$  reads “ $i$  orders to  $j$  in institution  $s$ , in virtue of condition  $cond$ , to perform action  $\alpha$ ”.

**Intentional interpretation** According to the FIPA-ACL semantics, a request is feasible only if the speaker does not believe the hearer to already intend to perform the commanded action, and does believe that the part of the feasibility preconditions of the commanded action that concerns him (*i.e.* that are his mental attitudes) are valid. Here we consider that when  $\alpha$  is an action of agent  $j$  then  $FP(\alpha)$  is of the form  $FP_i(\alpha) \wedge FP_{\neq i}(\alpha)$  where the former is “ $i$ ’s part of  $FP(\alpha)$ ” (similar to FIPA-ACL notation  $FP(\alpha)[i \setminus j]$ , that is the part of  $FP(\alpha)$  that are mental attitudes of agent  $i$ ). But we do not impose this constraint on the feasibility of  $\alpha$  as a feasibility precondition of the command, but rather as a permission precondition. So:

$$FP(command(i, j, s, \alpha, cond)) = \top$$

The rational effect of a command (*i.e.* the effect that  $i$  intends  $j$  to believe that  $i$  intends to achieve) is that  $j$  has performed the commanded action:

$$RE(command(i, j, s, \alpha, cond)) = done_{j:\alpha} \top$$

**Institutional interpretation** The permission precondition to command someone to perform an action is to be empowered to do so, *i.e.* to dispose of the institutional power to create the obligation to perform the commanded action by commanding it, and to validate the condition of this power, that is an explicit attribute  $cond$  of this directive speech act. An additional permission precondition is the constraint coming from FIPA feasibility precondition that we relaxed, that is that the part of the feasibility preconditions of  $\alpha$  that depends on  $i$  hold (one should not command someone to perform an action whose preconditions are made false by his own mental attitudes).

$$PermC(command(i, j, s, \alpha, cond)) = \\ cond \wedge power(i, s, cond, command(i, j, s, \alpha, cond), Odone_{j:\alpha} \top) \wedge FP_i(\alpha)$$

The explicit institutional effect of this power is to create two new institutional facts, corresponding to the obligation for  $j$  to perform  $\alpha$ , and the recording of  $j$ ’s knowledge of his obligation. Actually this obligation could exist before, and in this case the command corresponds to a notification; but it can also be created from scratch by the command (see the examples in the next paragraph for clarification).

$$EE(command(i, j, s, \alpha, cond)) = D_s Odone_{j:\alpha} \top \wedge D_s B_j Odone_{j:\alpha} \top$$

**Example** For example a parent can command his children to clean his room. In this case, the action becomes obligatory through the command, because of the parent’s authority over his son (that is the condition of his power). Similarly a professor commanding his students to make some schoolwork creates the obligation for them to do so, on the strength of his role of professor.

But an order does not necessarily *create* an obligation, and may just put in focus an existing one. For example a bailiff can be sent to officially command an uncooperative client to pay an invoice. In this case the obligation already exists (and is attested by the invoice) so the bailiff only reminds the client of it<sup>4</sup>. He is permitted to perform such a command in virtue of his role of bailiff and because he is sent by the provider.

#### 4.5 Declaratives: declare

This declarative speech act changes the institutional reality by creating a new institutional fact. The notation  $declare(i, j, s, cond, \varphi)$  reads “*i* declares to *j* in institution *s* that given condition *cond*, the fact  $\varphi$  is now established”. The condition usually bears upon the speaker’s role that empowers him to perform such a declaration.

**Intentional interpretation** This intentional interpretation is partly inspired from the intentional interpretation of an  $inform(i, j, s, D_s \varphi)$ . The feasibility precondition of a declaration is that the speaker does not believe the declared fact to be already established (indeed a declaration must create a **new** institutional fact). The rational effect (*i.e.* the intended effect) is the same as an inform about  $D_s \varphi$ , *i.e.* to make the hearer aware of this information. So:

$$FP(declare(i, j, s, cond, \varphi)) = \neg B_i D_s \varphi$$

$$RE(declare(i, j, s, cond, \varphi)) = B_j D_s \varphi$$

**Institutional interpretation** The permission precondition to perform  $declare(i, j, s, cond, \varphi)$  is that *i* really has the power to establish the declared fact  $\varphi$  by declaring it under the announced conditions *cond*. This power is locally granted by each specific institution to some agents depending on their role. For example the French republic grants the mayors the right to pronounce two people husband and wife, under the condition that they both consent to it. Thus an ordinary agent who is not mayor does not have this power, so that he is not allowed to pronounce marriages.

$$PermC(declare(i, j, s, cond, \varphi)) =$$

$$power(i, s, cond, declare(i, j, s, cond, \varphi), \varphi)$$

---

<sup>4</sup> Actually this seems to be a notification rather than a command, but the aim is to make the client behave, while the aim of a notification is only to make the receiver officially aware of what is notified. In further work we expect to study into more details the links between declarations, commands and notifications.

The explicit effect of a declaration is to store the declared fact in the institution, as well as the fact that the hearer is officially aware of this fact.

$$EE(\text{declare}(i, j, s, \text{cond}, \varphi)) = D_s \varphi \wedge D_s B_j D_s \varphi$$

This explicit effect is only obtained under the additional condition that *cond* is valid:

$$PowC(\text{declare}(i, j, s, \text{cond}, \varphi)) = \text{cond}$$

**Example** For example a country can declare war to another one, by the voice of its representative that is empowered to do so, and under some conditions like the agreement of some counsellors. A mayor is empowered by its country to pronounce weddings under some conditions that the people are of age and consenting.

Citizens have to declare their income to the public treasury in order to calculate the amount of tax that they will pay. This is a declaration since the effect is a new institutional fact officially establishing one's declared income as being believed by him. Any citizen is empowered to do so. Moreover the law imposes a constraint on the generated commitment, that is an obligation to believe this income to be true. Thereby if the declared income was false the citizen is liable for prosecution and sanctions.

#### 4.6 Example of reasoning with our action laws

This example is situated in the context of a B2B exchange (in institution *b2b*) between a buyer *b* and a seller *s*. The seller intends potential clients to know the prices of his products, *e.g.*  $I_s B_b p$ . With our relaxed feasibility precondition, he can use an assertive speech act whatever the context. Though if the buyer has already been informed of the prices before ( $D_{b2b} B_s B_b p$ ), the seller is not permitted to inform him again. Thus if he informs him anyway, according to  $IEL_\alpha$  he violates an obligation. This can be detected by other agents, and specific rules of the institution may specify sanctions to compensate this. Being aware of such pre-specified sanctions, an agent can deliberately choose to violate an obligation if the intended outcome (here, that clients be aware of the seller's offer) is more important than the incurred sanction. This shows the importance of having both intentional and institutional semantics of speech acts, to allow agents to reason about the relative importance of their goals and their obligations, in order to make an appropriate decision.

## 5 Detailed comparison with other work

In this section we compare our semantics of speech acts with those proposed by Fornara and Colombetti, and by Lorini *et al.* (that we have presented above).

## 5.1 Concept of commitment

We have shown before that what we mean by commitment in this work is a ratified mental attitude, *i.e.* a mental attitude (belief or intention) stored in the institution. This notion is similar to Fornara and Colombetti's commitment that is also a public concept, except that we have not made explicit its creditor. Actually the debtor is committed towards the whole institution, but an implicit creditor can sometimes be found in the content  $\varphi$  of the commitment. For example if agent  $i$  promises to  $j$  to pay him, he commits himself to a proposition involving agent  $j$ , expressing that  $j$  will be payed at some future instant. The creditor can sometimes be found in the sanction associated with the violation of the commitment, too; for example the obligation to pay damages to an agent.

Our notion of ratified mental attitude is also similar to Gaudou *et al.*'s notion of acceptance, because it must influence the agent's behaviour and utterances. Indeed, the agent's ratified mental attitudes are mental attitudes that he has expressed, that are stored in the institution, and to which he must conform while subsequently acting and speaking, even if they are not consistent with his real mental attitudes. For example an agent who promises that he has seen a given movie must then be able to talk about it in order to be consistent with his promise; if he is unable to narrate the end of the movie one can notice that he is contradicting his commitment.

## 5.2 Notion of institution

By institution we mean a set of rules and facts that are adopted by a group of agents (the members of the institution). This seems to be a more generic notion than Lorini *et al.*'s concept of informal institution, since it accounts for this particular kind of institutions but also for various other ones: laws of a country, rules of a game, contract between businesses, social norms of a culture... In particular it allows to have institutional rules that are ignored by the members of the institution, what is the case for law for example, since one cannot be aware of the whole set of laws of his country, while he is one of its citizens. Fornara and Colombetti do not make explicit the institutional context in which their speech acts are interpreted, so we believe that they also consider a kind of "ordinary communication" institutional context.

In our view informal institutions are described by a specific set of facts and rules, determining their specific functioning. In particular the fact that all agents must accept a fact for it to become institutional is a particular institutional law. In other kinds of institutions, facts must be adopted by a majority of members (voting to create a law or to elect the president for example), or the opinion of one single member can suffice (the referee is always right). Thus we cannot adopt such an hypothesis in our account. Indeed on the contrary we consider the generic interpretation of speech acts in any institution  $s$ . More specific rules can be additionally specified in each particular interpretation, but the object of this paper is to identify for each category of speech acts the features that are common to their institutional interpretations whatever the institutional context.

## 6 Conclusion

In this work we have provided an expressive logical framework blending the agents' mental attitudes (beliefs, intentions) with their social attitudes (obligations, institutional facts and powers...). To illustrate its expressivity, please notice that our framework allows to represent some forms of contrary-to-duty obligations-to-do. Such obligations take the form:

$$O_{after_\alpha} \perp \rightarrow after_\alpha O_{done_{repair_\alpha}} \top$$

where  $repair_\alpha$  is the contrary-to-duty obligation associated to the violation of the obligation to refrain from doing  $\alpha$ . This means that if it is forbidden to perform  $\alpha$ , then after  $\alpha$  it is obligatory to perform a repairing action  $repair_\alpha$ .

We have then used this framework to provide a semantics for an agent communication language based on FIPA-ACL but relaxing its widely criticised too constraining feasibility conditions, and adding permission preconditions. This way, agents can choose to perform forbidden speech acts but would then be liable to sanctions in the corresponding institution. Our ACL semantics also includes new speech acts (commissives and declaratives). It generalises existing approaches by unifying the intentional and institutional dimensions in one single framework, while strongly distinguishing them; moreover it allows to consider various kinds of institutional contexts; finally it provides action laws taking both dimensions into account.

In future work we intend to improve the institutional and intentional semantics of speech acts by accounting for deadlines. Various researchers [4, 12, 9] have shown that an important feature of obligations to perform an action is the deadline before which this action must be performed, that is essential to be able to assess the violation or fulfillment of such obligations. Though for the sake of simplicity we have omitted deadlines in this paper. An idea to manage them in future work could be to use existing formalisations of norms with deadlines, or to ground on linear temporal logic with until and since operators [5].

Finally we would like to mention that our framework for the institutional interpretation of speech acts has been successfully implemented into institutional agents that have been used in a prototype of industrial application: a multi-agent mediation platform for automated business to business exchanges [2].

## References

1. G. Boella, R. Damiano, J. Hulstijn, and L. van der Torre. Role-based semantics for agent communication: embedding of the 'mental attitudes' and 'social commitments' semantics. In *AAMAS'06*, Hakodate, Hokkaido, Japon, 2006.
2. F. Bourge, S. Picant, C. Adam, and V. Louis. A multi-agent mediation platform for automatic b2b exchanges. In *ESAW*, 2008. demonstration.
3. J. Broersen, M. Dastani, and L. van der Torre. Beliefs, obligations, intentions, and desires as components in an agent architecture. *International Journal of Intelligent Systems*, 20(9):893–919, 2005.

4. J. Broersen, F. Dignum, V. Dignum, and J.-J. C. Meyer. Designing a deontic logic of deadlines. In *DEON*, pages 43–56, 2004.
5. J. Burgess. Logic and time. *The Journal of Symbolic Logic*, 44(4):566–582, 1979.
6. J. Burgess. *Handbook of philosophical logic*, volume 7, chapter Basic tense logic, pages 1–42. Kluwer Academic Publishers, 2nd edition, 2002.
7. C. Castelfranchi. Commitments : From individual intentions to groups and organizations. In *ICMAS-95*, pages 41–48, San Francisco, 1995.
8. B. Chaib-draa, M.-A. Labrie, M. Bergeron, and P. Pasquier. Diagal : An agent communication language based on dialogues games and sustained by social commitments. *Journal of Autonomous Agent and Multi-Agent Systems*, 13:61–95, 2006.
9. R. Demolombe, P. Bretier, and V. Louis. Norms with deadlines in dynamic deontic logic. In *ECAI'06*, pages 751–752, 2006.
10. R. Demolombe and V. Louis. Norms, institutional power and roles: towards a logical framework. In *16th International Symposium on Methodologies for Intelligent Systems (ISMIS'06)*, volume LNAI 4203, pages 514–523. Springer, 2006.
11. R. Demolombe and V. Louis. Speech acts with institutional effects in agent societies. In *DEON'06*, 2006.
12. F. Dignum and R. Kuiper. Obligations and dense time for specifying deadlines. In *HICSS*, volume 5, pages 186–195. IEEE Computer Society, 1998.
13. F. Dignum, D. Morley, E. Sonenberg, and L. Cavedon. Towards socially sophisticated bdi agents. In *ICMAS'2000*, pages 111–118, 2000.
14. F. Dignum and H. Weigand. Communication and deontic logic. In R. Wieringa and R. Feenstra, editors, *Information Systems, correctness and reusability*, pages 242–260, Singapore, 1995. World Scientific.
15. T. Finin, R. Fritzson, D. McKay, and R. McEntire. KQML as an agent communication language. In *Int. conf. Information and knowledge management*, 1994.
16. FIPA. The foundation for intelligent physical agents. <http://www.fipa.org>.
17. N. Fornara and M. Colombetti. A commitment-based approach to agent communication. *Applied Artificial Intelligence*, 18(9-10):853–866, 2004.
18. B. Gaudou, A. Herzig, D. Longin, and M. Nickles. A new semantics for the fipa agent communication language based on social attitudes. In G. Brewka, S. Coradeschi, A. Perini, and P. Traverso, eds, *ECAI*, pp. 245–249. IOS Press, 2006.
19. B. Gaudou, D. Longin, E. Lorini, and L. Tummolini. Anchoring institutions in agents' attitudes: towards a logical framework for autonomous multi-agent systems. In *AAMAS'08*, 2008.
20. A. Jones and J. Carmo. *Handbook of philosophical logic*, volume 8, chapter Deontic Logic and Contrary-to-duties, pages 265–343. Kluwer Academic Publishers, 2002.
21. A. Jones and M. Sergot. A formal characterisation of institutionalised power. *Journal of the interest group in pure and applied logics*, 4(3), 1996.
22. E. Lorini, D. Longin, and B. Gaudou. The institutional dimension of speech acts: a logical approach based on the concept of acceptance. Research report, IRIT, 2008.
23. A. Rao and M. Georgeff. Modeling rational agents within a BDI-architecture. In *KR'91*, 1991.
24. D. Sadek. A study in the logic of intention. In *KR'92*, 1992.
25. J. R. Searle. Speech acts : An essay in the philosophy of language. In *Cambridge Univ. Press*, 1969.
26. M. Singh. An ontology for commitments in multiagent systems: Towards a unification of normative concepts. 1999.
27. M. P. Singh. A social semantics for agent communication languages. In F. Dignum and M. Greaves, eds., *Issues in Agent Communication*, pp. 31–45. Springer, 2000.
28. M. Wooldridge. *Reasoning about rational agents*. MIT Press, 2000.



# Tableaux for Acceptance Logic

Mathijs de Boer<sup>1</sup>, Andreas Herzig<sup>2</sup>, Tiago de Lima<sup>3</sup>, and Emiliano Lorini<sup>2</sup>

<sup>1</sup> University of Luxembourg, Luxembourg

<sup>2</sup> IRIT, University of Toulouse 3, France

<sup>3</sup> Eindhoven University of Technology, The Netherlands

**Abstract.** We continue the work initiated in [1–3], where the acceptance logic, a modal logic for modelling individual and collective acceptances was introduced. This logic is aimed at capturing the concept of acceptance *qua* member of an institution as the kind of attitude that agents are committed to when they are “functioning as members of an institution”. Acceptance logic can also be used to model judgement aggregation: it deals with how a collective acceptance of the members of an institution about a certain fact  $\varphi$  is created from the individual acceptances of the members of the institution. The contribution of this paper is to present a tableau method for the logic of acceptance. The method automatically decides whether a formula of the logic of acceptance is satisfiable thereby providing an automated reasoning procedure for judgement aggregation in the logic of acceptance.

**Key words:** Semantic tableaux method, acceptance logic, judgement aggregation, discursive dilemma

## 1 Introduction

The notion of ‘acceptance’ has been extensively studied in philosophy and social sciences where several authors have distinguished it from the classical notion of belief (see [4–6] for instance). Other authors have been interested in studying the foundational role of acceptance in the existence and in the dynamics of groups and institutions. It has been stressed in [7] (see also [8]) that the existence and the dynamics of an institution depend on the acceptances of the norms and the rules of the institution by the members of the institution. For example, for a certain norm to be a norm of institution  $x$ , all members of institution  $x$  must accept such norm to be valid. This relationship between acceptance and institutions was already emphasised in the philosophical doctrine of Legal Positivism [9]. According to Hart, the foundations of a normative system or institution consist of adherence to, or acceptance of, an ultimate rule of recognition by which the validity of any rule of the institution may be evaluated.<sup>4</sup>

In some recent works [1–3] we have presented a logical framework in which such relationship between acceptances and institutions can be formally studied.

---

<sup>4</sup> In Hart’s theory, the rule of recognition is the rule which specifies the ultimate criteria of validity in a legal system.

We conceive institutions as rule-governed social practices on the background of which the agents reason. For example, take the case of a game like Cluedo. The institutional context is the rule-governed social practice which the agents conform to in order to be competent players and on the background of which agents reason. In the context of Cluedo, an agent accepts that something has happened *qua* player of Cluedo (e.g., the agent accepts that Mrs. Red is the murderer *qua* player of Cluedo). Our logic is aimed at capturing the state of acceptance *qua* member of an institution as the kind of acceptance one is committed to when one is “functioning as a member of the institution” [7]. Moreover, it enables to formalise the concept of ‘collective acceptance’ of groups of agents. Following [10, 7], we conceive a collective acceptance held by a set of agents  $G$  *qua* members of a certain institution  $x$  as the kind of acceptance the agents in  $G$  are committed to when they are “functioning together as members of this institution”. For example, in the context of Greenpeace agents (collectively) accept that their mission is to protect the Earth *qua* members of Greenpeace. The state of acceptance *qua* members of Greenpeace is the kind of acceptance these agents are committed to when they are functioning together as members of Greenpeace. Thus, in our logical framework a collective acceptance by a set of agents  $G$  is based on the identification of the agents in  $G$  as members of a certain institution (or group, team, organisation, etc.) and on the fact that the agents in  $G$  recognise each other as members of the same institution (or group, team, organisation, etc.).

More recently [11], we have shown that our logic of acceptance can also be applied to modelling some interesting aspects of judgement aggregation. In the logic of acceptance the problem of judgement aggregation is a particular case of the problem of explaining how collective acceptance of the members of a certain group about a certain fact  $\varphi$  is created from the individual acceptances of the members of this group.

The contribution of this article is to present a tableau method for the logic of acceptance we introduced in [1, 2]. The method automatically decides whether a formula of the logic of acceptance is satisfiable thereby providing an automated reasoning procedure for making judgement aggregation in modal logic.

The remainder of the paper is organised as follows. First, in Section 2 we briefly present acceptance logic. Then, in Section 3 we present our tableau method. In Section 4 we apply it to a classical scenario in judgment aggregation, the so-called Discursive Dilemma [12, 13]. And finally, Section 5 concludes.

## 2 Acceptance logic

The logic AL (*Acceptance Logic*) was introduced in [1, 2]. It allows to express that some agents identify themselves as members of a certain institution and what (groups of) agents accept while functioning together as members of an institution. The principles of AL clarify the relationships between individual acceptance (acceptances of individual agents) and collective acceptance (acceptances of groups of agents).

## 2.1 Syntax

Assume a finite non-empty set  $X$  of labels denoting institutional contexts, a finite non-empty set  $N$  of labels denoting agents and a countable set  $P$  of atomic formulae. We use  $2^{N^*}$  to denote the set  $2^N \setminus \emptyset$ .

The language  $\mathcal{L}_{AL}$  of acceptance logic is the set of formulae  $\varphi$  defined by the following BNF:

$$\varphi ::= p \mid \neg\varphi \mid \varphi \vee \varphi \mid A_{G:x}\varphi$$

where  $p$  ranges over  $P$ ,  $G$  ranges over  $2^{N^*}$ , and  $x$  ranges over  $X$ .

The other classical Boolean connectives  $\wedge$ ,  $\rightarrow$ ,  $\leftrightarrow$ ,  $\top$  (tautology) and  $\perp$  (contradiction) are defined using  $\vee$  and  $\neg$  in the usual manner. And for simplicity, we write  $i:x$  instead of  $\{i\}:x$ .

The formula  $A_{G:x}\varphi$  reads ‘the agents in  $G$  accept that  $\varphi$  while functioning together as members of institution  $x$ ’. For example,  $A_{G:Greenpeace}protectEarth$  expresses that the agents in  $G$  accept that the mission of Greenpeace is to protect the Earth while functioning as activists in the context of Greenpeace; and  $A_{i:Catholic}PopeInfallible$  expresses that agent  $i$  accepts that the Pope is infallible while functioning as a Catholic in the context of the Catholic Church.

The same agent may accept contradictory propositions in two different contexts. For example, while functioning as a Catholic, agent  $i$  accepts that killing is forbidden, and while functioning as a soldier  $i$  accepts that killing is allowed.

The formula  $A_{G:x}\perp$  has to be read ‘agents in  $G$  are not functioning together as members of institution  $x$ ’. This means that we assume that functioning as a member of an institution is, at least in this minimal sense, a rational activity. Conversely,  $\neg A_{G:x}\perp$  has to be read ‘agents in  $G$  are functioning together as members of institution  $x$ ’. Thus,  $\neg A_{G:x}\perp \wedge A_{G:x}\varphi$  stands for ‘agents in  $G$  are functioning together as members of institution  $x$  and they accept that  $\varphi$  while functioning together as members of  $x$ ’ or simply ‘agents in  $G$  accept that  $\varphi$  *qua* members of institution  $x$ ’. This is a case of group acceptance. For the individual case, the formula  $\neg A_{i:x}\perp \wedge A_{i:x}\varphi$  has to be read ‘agent  $i$  accepts that  $\varphi$  *qua* member of institution  $x$ ’.

## 2.2 Semantics and axiomatisation

We use a standard possible worlds semantics. Let the set of all pairs of non-empty sets of agents and institutional contexts be  $\Delta = \{G:x \mid G \in 2^{N^*} \text{ and } x \in X\}$ . An acceptance model is a triple  $\langle W, \mathcal{A}, \mathcal{V} \rangle$  where:  $W$  is a non-empty set of possible worlds,  $\mathcal{A} : \Delta \rightarrow W \times W$  maps every  $G:x \in \Delta$  to a relation  $\mathcal{A}(G:x)$  between possible worlds in  $W$  and  $\mathcal{V} : P \rightarrow 2^W$  is valuation function associating a set of possible worlds  $\mathcal{V}(p) \subseteq W$  to each atomic formula  $p$  of  $P$ .

Instead of  $\mathcal{A}(G:x)$  we write  $\mathcal{A}_{G:x}$ , and we use  $\mathcal{A}_{G:x}(w)$  to denote the set  $\{w' \mid \langle w, w' \rangle \in \mathcal{A}_{G:x}\}$ .  $\mathcal{A}_{G:x}(w)$  is the set of worlds that is acceptable by the agents in  $G$  while functioning together as members of institution  $x$ .

Given  $M = \langle W, \mathcal{A}, \mathcal{V} \rangle$  and  $w \in W$ , the pair  $\langle M, w \rangle$  is a pointed acceptance model. The satisfaction relation  $\models$  between formulae of  $\mathcal{L}_{AL}$  and pointed acceptance models  $\langle M, w \rangle$  is defined as usual for atomic propositions, negation and

disjunction. The satisfaction relation for acceptance operators is the following:

$$M, w \models A_{G:x}\varphi \quad \text{iff} \quad M, w' \models \varphi \text{ for all } w' \in \mathcal{A}_{G:x}(w)$$

Validity of a formula  $\varphi$  (noted:  $\models \varphi$ ) is defined as usual.

The axiomatisation of AL is presented in Figure 1. As usual, the K-principles are the axioms and inference rules of the basic modal logic K.

$$\begin{array}{ll}
\text{(K)} & \text{All K-principles for the operators } A_{G:x} \\
\text{(4*)} & A_{G:x}\varphi \rightarrow A_{H:y}A_{G:x}\varphi \quad (\text{if } H \subseteq G) \\
\text{(5*)} & \neg A_{G:x}\varphi \rightarrow A_{H:y}\neg A_{G:x}\varphi \quad (\text{if } H \subseteq G) \\
\text{(Inc)} & (\neg A_{G:x}\perp \wedge A_{G:x}\varphi) \rightarrow A_{H:x}\varphi \quad (\text{if } H \subseteq G) \\
\text{(Una)} & A_{G:x}\left(\bigwedge_{i \in G} A_{i:x}\varphi \rightarrow \varphi\right)
\end{array}$$

**Fig. 1.** Axiomatisation of acceptance logic.

Axioms **4\*** and **5\*** are introspection axioms: when the agents in a set  $G$  function together as members of institution  $x$  then, for all  $y \in X$  and all  $H \subseteq G$ , the agents in  $H$  have access to all the facts that are accepted (or that are not accepted) by the agents in  $G$ . In particular, if the agents in  $G$  (do not) accept that  $\varphi$  while functioning together as members of institution  $x$  then, while functioning together as members of institution  $x$ , the agents of every subset  $H$  of  $G$  accept that agents in  $G$  (do not) accept that  $\varphi$ .

Axiom **Inc** says that, if the agents in  $G$  accept that  $\varphi$  *qua* members of institution  $x$  then every subset  $H$  of  $G$  accepts  $\varphi$  while functioning together as members of institution  $x$ . This means that what is accepted by the agents in  $G$  *qua* members of institution  $x$  are necessarily accepted by agents in all of  $G$ 's subsets with respect to the same institutional context  $x$ . Axiom **Inc** describes the *top down* process leading from  $G$ 's collective acceptance to the individual acceptances of  $G$ 's members.

Axiom **Una** expresses a unanimity principle according to which the agents in  $G$ , while functioning together as members of institution  $x$ , accept that if each of them individually accepts that  $\varphi$  while functioning as member of  $x$ , then  $\varphi$  is the case. This axiom describes the *bottom up* process leading from individual acceptances of the members of  $G$  to the collective acceptance of the group  $G$ .

In order to make our axioms valid we impose the following constraints on acceptance models, for any world  $w \in W$ , institutional context  $x \in X$ , and groups  $G, H \in 2^{N^*}$  such that  $H \subseteq G$ :

$$\begin{array}{ll}
\text{(C.4*)} & \text{if } w_2 \in \mathcal{A}_{H:y}(w_1) \text{ and } w_3 \in \mathcal{A}_{G:x}(w_2) \text{ then } w_3 \in \mathcal{A}_{G:x}(w_1); \\
\text{(C.5*)} & \text{if } w_2 \in \mathcal{A}_{H:y}(w_1) \text{ and } w_3 \in \mathcal{A}_{G:x}(w_1) \text{ then } w_3 \in \mathcal{A}_{G:x}(w_2); \\
\text{(C.Inc)} & \text{if } \mathcal{A}_{G:x}(w) \neq \emptyset \text{ then } \mathcal{A}_{H:x}(w) \subseteq \mathcal{A}_{G:x}(w); \\
\text{(C.Una)} & \text{if } w_2 \in \mathcal{A}_{G:x}(w_1) \text{ then } w_2 \in \bigcup_{i \in G} \mathcal{A}_{i:x}(w_2).
\end{array}$$

Axiom **4\*** corresponds to semantic constraint **C.4\***, Axiom **5\*** corresponds to **C.5\***, Axiom **Inc** to **C.Inc**, and **Una** to **C.Una** (in the sense of correspondence theory). We also note that **C.4\*** and **C.5\*** together are equivalent to the following semantic constraint: if  $w_2 \in \mathcal{A}_{H:y}(w_1)$  then  $\mathcal{A}_{G:x}(w_1) = \mathcal{A}_{G:x}(w_2)$ . Thus, the acceptance models considered here are exactly the same as proposed in [1–3]. The theorem below has been shown in [2].

**Theorem 1.** *The axiomatisation in Figure 1 is sound and complete with respect to the class of acceptance models satisfying constraints **C.4\***, **C.5\***, **C.Inc** and **C.Una**.*

### 2.3 Discussion about the monotonicity principle

One may think that the following principle would be desirable in Acceptance Logic:

$$\text{(Mon)} \quad \neg A_{G:x} \perp \rightarrow \neg A_{H:x} \perp \quad (\text{if } H \subseteq G)$$

which corresponds to semantic constraint: if  $\mathcal{A}_{G:x}(w) \neq \emptyset$  then  $\mathcal{A}_{H:x}(w) \neq \emptyset$ . Principle **Mon** expresses a property of monotonicity about institution membership. It was also discussed from a different perspective in our previous works on Acceptance Logic [1, 2].

We prefer not including **Mon** in the current version of AL because we are interested in strong notions of ‘constituted group’ and ‘group identification’ which are formally expressed by constructions  $\neg A_{G:x} \perp$ . As we said above,  $\neg A_{G:x} \perp$  means “the agents in  $G$  are functioning together as members of the institution  $x$ ” or, stated differently, “ $G$  constitutes a group of members of the institution  $x$ ”. We suppose here that the latter sentences just express that: every agent in  $G$  identifies himself as a member of institution  $x$  and recognizes  $G$  as a *group of members* of institution  $x$ . Under this assumption, **Mon** is not valid. The following example illustrates this point. Imagine that the eleven agents in  $\{1, 2, \dots, 11\}$  constitute a football team (i.e.  $\neg A_{\{1,2,\dots,11\}:team} \perp$ ). This means that every agent in  $\{1, 2, \dots, 11\}$  identifies himself as a member of the football team and recognizes  $\{1, 2, \dots, 11\}$  as a football team. This does not entail that  $\{1, 2, \dots, 10\}$  constitute a football team (i.e.  $\neg A_{\{1,2,\dots,10\}:team} \perp$ ). Indeed, it is not the case that every agent in  $\{1, 2, \dots, 10\}$  recognizes  $\{1, 2, \dots, 10\}$  as a football team. (Only ten players do not constitute a football team).

It worth noting that **Mon** becomes a reasonable principle under a different reading of the construction  $\neg A_{G:x} \perp$ . Namely, suppose that  $\neg A_{G:x} \perp$  just means: every agent in  $G$  identifies himself as a member of institution  $x$  and recognises every agent in the set of agents  $G$  as a member of institution  $x$ . Under this assumption,  $\neg A_{G:x} \perp$  should imply  $\neg A_{H:x} \perp$ , for  $H \subseteq G$ .

## 3 The tableau method

In this section we present a proof method for AL that uses semantic tableaux. As a typical tableaux method, given a formula  $\varphi$ , it systematically tries to construct a model for it. When it fails,  $\varphi$  is inconsistent and thus, its negation is valid.

Each formula in the tableau is prefixed by a natural number that stands for a possible world in the model under construction, similar to the notion used by Fitting ([14, Chapter 8]).

**Definition 1 (Labelled formula).** A labelled formula is a pair of the form  $\langle n, \varphi \rangle$  such that  $n \in \mathbb{N}$  and  $\varphi \in \mathcal{L}_{AL}$ .

Our method also builds the relations  $\mathcal{A}_{G:x}$  between possible worlds that form the model. These relations are represented in tableau by means of arrows between possible worlds, which are represented as triples of the form  $\langle G:x, n, n' \rangle$ . That is, the tableau contains a set of labelled formulae and also a set of such triples, as defined in the sequel.

**Definition 2 (Branch).** A branch is a pair of the form  $\langle L, S \rangle$  such that  $L$  is a set of labelled formulae and  $S \subseteq (\Delta \times \mathbb{N} \times \mathbb{N})$ .

**Definition 3 (Tableau).** Let  $\varphi \in \mathcal{L}_{AL}$ . A tableau for  $\varphi$  is a set of branches  $T$  inductively defined as follows:

- $T = \{\langle \{0, \varphi\}, \emptyset \rangle\}$ . This is called the initial tableau for  $\varphi$ .
- $T' = (T \setminus \{\langle L, S \rangle\}) \cup \{\langle L'_1, S'_1 \rangle, \langle L'_2, S'_2 \rangle, \dots, \langle L'_n, S'_n \rangle\}$ , where  $T$  is a tableau for  $\varphi$  containing the branch  $\langle L, S \rangle$  and each  $\langle L'_i, S'_i \rangle$  is a branch generated by one of the tableau rules defined below: (A more standard presentation of some of these tableau rules is given in Figure 2.)
  - (**R.** $\neg$ ) If  $\langle n, \neg\varphi \rangle \in L$  then generate  $L'_1 = L \cup \{\langle n, \varphi \rangle\}$  and  $S'_1 = S$ .
  - (**R.** $\wedge$ ) If  $\langle n, \varphi_1 \wedge \varphi_2 \rangle \in L$  then generate  $L'_1 = L \cup \{\langle n, \varphi_1 \rangle, \langle n, \varphi_2 \rangle\}$  and  $S'_1 = S$ .
  - (**R.** $\vee$ ) If  $\langle n, \neg(\varphi_1 \wedge \varphi_2) \rangle \in L$  then generate  $L'_1 = L \cup \{\langle n, \neg\varphi_1 \rangle\}$  and  $S'_1 = S$ ; and also  $L'_2 = L \cup \{\langle n, \neg\varphi_2 \rangle\}$  and  $S'_2 = S$ .
  - (**R.** $\Box$ ) if  $\langle n, \mathcal{A}_{G:x}\varphi \rangle \in L$  and  $\langle G:x, n, n' \rangle \in S$  then generate  $L'_1 = L \cup \{\langle n', \varphi \rangle\}$  and  $S'_1 = S$ .
  - (**R.** $\Diamond$ ) if  $\langle n, \neg\mathcal{A}_{G:x}\varphi \rangle \in L$  then generate  $L'_1 = L \cup \{\langle n', \neg\varphi \rangle\}$  and  $S'_1 = S \cup \{\langle G:x, n, n' \rangle\}$ , for some  $n'$  that does not occur in  $L$ .
  - (**R.**4\*) if  $H \subseteq G$  and  $\langle H:y, n, n' \rangle, \langle G:x, n', n'' \rangle \in S$  then generate  $L'_1 = L$  and  $S'_1 = S \cup \{\langle G:x, n, n'' \rangle\}$ .
  - (**R.**5\*) If  $H \subseteq G$  and  $\langle H:y, n, n' \rangle, \langle G:x, n, n'' \rangle \in S$  then generate  $L'_1 = L$  and  $S'_1 = S \cup \{\langle G:x, n', n'' \rangle\}$ .
  - (**R.**Inc) If  $H \subseteq G$  and  $\langle H:x, n, n' \rangle \in S$  then generate  $L'_1 = L$  and  $S'_1 = S \cup \{\langle G:x, n, n' \rangle\}$ ; and also  $L'_2 = L \cup \{\langle n, \mathcal{A}_{G:x}\perp \rangle\}$  and  $S'_2 = S$ .
  - (**R.**Una) If  $G = \{i_1, \dots, i_k\}$  and  $\langle G:x, n, n' \rangle \in S$  then generate  $L'_1 = L$  and  $S'_1 = S \cup \{\langle \{i_1\}:x, n', n' \rangle\}$ ;  $L'_2 = L$  and  $S'_2 = S \cup \{\langle \{i_2\}:x, n', n' \rangle\}$ ;  $\dots$ ;  $L'_k = L$  and  $S'_k = S \cup \{\langle \{i_k\}:x, n', n' \rangle\}$ .

Rules **R.** $\neg$ , **R.** $\wedge$  and **R.** $\vee$  work exactly as for classical propositional logic. Rules **R.** $\Box$  and **R.** $\Diamond$  work exactly as for modal logic K, as proposed, e.g., in [14]. The other rules are meant to ensure that the model to be created will be an acceptance model. Tableau Rule **R.**4\* implements semantic constraint **C.**4\*, Tableau Rule **R.**5\* implements constraint **C.**5\*, Rule **R.**Inc implements **C.**Inc and **R.**Una implements **C.**Una.

(R.□)	$\frac{\langle n, \mathbf{A}_{G:x}\varphi \rangle; \langle G:x, n, n' \rangle}{\langle n', \varphi \rangle;}$	
(R.◇)	$\frac{\langle n, \neg \mathbf{A}_{G:x}\varphi \rangle;}{\langle n', \neg\varphi \rangle; \langle G:x, n, n' \rangle}$	for a new $n'$
(R.4*)	$\frac{; \langle H:y, n, n' \rangle, \langle G:x, n', n'' \rangle}{; \langle G:x, n, n'' \rangle}$	where $H \subseteq G$
(R.5*)	$\frac{; \langle H:y, n, n' \rangle, \langle G:x, n, n'' \rangle}{; \langle G:x, n', n'' \rangle}$	where $H \subseteq G$
(R.Inc)	$\frac{; \langle H:x, n, n' \rangle}{; \langle G:x, n, n' \rangle   \langle n, \mathbf{A}_{G:x}\perp \rangle;}$	where $H \subseteq G$
(R.Una)	$\frac{; \langle G:x, n, n' \rangle}{; \langle \{i_1\}:x, n', n' \rangle   \dots   \langle \{i_k\}:x, n', n' \rangle}$	where $G = \{i_1, \dots, i_k\}$

**Fig. 2.** Tableau rules

**Definition 4 (Closed tableau).** *The set of labelled formulae  $L$  is closed if and only if  $\{\langle n, \varphi \rangle, \langle n, \neg\varphi \rangle\} \subseteq L$ , for some  $n$  and  $\varphi$ . A branch is closed if and only if its set of labelled formulae is closed. A tableau is closed if and only if all its branches are closed. A tableau is open if and only if it is not closed.*

*Example 1.* Now, let us see how the method can be used to show that the formula:  $(\mathbf{A}_{ij:x}\mathbf{A}_{i:x}p \wedge \mathbf{A}_{ij:x}\mathbf{A}_{j:x}p) \rightarrow \mathbf{A}_{ij:x}p$  is valid in acceptance logic. As we will see, if there is a closed tableau for  $\varphi$ , then no model satisfies  $\varphi$ , which means that  $\neg\varphi$  is valid. Therefore, if we provide a closed tableau for the negation of our formula above, we show its validity.

Such a tableau is given in Figure 3. Each line of the figure displays either a labelled formula, an arrow or both, which are the elements of the tableau branches. The number in parentheses on the left is used to identify the line. On the right, also in parentheses, we find the rule that generated that line, and what lines has been used in the application of such rule. For example, the labelled formula and arrow in line 6 have been generated by the application of **R.◇**, using the labelled formula in line 3. Also, lines 9 and 11 have been generated by the application of **R.Una**, using the arrow in line 6.

The input formula is in line 1 (we spelt out the abbreviations). The construction of the closed tableau started with the initial tableau for the input formula. The latter corresponds to line 1 alone. Then, a new tableau has been generated by the application of **R.∧** using the labelled formula in line 1. The latter corresponds to lines 1 and 2 together, and so on. When **R.Una** has been applied

using the arrow in line 6, it generated two branches. This is represented in the figure by the vertical line dividing the tableau in two parts after line 8.

(1) $0, (A_{ij:x}A_{i:x}p \wedge A_{ij:x}A_{j:x}p) \wedge \neg A_{ij:x}p$	
(2) $0, A_{ij:x}A_{i:x}p \wedge A_{ij:x}A_{j:x}p$	(R. $\wedge$ : 1)
(3) $0, \neg A_{ij:x}p$	(R. $\wedge$ : 1)
(4) $0, A_{ij:x}A_{i:x}p$	(R. $\wedge$ : 2)
(5) $0, A_{ij:x}A_{j:x}p$	(R. $\wedge$ : 2)
(6) $1, \neg p$	$ij:x, 0, 1$ (R. $\diamond$ : 3)
(7) $1, A_{i:x}p$	(R. $\square$ : 4, 6)
(8) $1, A_{j:x}p$	(R. $\square$ : 5, 6)
(9) $i:x, 1, 1$ (R.Una : 6)	(11) $j:x, 1, 1$ (R.Una : 6)
(10) $1, p$ closed	(12) $1, p$ closed
(6, 10)	(6, 12)

**Fig. 3.** A tableau for Example 1

*Example 2.* On the other hand, if no closed tableau for  $\varphi$  exists, then  $\varphi$  is satisfiable. Let us see what happens when we try to generate a closed tableau for the formula:  $\neg A_{i:x}p \wedge A_{ij:x} \neg A_{i:x}p$  which is satisfiable. This is done in Figure 4. Note that one of the branches is closed. On the other hand, no rule can be applied in the rightmost branch, which means that this tableau will remain open. On the leftmost branch we have a rather different phenomenon. We can continue applying the same set of rules indefinitely. This will generate more branches that can be closed, but we can never close all of them. This means that we can also consider such a branch as an open one.

We proceed by proving soundness of the method. But first we need yet another definition.

**Definition 5 (Satisfiable branch).** *The branch  $b = \langle L, S \rangle$  is satisfiable if and only if there exists an acceptance model  $M = \langle W, \mathcal{A}, \mathcal{V} \rangle$  and a function  $f : \mathbb{N} \rightarrow W$  such that:*

1.  $\langle f(n), f(n') \rangle \in \mathcal{A}_{G:x}$ , for all  $\langle G:x, n, n' \rangle \in S$ ; and
2.  $M, f(n) \models \varphi$ , for all  $\langle n, \varphi \rangle \in L$ .

**Theorem 2 (Soundness).** *If there is a closed tableau for  $\neg\varphi$  then  $\varphi$  is valid.*

*Proof.* We show that if  $\varphi$  is satisfiable then there is no closed tableau for  $\varphi$ . It is enough to show that all tableau rules preserve satisfiability. That is, it is enough to show that: if the branch  $b = \langle L, S \rangle$  is satisfiable then the set of branches  $B = \{b'_1, \dots, b'_k\}$  generated by any tableau rule contains a satisfiable branch.



(1) $0, \neg A_{i:x}p \wedge A_{ij:x} \neg A_{i:x}p$	
(2) $0, \neg A_{i:x}p$	( <b>R.<math>\wedge</math></b> : 1)
(3) $0, A_{ij:x} \neg A_{i:x}p$	( <b>R.<math>\wedge</math></b> : 1)
(4) $1, \neg p$	$i:x, 0, 1$ ( <b>R.<math>\diamond</math></b> : 2)
(5) $ij:x, 0, 1$ ( <b>R.Inc</b> : 4)	(16) $0, A_{ij:x} \perp$ ( <b>R.Inc</b> : 4)
(6) $1, \neg A_{i:x}p$	( <b>R.<math>\square</math></b> : 3, 5)
(7) $ij:x, 1, 1$ ( <b>R.5*</b> : 5)	
(8) $2, \neg p$	$i:x, 1, 2$ ( <b>R.<math>\diamond</math></b> : 6)
(9) $ij:x, 1, 2$ ( <b>R.Inc</b> : 8)	(14) $1, A_{ij:x} \perp$ ( <b>R.Inc</b> : 8)
(10) $ij:x, 0, 2$ ( <b>R.4*</b> : 4, 9)	(15) $1, \perp$ ( <b>R.<math>\square</math></b> : 14, 7)
(11) $2, \neg A_{i:x}p$	closed (15)
(12) $ij:x, 2, 2$ ( <b>R.5*</b> : 9)	
(13) $3, \neg p$	$i:x, 2, 3$ ( <b>R.<math>\diamond</math></b> : 11)
⋮	

**Fig. 4.** A tableau for Example 2

Indeed, and to see that it is enough, suppose that  $b'_i$  is satisfiable and closed. Then  $L'_i$  contains two labelled formulae  $\langle n, \varphi \rangle$  and  $\langle n, \neg \varphi \rangle$ . Because  $b'_i$  is satisfiable, there exists an acceptance model  $M$  and a function  $f$  such that  $M, f(n) \models \varphi$  and  $M, f(n) \models \neg \varphi$ , which is a contradiction.

Now, suppose that the branch  $b$  is satisfiable. The proof that the rules **R. $\neg$**  and **R. $\wedge$**  preserve satisfiability is straightforward and thus, left as an exercise to the reader. We proceed by showing that the modal rules preserve satisfiability.

**R. $\square$** :  $M, f(n) \models A_{G:x} \varphi$  and  $\langle f(n), f(n') \rangle \in \mathcal{A}_{G:x}$  (by hypothesis). Then  $M, f(n') \models \varphi$  (by definition).

**R. $\diamond$** :  $M, f(n) \models \neg A_{G:x} \varphi$  (by hypothesis). Then there exists  $w' \in \mathcal{A}_{G:x}(f(n))$  such that  $M, w' \models \neg \varphi$  (again, by definition). Now, consider the function  $f' : \mathbb{N} \rightarrow W$  such that  $f'(n) = f(n)$ , for all  $n$  occurring in  $L$ , and  $f'(n') = w'$ . Then  $M, f'(n'') \models \varphi''$ , for all  $\langle n'', \varphi'' \rangle \in L$  (because  $n'$  does not occur in  $L$ ), and  $M, f'(n') \models \neg \varphi$ .

**R.4\***: Let  $H \subseteq G$ .  $\langle f(n), f(n') \rangle \in \mathcal{A}_{H:y}$  and  $\langle f(n'), f(n'') \rangle \in \mathcal{A}_{G:x}$  (by hypothesis). Then  $\langle f(n), f(n'') \rangle \in \mathcal{A}_{G:x}$ , since  $M$  is an acceptance model respecting **C.4\***. Then, the branch  $\langle L'_1, S'_1 \rangle$  is satisfiable.

**R.5\***: Let  $H \subseteq G$ .  $\langle f(n), f(n') \rangle \in \mathcal{A}_{H:x}$  and  $\langle f(n), f(n'') \rangle \in \mathcal{A}_{G:x}$  (by hypothesis). Then,  $\langle f(n'), f(n'') \rangle \in \mathcal{A}_{G:x}$ , since  $M$  is an acceptance model respecting **C.5\***. Then, the branch  $\langle L'_1, S'_1 \rangle$  is satisfiable.

**R.Inc**: Let  $H \subseteq G$ .  $M, f(n) \models A_{G:x} \varphi$  and  $\langle f(n), f(n') \rangle \in \mathcal{A}_{H:x}$  (by hypothesis). Note that  $M, f(n') \not\models \perp$  (because we assume that the branch is satisfiable) then  $M, f(n) \models \neg A_{H:x} \perp$ . The latter implies  $\langle f(n), f(n') \rangle \in \mathcal{A}_{G:x}$  or  $\mathcal{A}_{G:x}(f(n)) = \emptyset$ , since  $M$  is an acceptance model respecting **C.Inc**. Therefore, one of the branches generated by **R.Inc** is satisfiable.

**R.Una**: Let  $G = \{i_1, \dots, i_k\}$ .  $\langle f(n), f(n') \rangle \in \mathcal{A}_{G:x}$  (by hypothesis). Then,  $\langle f(n'), f(n') \rangle \in \mathcal{A}_{i_j:x}$  for some  $1 \leq j \leq k$ , since  $M$  is an acceptance model

respecting **C.Una**. Therefore, one of the branches generated by **R.Una** is satisfiable.  $\square$

In the sequel we prove completeness. First though, we need another auxiliary definition.

**Definition 6 (Saturated tableau).** *Let  $T$  be a tableau for  $\varphi$  containing the branch  $b = \langle L, S \rangle$ . The branch  $b$  is ‘saturated under the tableau rule  $\rho$ ’ if and only if  $L$  contains  $L'_i$  and  $S$  contains  $S'_i$  for some branch  $\langle L'_i, S'_i \rangle$  generated by the application of the rule  $\rho$  to  $b$ . A branch is (simply) ‘saturated’ if and only if it is saturated under all tableau rules. A tableau is saturated if and only if all its branches are saturated.*

**Theorem 3 (Completeness).** *If  $\varphi$  is valid then there exists a closed tableau for  $\neg\varphi$ .*

*Proof.* Let  $\langle L, S \rangle$  be an open and saturated branch of the tableau. We build a model  $M = \langle W, \mathcal{A}, \mathcal{V} \rangle$  such that  $W = \{n \mid \langle n, \varphi \rangle \in L\}$ , for some  $\varphi$ ;  $\mathcal{A}_{G:x}(n) = \{n' \mid \langle G:x, n, n' \rangle \in S\}$ ; and  $\mathcal{V}(p) = \{n \mid \langle n, p \rangle \in L\}$ .

Clearly,  $M$  is an acceptance model, because the branch is saturated under **R.4\***, **R.5\***, **R.Inc** and **R.Una**.

Now, we show that for all  $\langle n, \varphi \rangle \in L$ ,  $M, n \models \varphi$ . It is done by induction on the structure of  $\varphi$ .

There are two cases in the induction base. (1)  $\varphi = p$ , i.e.,  $\langle n, p \rangle \in L$ . Then  $n \in \mathcal{V}(p)$ , iff  $M, n \models p$  (by definition). (2)  $\varphi = \neg p$ , i.e.,  $\langle n, \neg p \rangle \in L$ . Then  $\langle n, p \rangle \notin L$ , because  $L$  is open. Then  $n \notin \mathcal{V}(p)$ , iff  $M, n \models \neg p$  (by definition).

There are five cases in the induction step. (1)  $\varphi = \neg\neg\varphi_1$ , i.e.,  $\langle n, \neg\neg\varphi_1 \rangle \in L$ . Then  $\langle n, \varphi_1 \rangle \in L$ , because  $L$  is saturated under **R. $\neg$** . Then  $M, n \models \varphi_1$ , by induction hypothesis, iff  $M, n \models \neg\neg\varphi_1$ , by definition. Cases (2)  $\varphi = \varphi_1 \wedge \varphi_2$  and (3)  $\varphi = \neg(\varphi_1 \wedge \varphi_2)$ , are shown analogously as case (1) using rules **R. $\wedge$**  and **R. $\vee$** , respectively. They are left as an exercise to the reader. (4)  $\langle n, A_{G:x}\varphi \rangle \in L$ . Then for all  $n' \in W$ , if  $\langle G:x, n, n' \rangle \in S$ , then  $\langle n', \varphi \rangle \in L$  (because  $L$  is saturated under rule **R. $\Box$** ). Then for all  $n' \in W$ , if  $n' \in \mathcal{A}_{G:x}(n)$ , then  $M, n' \models \varphi$ , by induction hypothesis. Therefore,  $M, n \models A_{G:x}\varphi$ . (5)  $\langle n, \neg A_{G:x}\varphi \rangle \in L$ . Then there is  $n' \in W$  such that  $\langle G:x, n, n' \rangle \in S$  and  $\langle n', \neg\varphi \rangle \in L$  (because  $L$  is saturated under rule **R. $\Diamond$** ). Then there is  $n' \in W$  such that  $n' \in \mathcal{A}_{G:x}(n)$  and  $M, n' \models \neg\varphi$ , by induction hypothesis. Therefore,  $M, n \models \neg A_{G:x}\varphi$ .  $\square$

**Theorem 4 (Termination).** *There exists an implementation of the tableau method that halts for every input formula  $\varphi$ .*

*Proof (Sketch).* We assume an implementation of the tableau method that employs a ‘loop-test’. That is, a procedure that, once the latest generated tableau is saturated under all rules but **R. $\Diamond$** , verifies whether the application of the latter rule to the witness formula  $\psi$  will generate a world such that the set of labelled formulae having this world as label and the set of arrows involving this world will be included in the respective sets for a different world already present in the

branch. If it is the case, then  $\psi$  is marked and **R. $\diamond$**  will not be applied using this formula any more.

Then, the argument for termination goes as for logic S4 (as used, e.g., in [15, 16]): the formulae generated by the rules are in a finite set  $S$ , and therefore only a finite number of different nodes can be generated by the tableaux procedure. The only difference is that here the finite set  $S$  is not just the set of sub-formulae of the initial formula  $\varphi$  as for S4, but its closure, which is the set of sub-formulae of  $\varphi$  union the set of  $A_{G;x}\perp$  such that  $G$  and  $x$  occur in  $\varphi$  (due to rule **R.Inc**). Moreover, the set of labels of every arrow is finite.  $\square$

## 4 An example: The discursive dilemma

In the recent years many researchers in philosophy, computer science and political sciences have been working on the issue of judgement aggregation (e.g., [17–22]). The problem is: How can a group of individuals aggregate the group members’ individual judgements on some interconnected propositions into corresponding collective judgements on these propositions? Such problems occur in different social and legal contexts like committees, legislatures, judiciaries and expert panels.

Our logic of acceptance is a formal framework in which some important aspects of judgement aggregation can be modelled. Moreover, the tableau method for the logic of acceptance presented in Section 3 provides an interesting solution for making automated reasoning about judgement aggregation.

In the logic of acceptance the problem of judgement aggregation is a particular case of the problem of explaining how collective acceptance of the members of a certain group in an institutional context  $x$  about a certain fact  $\varphi$  is created from the individual acceptances in  $x$  of the members of the same group.

We here consider a well-known problem in judgement aggregation called ‘doctrinal paradox’ or ‘discursive dilemma’ [12, 13]. The scenario of the discursive dilemma is a three-member court which has to judge whether a defendant is liable for a breach of contract. According to the legal doctrine, the defendant is liable (*lia*) if and only if he did a certain action (*act*) and he had a contractual obligation not to do this action (*obl*). This is expressed in propositional logic by the connection rule  $lia \leftrightarrow (act \wedge obl)$ . The three judges use majority rule to decide on this issue. The opinions of the three judges are given in Table 1.

	<i>act</i>	<i>obl</i>	$lia \leftrightarrow (act \wedge obl)$	<i>lia</i>
Judge 1	yes	yes	yes	yes
Judge 2	yes	no	yes	no
Judge 3	no	yes	yes	no
Majority	yes	yes	yes	no

**Table 1.** Discursive dilemma

It is supposed that all the judges accept the rule  $lia \leftrightarrow (act \wedge obl)$ . Judge 1 accepts both  $act$  and  $obl$  and, by the connection rule, he accepts  $lia$ . Judge 2 accepts  $act$  and rejects  $obl$  and, by the connection rule, he rejects  $lia$ . Finally, judge 3 rejects  $act$  and accepts  $obl$  and, by the connection rule, he rejects  $lia$ . If the three judges apply a majority rule on each proposition then they face a paradox. Indeed, a majority accepts  $act$ , a majority accepts  $obl$ , a majority accepts the connection rule  $lia \leftrightarrow (act \wedge obl)$ . But the majority rejects  $lia$ . Thus, when majority voting is applied to each single proposition it yields an inconsistent collective set of judgements (see the last row in Table 1). Note that this inconsistency occurs even though the sets of judgements of the individual judges are all consistent.

Let us now show how the discursive dilemma can be formalised in the logic of acceptance.

We first suppose that 1, 2 and 3 *qua* judges of the court accept the connection rule:

$$\neg A_{123:c} \perp \quad (1)$$

$$A_{123:c}(lia \leftrightarrow (act \wedge obl)) \quad (2)$$

Then, we have that judge 1 announces that, *qua* judge of the court, he accepts  $act \wedge obl$ . Judge 2 announces that, *qua* judge of the court, he accepts  $act \wedge \neg obl$ . Judge 3 announces that, *qua* judge of the court, he accepts  $\neg act \wedge obl$ . This has the following effect:

$$A_{123:c} A_{1:c}(act \wedge obl) \quad (3)$$

$$A_{123:c} A_{2:c}(act \wedge \neg obl) \quad (4)$$

$$A_{123:c} A_{3:c}(\neg act \wedge obl) \quad (5)$$

Finally, the three judges use a majority principle for each proposition  $act$ ,  $obl$  and  $lia$ . This majority principle is formally expressed by the following six hypotheses:

$$\begin{aligned} \mathbf{Maj} = \{ & A_{123:c} \bigwedge_{i,j \in \{1,2,3\}, i \neq j} ((A_{i:c} act \wedge A_{j:c} act) \rightarrow act), \\ & A_{123:c} \bigwedge_{i,j \in \{1,2,3\}, i \neq j} ((A_{i:c} \neg act \wedge A_{j:c} \neg act) \rightarrow \neg act), \\ & A_{123:c} \bigwedge_{i,j \in \{1,2,3\}, i \neq j} ((A_{i:c} obl \wedge A_{j:c} obl) \rightarrow obl), \\ & A_{123:c} \bigwedge_{i,j \in \{1,2,3\}, i \neq j} ((A_{i:c} \neg obl \wedge A_{j:c} \neg obl) \rightarrow \neg obl), \\ & A_{123:c} \bigwedge_{i,j \in \{1,2,3\}, i \neq j} ((A_{i:c} lia \wedge A_{j:c} lia) \rightarrow lia), \\ & A_{123:c} \bigwedge_{i,j \in \{1,2,3\}, i \neq j} ((A_{i:c} \neg lia \wedge A_{j:c} \neg lia) \rightarrow \neg lia) \} \end{aligned}$$

It is possible to prove that 1, 2, 3, 4 and 5 together with **Maj** lead to a contradiction using the axiomatisation of acceptance logic. Indeed, from hypotheses 3, 4 and 5 we infer

$$A_{123:c}(A_{1:c}act \wedge A_{2:c}act) \wedge A_{123:c}(A_{1:c}obl \wedge A_{3:c}obl).$$

By the first and third hypotheses in **Maj**, the latter implies

$$A_{123:c}act \wedge A_{123:c}obl$$

and by hypothesis 2 and standard modal principles the latter implies

$$A_{123:c}lia.$$

From hypotheses 1 and 2, by axioms **4\*** and **5\***, we can infer

$$A_{123:c}(A_{123:c}(lia \leftrightarrow (act \wedge obl)) \wedge \neg A_{123:c}\perp).$$

By Axiom **Inc** and standard modal principles, the latter implies

$$A_{123:c}(A_{1:c}(lia \leftrightarrow (act \wedge obl)) \wedge A_{2:c}(lia \leftrightarrow (act \wedge obl)) \wedge A_{3:c}(lia \leftrightarrow (act \wedge obl))).$$

From this, by hypotheses 3, 4 and 5 and standard modal principles, we can infer

$$A_{123:c}(A_{1:c}lia \wedge A_{2:c}\neg lia \wedge A_{3:c}\neg lia)$$

and by the sixth hypothesis in **Maj** the latter implies

$$A_{123:c}\neg lia.$$

Thus, we have  $A_{123:c}\perp$ , and by hypothesis 1 we can infer  $\perp$ .

However, if one slightly changes the hypotheses, the proof may change completely. This means that the axiomatisation does not provide a straightforward way to automatise the process of finding such kind of inconsistency. On the other hand, a tableau method is meant to provide a systematic way to search for models and, thereby, an easy way to automatise the process of deciding whether a formula is inconsistent or not.

As an illustration, we show a closed tableau for this example (schematically) in Figure 5. On that figure, the lines identified by H.1–5 correspond to hypotheses 1–5 above, respectively. The line identified by H.6 corresponds to the first hypothesis in **Maj**, H.7 to the third hypothesis, and H.8 to the sixth. The branch containing lines 13 and 14 is already closed. The other branches are not completely displayed. But it is easy to see that all branches generated from line 12 on can be closed, since lines 8, 9 and 12 together are inconsistent in propositional logic. The branches generated from line 15 on can be closed in an analogous way as the latter. And the branches generated from line 5 on need more effort, but they can be closed too, in a similar way as the latter by using lines H.2, H.3, H.6 and H.7.

## 5 Conclusion

The contribution of this paper is a semantic tableau method for acceptance logic. The method consists in a procedure to check satisfiability of formulae, that can be easily automatized. Given that acceptance logic can be used to formalise some aspects of judgement aggregation, our method also provides an automated reasoning procedure for making judgement aggregation in modal logic.

It is to be noted that, differently from [21, 22], in which logical approaches specialised for judgement aggregation have been proposed, in the logic of accep-

(H.1) 0, $\neg A_{123:c} \perp$			
(H.2) 0, $A_{123:c}(bia \leftrightarrow (act \wedge obl))$			
(H.3) 0, $A_{123:c}A_{1:c}(act \wedge obl)$			
(H.4) 0, $A_{123:c}A_{2:c}(act \wedge \neg obl)$			
(H.5) 0, $A_{123:c}A_{3:c}(\neg act \wedge obl)$			
(H.6) 0, $A_{123:c}(A_{1:c}act \wedge A_{3:c}act) \rightarrow act$			
(H.7) 0, $A_{123:c}(A_{1:c}obl \wedge A_{2:c}obl) \rightarrow obl$			
(H.8) 0, $A_{123:c}(A_{2:c}\neg bia \wedge A_{3:c}\neg bia) \rightarrow \neg bia$			
(1) 1, $\neg \perp$			
(2) 1, $A_{2:c}(act \wedge \neg obl)$		123:c, 0, 1 (R.◇ : H.1)	
(3) 1, $(A_{2:c}\neg bia \wedge A_{3:c}\neg bia) \rightarrow \neg bia$		(R.□ : H.4, 1)	
(4)		(R.□ : H.8, 1)	
		123:c, 1, 1 (R.5* : 1)	
(5) 1, $\neg bia$ (R.V : 3)			
⋮			
closed (H.2, H.3, H.6, H.7)			
	(6) 1, $\neg(A_{2:c}\neg bia \wedge A_{3:c}\neg bia)$ (R.V : 3)		
	(7) 1, $\neg A_{2:c}\neg bia$ (R.V : 6)		
	(8) 2, $\neg bia$ 2:c, 1, 2 (R.◇ : 7)		
	(9) 2, $act \wedge \neg obl$ (R.□ : 2, 8)		
(10)	123:c, 1, 2 (R.Inc : 8)	(13) 1, $A_{123:c} \perp$ (R.Inc : 8)	
(11)	123:c, 0, 2 (R.4* : 1, 10)	(14) 1, $\perp$ (R.□ : 13, 4)	
(12) 2, $bia \leftrightarrow (act \wedge obl)$	(R.□ : H.2, 11)	closed	(14)
⋮			
closed	(8, 9, 12)		
		(15) 1, $\neg A_{3:c}\neg bia$ (R.V : 13)	
		⋮	
		closed (H.2, H.5, 15)	

Fig. 5. A closed tableau for the discursive dilemma

tance judgement aggregation is just an application. We have shown in [2, 3] that the logic of acceptance is a much general formal framework in which the static and dynamic aspects of institutions can be studied (e.g., static and dynamic aspects of social roles, norms and rules).

Moreover, this method uses a rather different set of rules. For instance, rules for axioms **4\*** and **5\*** (which are nothing but a variation of the usual axioms 4 and 5) are similar to the so-called “structural rules” proposed in [23]. Our tableau rules are also modular. This means that it is possible to provide a sound, complete and terminating satisfiability checking method for a logic without one of the semantic constraints, by just removing the corresponding tableau rule. The addition of rules is also possible. For example, one could add the tableau rule:

$$(\mathbf{R.Mon}) \quad \frac{}{\begin{array}{l} ; \langle G:x, n, n' \rangle \\ ; \langle H:x, n, n' \rangle \end{array}}$$

which corresponds to Axiom **Mon**, thus, obtaining a tableau method for the acceptance logic with Axiom **Mon** proposed in [1, 2].

As possible future works, we intend to investigate computational complexity, and also possible extensions of our method able to address acceptance logic with dynamic operators, such as the logic proposed in [24].

## 6 Acknowledgements

The contribution by Tiago de Lima is part of the research program Moral Responsibility in R&D Networks, supported by the Netherlands Organisation for Scientific Research (NWO), under grant number 360-20-160.

## References

1. Gaudou, B., Longin, D., Lorini, E., Tummolini, L.: Anchoring institutions in agents’ attitudes: Towards a logical framework for autonomous MAS. In Padgham, L., Parkes, D.C., eds.: Proceedings of AAMAS’08. (2008) 728–735
2. Lorini, E., Longin, D., Gaudou, B., Herzig, A.: The logic of acceptance: Grounding institutions on agents’ attitudes. *Journal of Logic and Computation* (2009) doi:10.1093/logcom/exn103.
3. Lorini, E., Longin, D.: A logical account of institutions: From acceptances to norms via legislators. In Brewka, G., Lang, J., eds.: Proceedings of KR’08. (2008) 38–48
4. Bratman, M.E.: Practical reasoning and acceptance in context. *Mind* **101**(401) (1992) 1–15
5. Cohen, L.J.: An essay on belief and acceptance. Oxford University Press, New York, USA (1992)
6. Tuomela, R.: Belief versus Acceptance. *Philosophical Explorations* **2** (2000) 122–137
7. Tuomela, R.: The Philosophy of Social Practices: A Collective Acceptance View. Cambridge University Press, Cambridge (2002)
8. Boella, G., van der Torre, L.: Norm negotiation in multiagent systems. *International Journal of Cooperative Information Systems* **16**(1) (2007) 97–122

9. Hart, H.L.A.: The concept of law. Clarendon Press, Oxford (1992) new edition.
10. Gilbert, M.: On Social Facts. Routledge, London and New York (1989)
11. Herzig, A., de Lima, T., Lorini, E.: On the dynamics of institutional agreements. Manuscript (2009)
12. Pettit, P.: Deliberative democracy and the discursive dilemma. *Philosophical Issues* **11** (2001) 268–99
13. Kornhauser, L.A., Sager, L.G.: Unpacking the court. *Yale Law Journal* **96** (1986) 82–117
14. Fitting, M.: Proof Methods for Modal and Intuitionistic Logics. Springer (1983)
15. Halpern, J., Moses, Y.: A guide to completeness and complexity for modal logics of knowledge and belief. *Artificial Intelligence* **54** (1992) 311–379
16. Goré, R.: Tableau methods for modal and temporal logics. In D’Agostino, M., Gabbay, D.M., Hahnle, R., Posegga, J., eds.: *Handbook of Tableau Methods*. Springer (1999) 297–396
17. Pauly, M., van Hees, M.: Logical constraints on judgment aggregation. *Journal of Philosophical Logic* **35**(6) (2006) 569–585
18. List, C., Pettit, P.: Aggregating sets of judgments: An impossibility result. *Economics and Philosophy* **18** (2002) 89–110
19. Goldman, A.: Group knowledge versus group rationality: Two approaches to social epistemology. *Episteme* **1**(1) (2004) 11–22
20. List, C.: Group knowledge and group rationality: A judgment aggregation perspective. *Episteme* **2**(1) (2005) 25–38
21. Ágotnes, T., van der Hoek, W., Wooldridge, M.: Reasoning about judgment and preference aggregation. In: *Proceedings of AAMAS’07*. (2007) 566–573
22. Pauly, M.: Axiomatizing collective judgment sets in a minimal logical language. *Synthese* **158** (2007) 233–250
23. Castilho, M.A., Fariñas del Cerro, L., Gasquet, O., Herzig, A.: Modal tableaux with propagation rules and structural rules. *Fundamenta Informaticae* **20** (1998) 1–17
24. Herzig, A., de Lima, T., Lorini, E.: What do we accept after an announcement? In Meyer, J.-J.Ch., Broersen, J. eds.: *Pre-proceedings of the KR’08-Workshop KRAMAS*, pages 81–94, (Available at: <http://www.cs.uu.nl/events/kramas2008/PreProceedingsKRAMAS2008.pdf>) (2008)



# Ontology and Time Evolution of Obligations and Prohibitions using Semantic Web Technology

Nicoletta Fornara<sup>1</sup> and Marco Colombetti<sup>1,2</sup>

<sup>1</sup> Università della Svizzera italiana, via G. Buffi 13, 6900 Lugano, Switzerland  
{nicoletta.fornara, marco.colombetti}@lu.unisi.ch,

<sup>2</sup> Politecnico di Milano, piazza Leonardo Da Vinci 32, Milano, Italy  
marco.colombetti@polimi.it

**Abstract.** The specification and monitoring of conditional obligations and prohibitions with starting points and deadlines is a crucial aspect in the design of open interaction systems. In this paper we regard such obligations and prohibitions as cases of social commitment, and propose to model them in OWL, the logical language recommended by the W3C for Semantic Web applications. In particular we propose an application-independent ontology of the notions of social commitment, temporal proposition, event, agent, role and norms that can be used in the specification of any open interaction system. We then delineate a hybrid solution that uses the ontology, SWRL rules, and a Java program to dynamically monitor or simulate the temporal evolution of social commitments, due to the elapsing of time and to the actions performed by the agents interacting within the system.

## 1 Introduction

The specification of *open interaction systems*, where heterogeneous, autonomous, and self-interested agents can interact by entering and leaving dynamically the system, is widely recognized to be a crucial issue in the development of distributed applications on the Internet, like e-commerce applications, or collaborative applications for the automatic creation of virtual organizations. An important aspect of the specification of open systems is the possibility to define the actions that agents should or should not perform in a given interval of time, that is, the possibility to define social commitments with starting points and deadlines, and to monitor and react to their fulfilment or violation.

As we discussed in our previous works [9, 10, 8] in our OCeAN meta-model for the specification of artificial institutions, commitments for the interacting agents can be created by the activation of norms associated to the agents' roles, or by the performance of agent communicative acts, like promises. In this paper we explore how to use OWL (in its OWL 2 DL version<sup>3</sup>), the logical language recommended by W3C for Semantic Web applications, to specify the deontic part of the OCeAN meta-model. More precisely, we show how it is possible to

<sup>3</sup> [http://www.w3.org/2007/OWL/wiki/OWL/Working\\_Group/](http://www.w3.org/2007/OWL/wiki/OWL/Working_Group/)

specify social commitment to express conditioned obligations and prohibitions on time intervals, in OWL.

There are many advantages in using a decidable logical language like OWL to specify an open interaction system, and in particular that: (i) Semantic Web technologies are increasingly becoming a standard for Internet applications; (ii) the language is supported by reasoners (like Fact++<sup>4</sup>, Pellet<sup>5</sup>, and the rule reasoner of the Jena Semantic Web framework<sup>6</sup>) that are more efficient than available alternatives (like the Discrete Event Calculus Reasoner<sup>7</sup>); (iii) it is possible to achieve a high degree of interoperability of data and applications, which is indeed a crucial precondition for the development of open systems.

The idea of using OWL for modelling and monitoring the dynamic evolution of open artificial institutions can be developed following different approaches. A first option would be to implement an institutional model in a object oriented language like Java, and use OWL only to specify the ontology of the content of communicative acts and norms. As a result reasoning may be used to deduce, for example, that the performance of a certain act implies the performance of another act, and thus the fulfillment of a given commitment. An alternative approach, which we investigate in this paper, consists in using OWL to express, as far as possible, the normative component of the OCeAN meta-model. As we shall see, this requires the use of SWRL (Semantic Web Rule Language<sup>8</sup>) and Java code to overcome certain expressiveness limitations of OWL. Indeed, with both OWL 1 (the current standard) and OWL 2 there are at least two major problems:

- The treatment of time. OWL has no temporal operators; on some occasions it is possible to bypass the problem by using SWRL rules with temporal built-ins, but in any case this does not provide full temporal reasoning capabilities.
- The open-world assumption. In many applications, nor being able to infer that an action has been performed is sufficient evidence that the action has not been performed; one would then like to infer, for example, that an obligation to perform the action has been violated. As standard OWL reasoning is carried out under the open world assumption, inferences of this type cannot be drawn. However, it is often possible to simulate a closed world assumption by adding closure axioms to an ontology.

The main contribution of this paper, with respect to our previous works, is to show how obligations and prohibitions can be formalized in OWL and SWRL for monitoring and simulation purposes with significant performance improvements with respect to the solution based on the Event Calculus that we presented elsewhere [8]. Another contribution of this work is a hybrid solution of the problem of monitoring the temporal evolution of obligations and prohibition, based on a

---

<sup>4</sup> <http://owl.man.ac.uk/factplusplus/>

<sup>5</sup> <http://clarkparsia.com/pellet>

<sup>6</sup> <http://jena.sourceforge.net/inference/>

<sup>7</sup> <http://decreasoner.sourceforge.net>

<sup>8</sup> <http://www.w3.org/Submission/SWRL/>

higher ontology of interaction, a set of SWRL rules, and a Java program implemented using suitable OWL libraries (like the Jena Semantic Web Framework for Java<sup>9</sup> or OWL API<sup>10</sup>).

The paper is organized as follows. In the next section we briefly introduce OWL and SWRL, that is, the Semantic Web languages that we use to formally specify the normative component of an open interaction system. In Section 3 we specify the algorithms that we plan to use to simulate or monitor the temporal evolution of an interaction system. Then in Section 4 we define the classes, properties, axioms, and rules that we take to underlie the normative specification of every interaction system. In Section 5 we present an actual example of a system specified using our meta-model. Finally in Section 6 we compare our approach with other proposals and draw some conclusions.

## 2 OWL and SWRL

OWL is a practical realization of a Description Logic system known as *SHOIN(D)*. It allows one to define *classes* (also called *concepts* in the DL literature), *properties* (also called *roles*), and *individuals*. An OWL ontology consists of: a set of class axioms to describe classes, which constitute the *Terminological Box (TBox)*; a set of property axioms to describe properties, which constitute a *Role Box (RBox)*; and a collection of assertions to describe individuals, which constitute an *Assertion Box (ABox)*.

Classes can be viewed as formal descriptions of sets of objects (taken from a nonempty universe), and individuals can be viewed as names of objects of the universe. Properties can be either *object properties* or *data properties*. The former describe binary relations between objects of the universe; the latter, binary relationships between objects and data values (taken from ML Schema datatypes).

A class is either a *basic class* (i.e., an atomic class name) or a *complex class* build through a number of available *constructors* that express Boolean operations and different types of restrictions on the members of the class.

Through *class axioms* one may specify subclass or equivalence relationships between classes, and that certain classes are disjoint. *Property axioms* allow one to specify that a given property is a subproperty of another property, that a property is the inverse of another property, or that a property is functional or transitive. Finally, *assertions* allow one to specify that an individual belongs to a class, that an individual is related to another individual through an object property, that an individual is related to a data value through a data property, or that two individuals are equal or different.

OWL can be regarded as a decidable fragment of First Order Logic (FOL). The price one pays for decidability, which is considered as an essential precondition for exploiting reasoning in practical applications, is limited expressiveness.

---

<sup>9</sup> <http://jena.sourceforge.net/>

<sup>10</sup> <http://owlapi.sourceforge.net/>

Even in OWL 2 (the more expressive version currently under specification) certain useful first-order statements cannot be formalized.

Recently certain OWL reasoners, like Pellet, have been extended to deal with SWRL rules. SWRL is a Datalog-like language, in which certain universally quantified conditional axioms (called *rules*) can be stated. To preserve decidability, however, rules have to be used in the *safe mode*, which means that before being exploited in a reasoning process all their variables must be instantiated by pre-existing individuals. An important aspect of SWRL is the possibility of including *built-ins*, that is, Boolean functions that perform operations on data values and return a truth value.

### Conventions

In what follows we use the notation  $p : C \rightarrow_O D$  to specify an *object property*  $p$  (not necessarily a function) with class  $C$  as domain and class  $D$  as range, and the notation  $q : C \rightarrow_D T$  to specify a *data property*  $q$  with class  $C$  as domain and the datatype  $T$  as range. We use capital initials for classes, and lower case initials for properties and individuals.

## 3 Specification and simulation or monitoring of an open interaction system

Our approach is to model an open interaction system using one or more artificial institutions. The definition of a specific artificial institution consists of: (i) a first component, called meta-model, which includes the definition of basic entities common to the specification of every institution, like the concepts of temporal proposition, commitment, institutional power, role, and norm, and the actions necessary for exchanging messages; (ii) a second component, pertaining to the institution in exam, which includes the definition of specific powers and norms that apply to the agents playing roles in the institution, and the definition of the concepts pertaining to the domain of the interaction (for example the actions of paying or delivering a product, bidding in an auction, etc.).

We start from the specification of a system, formalized as an application-independent OWL ontology (including a TBox, an RBox, and an ABox as detailed in Section 4). We then add an application-dependent ontology (as exemplified in Section 5) and use a Java program to let such ABox evolve in time, with the goal of *monitoring* the fulfilment or violation of obligations and prohibitions, and of *simulating* the evolution of the state of the system against different possible history of events.

In particular, when the system is used for monitoring, a Java program updates the state of the system, that is, it updates the ABox with new assertions to model the elapsing of time, to allow for closed-world reasoning on certain classes, and to model the actions performed by the interacting agents. When such updating is completed, a reasoner can be used to deduce the state of obligations and prohibitions. After that, when the ontology has reached a stable state (in the sense that all closed-world reasoning has been completed), the agents may

perform queries to know what are their pending obligations or prohibitions or to react to their violation or fulfillment. We assume that the events or actions that happen between two phases of update (that is, between two discrete instant of time) are queued in the data structure `ActionQueue` to be managed by the Java program subsequently.

When the system is used for simulation, the set of events that happen at run-time are known since the beginning, and are represented in the initial version of the ABox. In such a case the Java program simply updates the state of the system to represent the elapsing of time and to allow closed-world reasoning on certain classes; then the reasoner deduces the state of obligations and prohibitions at each time instant.

### Temporal evolution of the ontology

An external Java program is used to model the elapsing of time, the actions performed by the interacting agents at run-time (in the monitoring usage), and to allow for closed-world reasoning on certain classes (see Section 4.1 for details). The program performs the following operations:

1. initialize the simulation/monitoring time  $t$  equal to 0 and close the extensions of the classes  $C$  on which it is necessary to perform closed-world reasoning (by asserting that the class  $KC$  is equivalent to the enumeration of all individuals that can be proved to be members of the class  $C$  retrieved with the `retrieve(C)` query command);
2. insert in the ABox the assertion `happensAt(elapse, t)`;
3. insert in the ABox the events or actions that happen in the system between  $t - 1$  and  $t$  and that are cached in a queue (this involves creating new individuals of the class `Event`);
4. run a reasoner (more specifically, Pellet 2.0) to deduce all assertions that can be inferred (truth values of temporal propositions, states of commitments, etc.);
5. update the closure of the relevant classes;
6. increment the time of simulation  $t$  by 1 and go to the point 2.

After point 5, given that the ontology has reached a stable state it is possible to let agents perform queries about pending, fulfilled, or violated commitments in order to plan their subsequent actions. When the ontology is used for monitoring purposes, and given that internal time (i.e., the time as represented in the ontology) is discrete, it is necessary to wait the actual number of seconds that elapse between two internal instants.

The corresponding Java pseudo code is as follows:

```
t=0
for each class C that has to be closed
  assert KC ≡ {i1,...in} with {i1,...in} = retrieve(C)
while t<timeSimulation {
  assert happensAt(elapse,t)
  for each event en in ActionQueue
    assert happensAt(en,t)
```

```

run Pellet reasoner
for each class  $C$  that has to be closed
  remove equivalent class axioms of class KC
  assert  $KC \equiv \{i_1, \dots, i_n\}$  with  $\{i_1, \dots, i_n\} = \text{retrieve}(C)$ 
run agents queries
t=t+1
}

```

## 4 The ontology of obligations and prohibitions

In this section we present the TBox, the RBox, and part of the ABox that have to be included in the ontology of any interaction system modelled using the OCeAN concepts of temporal proposition, commitment, role, and norm. In particular we specify the classes, the properties and the axioms for modelling those concepts and introduce some SWRL rules to deduce the truth value of temporal propositions. Social commitments are a crucial concept in our approach because they are used to model obligations and prohibitions due either to the activation of norms or created by the performance of communicative acts, like promises. Thanks to their evolution in time, commitments can be used to monitor the behavior of autonomous agents by detecting their *violation* or *fulfilment*, as a precondition for reacting with suitable *passive or active sanctions* or with a *reward* [8].

Some general classes of our ontology are used as domain or range of the properties used to describe temporal propositions and commitments; they are class *Event*, class *Action* and class *Agent*. In particular, an event may have as a property its *time* of occurrence. Class *Action* is a subclass of *Event*, and has a further property used to represent the *actor* of the action. Such properties are defined as follows:

$$\begin{aligned}
&Event \sqcap Agent \sqsubseteq \perp; \quad Action \sqsubseteq Event; \\
&hasActor : Action \rightarrow_O Agent; \\
&happensAt : Event \rightarrow_D integer;
\end{aligned}$$

To represent the elapsing of time we introduce in the ABox the individual *elapse*, that is asserted to be a member of class *Event*:  $Event(elapse)$ .

### 4.1 Temporal propositions

Temporal propositions are used to represent the *content* and *condition* of social commitments. They are a construct used to relate in two different ways a *proposition* to an *interval of time*. In the current OWL specification, we distinguish between *positive temporal propositions* used in commitments to represent obligations (when an action has to be performed within a given interval of time), and *negative temporal propositions* used to model prohibitions (when an action must not be performed during a predefined interval of time).

The classes necessary to model temporal propositions are *TemporalProp*, with the two subclasses *TPPos* and *TPNeg* used to distinguish between positive and

negative temporal propositions. The classes *IsTrue* and *IsFalse* are used to model the truth values of temporal propositions. All this is specified by the following axioms:

$$\begin{aligned} &TemporalProp \sqcap Agent \sqsubseteq \perp; TemporalProp \sqcap Event \sqsubseteq \perp; \\ &TPPos \sqsubseteq TemporalProp; TPNeg \sqsubseteq TemporalProp; \\ &TPPos \sqcap TPNeg \sqsubseteq \perp; \\ &TemporalProp \equiv TPos \sqcup TPNeg; \\ &IsTrue \sqsubseteq TemporalProp; IsFalse \sqsubseteq TemporalProp; \\ &IsTrue \sqcap IsFalse \sqsubseteq \perp; \end{aligned}$$

The class *TemporalProp* is the domain of the following object and data properties:

$$\begin{aligned} &hasAction : TemporalProp \rightarrow_O Action; \\ &hasStart : TemporalProp \rightarrow_D integer; \\ &hasEnd : TemporalProp \rightarrow_D integer; \\ &TemporalProp \sqsubseteq = 1 hasAction \sqcap = 1 hasStart \sqcap = 1 hasEnd \end{aligned}$$

The classes *IsTrue* and *IsFalse* are used to keep track of the truth value of temporal propositions by means of two SWRL rules, that are different on the basis of the type of temporal proposition. A positive temporal proposition (i.e., a member of class *TPPos*) is used to represent an obligation to do something in a given interval of time, with starting points  $t_{start}$  and deadline  $t_{end}$ . We therefore introduce a rule that deduces that the truth value of the temporal proposition is true (i.e., the temporal proposition becomes member of the class *IsTrue*) if the action associated to the temporal proposition is performed between the  $t_{start}$  (inclusive) and the  $t_{end}$  (exclusive) of interval of time associated to the same proposition. In the following SWRL rule we use two built-ins to compare the current time with the interval of time associated to the temporal proposition:

**RuleTPPos1:**

$$\begin{aligned} &happensAt(elapse, ?t) \wedge happensAt(?a, ?t) \wedge TPos(?tp) \wedge hasAction(?tp, ?a) \wedge \\ &hasStart(?tp, ?ts) \wedge hasEnd(?tp, ?te) \wedge swrlb:lessThanOrEqual(?ts, ?t) \wedge \\ &swrlb:lessThan(?t, ?te) \rightarrow IsTrue(?tp) \end{aligned}$$

We then have to define a rule that, when the time  $t_{end}$  of a positive temporal proposition elapses, and such a temporal proposition is not true, deduces that the temporal proposition is member of the class *IsFalse*. Here closed-world reasoning comes into play, because we cannot assume the ABox to contain an explicit assertion that an action has not been performed: rather, we want to deduce that an action has not been performed by the lack of an assertion that it has been performed. Clearly, an SWRL rule like

$$\begin{aligned} &happensAt(elapse, ?te) \wedge hasEnd(?tp, ?te) \wedge TPos(?tp) \wedge (not IsTrue)(?tp) \\ &\rightarrow IsFalse(?tp) \end{aligned}$$

would not work, given that OWL/SWRL reasoners operate under the open world assumption. This means that the conclusion that a temporal proposition is false can only be reached for those propositions that can be definitely proved

not to be members of *IsTrue*. On the contrary, if a temporal proposition is not deduced to be *IsTrue* by RuleTp1, even if its deadline has been reached it will not be deduced to be *IsFalse*.

To solve this problem we first assume that our ABox contains complete information on the performance of actions. This allows us to adopt a closed-world perspective as far as the performance of actions is concerned. More specifically, we assume that the program specified in Section 3 will always update the ABox when an event has happened (i.e. the program can only inset in the ABox the information that an event has happened at current time  $t$ ); we then want to deduce that all temporal propositions, that cannot any longer become true because their deadline has elapsed, are false.

To get this result we need to perform some form of closed world reasoning on class *IsTrue*. As stated in [14] “the DL *ALCK* [5] adds a non-monotonic **K** operator (which is a kind of necessity operator) to the DL *ALC* to provide the ability to “turn on” the Closed World Assumption (CWA) when needed. The reasoning support for *ALCK* language has been implemented in Pellet to answer CWA queries that use the **K** operator”. However, our ontology uses a more expressive DL than *ALC*; moreover, the use of the **K** operator in SWRL rules is not supported.

We therefore take a different approach, based on an explicit *closure* of class *IsTrue*. More precisely, we introduce a new class, *KIsTrue*, which is meant to contain all temporal propositions that, at a given time, are known to be true. Class *KIsTrue* therefore represents, at any given instant, the explicit closure of class *IsTrue*. Given its intended meaning, class *KIsTrue* has to be a subclass of *IsTrue* (and, as a consequence, of *TemporalProp*):

$$KIsTrue \sqsubseteq IsTrue$$

To maintain class *KIsTrue* as the closure of class *IsTrue*, we define it periodically as equivalent to the enumeration of all individuals that can be proved to be members of *IsTrue*. This can be done by the Java program used to update the ABox to keep trace of the elapsing of time (described in Section 3) by executing the operations described in the following pseudo-code:

```
assert KIsTrue ≡ {tp1, ...tpn} with {tp1, ...tpn} = retrieve(IsTrue)
```

We now introduce a new class, *NotKIsTrue*, which is intended to contain all temporal propositions whose deadline is elapsed, and that are not members of *KIsTrue*. Such a class is defined as the difference between the set of all individuals that belong to *TemporalProp*, and the set of all those individuals that are members of *KIsTrue*:

$$NotKIsTrue \equiv TemporalProp \sqcap \neg KIsTrue$$

We are now ready to write a rule to deduce that the truth value of a positive temporal proposition is false if the deadline of the temporal proposition has elapsed, and it is not known that the associated action has been performed:

**RuleTPPos2:**

$$happensAt(elapse, ?te) \wedge hasEnd(?tp, ?te) \wedge TPPos(?tp) \wedge NotKIsTrue(?tp) \rightarrow IsFalse(?tp)$$



We now turn to negative temporal propositions, that is, temporal propositions that are members of the class *TPNeg* and are used to represent the prohibition to do something in a given interval of time. Such propositions belong to class *IsFalse* when the associated action is performed in the interval between  $t_{start}$  (inclusive) and  $t_{end}$  (exclusive). This can be deduced by the following rule:

**RuleTPNeg1:**

$$\begin{aligned} & happensAt(elapse, ?t) \wedge happensAt(?a, ?t) \wedge TPNeg(?tp) \wedge hasAction(?tp, ?a) \wedge \\ & hasStart(?tp, ?ts) \wedge hasEnd(?tp, ?te) \wedge swrlb:lessThanOrEqual(?ts, ?t) \wedge \\ & swrlb:lessThan(?t, ?te) \rightarrow IsFalse(?tp) \end{aligned}$$

Similarly to what we did for RuleTPPos2, we now use the closure of class *IsFalse*, that we call *KIsFalse*, to deduce that a negative temporal proposition *IsTrue* when its  $t_{end}$  has been reached and it has not yet been deduced that the proposition *IsFalse*:

$$\begin{aligned} & KIsFalse \sqsubseteq IsFalse \\ & NotKIsFalse \equiv TemporalProp \sqcap \neg KIsFalse \end{aligned}$$

**RuleTPNeg2:**

$$\begin{aligned} & happensAt(elapse, ?te) \wedge hasEnd(?tp, ?te) \wedge TPNeg(?tp) \wedge NotKIsFalse(?tp) \\ & \rightarrow IsTrue(?tp) \end{aligned}$$

## 4.2 Commitment

In the *OCeAN* meta-model of artificial institutions, commitments are used to model a social relation between a *debtor* a *creditor*, about a certain *content* and under a *condition*. Our idea is that by means of the performance of communicative acts, or due to the activation of norms, certain agents become committed with respect to another agent to perform a certain action within a given deadline (an *obligation*), or not to perform a given action during a given interval of time (a *prohibition*). Such commitments can be conditional on the truth of some proposition. In our model we assume that if an action is neither obligatory nor prohibited, then it is *permitted*.

In order to detect and react to commitment violation and fulfilment we need to deduce a commitments *state* (in our previous works [8] we also introduced precommitments to define the semantics of requests, but this is not relevant in the current work). We introduce in the ontology the class *Commitment*, disjoint from *Event*, *Agent* and *TemporalProp*.

$$\begin{aligned} & Commitment \sqcap Agent \sqsubseteq \perp; Commitment \sqcap Event \sqsubseteq \perp; \\ & Commitment \sqcap TemporalProp \sqsubseteq \perp; \end{aligned}$$

The *Commitment* class is the domain of the following object properties:

$$\begin{aligned} & hasDebtor : Commitment \rightarrow_O Agent; \\ & hasCreditor : Commitment \rightarrow_O Agent; \\ & hasContent : Commitment \rightarrow_O TemporalProp; \\ & hasCondition : Commitment \rightarrow_O TemporalProp; \end{aligned}$$

$hasSource : Commitment \rightarrow_O Norm;$   
 $Commitment \sqsubseteq \exists hasDebtor \sqcap \exists hasCreditor \sqcap =1hasContent \sqcap =1hasCondition;$

The *hasSource* property is used to keep trace of the norm that generated a commitment, as explained in Section 4.3. Obviously the debtor of a commitment has to be the actor of the action to which it is committed, as expressed by the following axiom:

$hasContent \circ hasAction \circ hasActor \sqsubseteq hasDebtor$

In some situations it is necessary to create unconditional commitments. To avoid writing different rules for conditional and for unconditional commitments, we introduce a temporal proposition individual, *tpTrue*, whose truth value is initially true; that is, we assert:  $IsTrue(tpTrue)$ . An unconditional commitment is then defined as a conditional commitment whose condition is *tpTrue*.

Our next problem is deducing whether a given commitment is:

- *pending*, when its condition is satisfied but its content is not known to be *IsTrue* or to be *IsFalse*;
- *fulfilled*, when its content is known to be *IsTrue*;
- *violated*, when its content is known to be *IsFalse* and its condition is known to be *IsTrue*.

Knowing the state of a commitment may be important for the interacting agents to plan their actions on the basis of the advantages of fulfilling certain commitments. We therefore introduce classes *IsPending*, *IsFulfilled*, and *IsViolated*, defined by the following axioms:

$IsFulfilled \sqcap IsViolated \sqsubseteq \perp;$   
 $IsPending \sqsubseteq Commitment; IsFulfilled \sqsubseteq Commitment;$   
 $IsViolated \sqsubseteq Commitment;$

We define the following axiom to deduce that a commitment is member of the class *IsPending*:

**Axiom1:**

$IsPending \equiv (\exists hasContent.NotKIsTrue) \sqcap (\exists hasContent.NotKIsFalse) \sqcap (\exists hasCondition.IsTrue)$

Note that as classes *NotKIsTrue* and *NotKIsFalse* are updated after running the reasoner, as soon as the *content* of a commitment becomes true the commitment is member of both class *IsPending* and class *IsFulfilled*.

Lists of fulfilled and of violated commitments can be obtained by retrieving the individuals that are respectively members of class *IsFulfilled* or *IsViolated*, defined by the following axioms:

**Axiom2:**

$IsFulfilled \equiv \exists hasContent.IsTrue$

**Axiom3:**

$IsViolated \equiv (\exists hasContent.IsFalse) \sqcap (\exists hasCondition.IsTrue)$

### 4.3 Norms and Roles

In OCeAN, norms are introduced to model obligations and prohibitions that, contrary to those created at run time by the performance of communicative acts, are implied by an institutional setting and can be specified at design time. For example, norms can be used to state the rules of an interaction protocol, like the protocol of a specific type of auction, or the rules of a seller-buyer interaction. Given that norms are usually specified at design time, when it is impossible to know which agents will actually interact in the system, one of their distinctive features is that they have to be expressed in term of the *roles* played by the agents. Therefore at run-time, when a norm becomes *active* (i.e., when its activating event happens), the actual debtor and creditor of the obligation or prohibition generated by the norm have to be computed on the basis of the roles played by the agents in the system at that moment.

Another important aspect of norms is that to enforce their fulfillment in an open system, it must be possible to specify *sanctions* or *rewards*. In [7] we suggested that a satisfactory model of sanctions has to distinguish between two different type of actions: the action that the violator of a norm has to perform to extinguish its violation (which we call *active sanction*), and the action that the agent in charge of norm enforcement may perform to deter agents from violating the norm (which we call *passive sanction*). Active sanctions can be represented in our model through a temporal proposition, whereas passive sanctions can be represented as new specific powers that the agent entitled to enforce the norm acquires when a norm is violated. As far as passive sanctions are concerned, another norm (that in [13] is called *enforcement norm*) may oblige the enforcer to punish the violation. Due to space limitations, in this paper we do not model the notion of *power*; thus passive sanctions are not treated in this paper. An obligation or prohibition generated by a norm can in turn violated; it will therefore be necessary to monitor the fulfillment or violation of such obligations or prohibition to punish the violation.

#### Role

Typically, artificial institutions provide for different roles. In a run of an auction, for example, we may have the roles of auctioneer and of participant; in a company, like an auction house, we may have the roles of boss or employee; and so on. More generally, also the debtor and the creditor of a commitment may be regarded as roles. Coherently with these examples, a role is identified by a label (like *auctioneer*, *participant*, etc.) and by the institutional entity that provides for the role. Such an institutional entity may be an organization (like an auction house), an institutional activity (like a run of an auction), or an institutional relationship (like a commitment). For example an agent may be the *auctioneer* of run 01 of a given auction, or an *employee* of IBM, or the *creditor* of a specific commitment.

We introduce class *Role* to represent the set of possible labels that representing roles and class *InstEntity* to represent the institutional entity within which a given role is played. Elements of class *AgentInRole* are used to reify the fact

that an agent plays a given role in a given institutional entity. Those classes are related by the following object properties:

$isPlayedBy : AgentInRole \rightarrow_O Agent;$   
 $hasRole : AgentInRole \rightarrow_O Role;$   
 $isIn : AgentInRole \rightarrow_O InstEntity;$

### Norm

Summarizing, a norm has: a *content* and a *condition*, modelled using temporal propositions; a *debtor* and a *creditor*, expressed in term of roles; an *activating event*; and a collection of *active* and *passive sanctions*. Norms are represented in our ontology using class *Norm* and the following object properties:

$hasRoleDebtor : Norm \rightarrow_O Role;$   $hasRoleCreditor : Norm \rightarrow_O Role;$   
 $hasNContent : Norm \rightarrow_O TemporalProp;$   
 $hasNCondition : Norm \rightarrow_O TemporalProp;$   
 $hasActivation : Norm \rightarrow_O Event;$   
 $hasASanction : Norm \rightarrow_O TemporalProp;$   
 $hasPSanction : Norm \rightarrow_O Power;$

When a norm is activated it is necessary to create as many commitments as there are agents playing the role associated to the debtor property of the norm. For example, the activation of a norm that applies to all the agents playing the role of *participant* of an auction, creates a commitment for each participant currently taking part to the auction. The creditors of these commitments are the agents that play the role reported in the creditor property of the norm. All these commitments have to be related by the *hasSource* object property (defined in Section 4.2) to the norm that generated them; this is important to know which norm generated a commitment and what sanctions apply for the violation of such commitment.

As every commitment is an individual of the ontology, the activation of a norm involves the generation of new individuals. However, the creation of new individuals in an *ABox* cannot be performed using OWL or SWRL. There are at least two possible solutions to this problem, which we plan to investigate in our future work. The first consists in defining a set of axioms in the ontology that allows the reasoner to deduce the existence of those commitments as anonymous objects with certain properties. With this solution, an agent that needs to know its pending commitments instead of simply retrieving the corresponding individuals will have to retrieve their contents, conditions and debtors. Another possible solution consists in defining a new built-in that makes it possible for SWRL rules to create new individuals as members of certain classes and with given properties. A similar problem will have to be solved to manage the creation of a sanctioning commitment generated by the violation of a commitment related to a norm, which has as content the temporal proposition associated to the active sanction of the norm.

## 5 Example

In this section we show how it is possible to specify the state of an interaction system and to simulate or monitor its evolution in time. To do so it is necessary to integrate the ontology defined in the previous sections with an application-dependent ontology, and to insert a set of individuals for representing commitments and temporal propositions in the ABox. In a real application these commitments and their temporal propositions will be created by the performance of communicative acts (defined in the *OCeAN* agent communication library [8]) or by the activation of norms. If the system is used for monitoring purposes, we assume that there is a way of mapping the actions that are actually executed onto their counterparts in the ontology.

Here we describe an example of interaction where a seller agent, Bob, promises to deliver a product (a book) to a buyer agent, Ann, on condition that the buyer agent pays a certain amount of money for the product. We also represent the prohibition for the seller to deliver a different product (a CD). Different possible evolution of the state of the interaction are possible on the basis of the agents' actions.

The ontology described in the previous sections has to be integrated as follows: *pay* and *deliver* are two different types of actions; both of them have a receiver and an object; the *pay* action also has an amount of money. In a more realistic application these concepts would be described in a detailed domain-dependent ontology.

The agents are represented with the following assertions:

*Agent(ann); Agent(bob); ≠ (ann, bob);*

The actions that we are interested to model in the ontology are represented by the following assertions:

*Action(payBook1); Action(deliverBook1); Action(deliverCD1);  
hasActor(payBook1, ann); hasActor(deliverBook1, bob);  
hasActor(deliverCD1, bob);  
≠ (payBook1, deliverBook1, deliverCD1, elapse);*

Temporal propositions are represented by the following assertions:

*TPPos(tpPayBook1); TPPos(tpDeliverBook1); TPNeg(tpNotDeliverCD1);  
hasAction(tpPayBook1, payBook1); hasStart(tpPayBook1, 1);  
hasEnd(tpPayBook1, 3);  
hasAction(tpDeliverBook1, deliverBook1); hasStart(tpDeliverBook1, 1);  
hasEnd(tpDeliverBook1, 2);  
hasAction(tpNotDeliverCD1, deliverCD1); hasStart(tpNotDeliverCD1, 0);  
hasEnd(tpNotDeliverCD1, 3);  
≠ (tpPayBook1, tpDeliverBook1, tpTrue); ≠ (tpNotDeliverCD1, tpTrue);*

Commitments are represented by the following assertions:

*Commitment(c1); Commitment(c2); Commitment(c3);  
hasDebtor(c1, ann); hasCreditor(c1, bob);  
hasContent(c1, tpPayBook1); hasCondition(c1, tpDeliverBook1);*

*hasDebtor(c2, bob); hasCreditor(c2, ann);*  
*hasContent(c2, tpDeliverBook1); hasCondition(c2, tpTrue);*  
*hasDebtor(c3, bob); hasCreditor(c3, ann);*  
*hasContent(c3, tpNotDeliverCD1); hasCondition(c3, tpTrue);*  
 $\neq (c1, c2, c3)$ ;

The history of the system is represented by the following assertions:

*happensAt(deliverBook1, 1)*

We created the ontology of the interaction system with the free, open source ontology editor Protege 4.0 beta<sup>11</sup>. As this version of Protege does not support the editing of SWRL rules, we created them with Protege 3.4 and inserted their RDF/ ML code in the ontology file.

In Table 1 we report the evolution of the ontology ABox in time, with particular regard to the truth value of the temporal propositions and the state of commitments. As the extension of classes *KIsTrue* and *KIsFalse* is computed by an external program, when the reasoner runs their extensions are specified in the axiom relative to the previous state. In the table we abbreviate the assertion *happensAt(elapse, n)* with the expression  $t = n$ .

<i>time</i>	$t = 0$	$t = 1$	$t = 2$	$t = 3$
<i>tpPayBook1</i> [1, 3]				<i>IsFalse</i>
<i>tpDeliverBook1</i> [1, 2]		<i>IsTrue</i>	<i>IsTrue</i>	<i>IsTrue</i>
<i>tpNotDeliverCD1</i> [0, 3]				<i>IsTrue</i>
<i>c1(ann, bob, tpPayBook1, tpDeliverBook1)</i>		<i>IsPending</i>	<i>IsPending</i>	<i>IsViolated</i>
<i>c2(bob, ann, tpDeliverBook1, tpTrue)</i>	<i>IsPending</i>	<i>IsFulfilled</i>	<i>IsFulfilled</i>	<i>IsFulfilled</i>
<i>c3(bob, ann, tpNotDeliverCD1, tpTrue)</i>	<i>IsPending</i>	<i>IsPending</i>	<i>IsPending</i>	<i>IsFulfilled</i>

Classes updated by the external program

<i>KIsTrue</i>	{ <i>tpTrue</i> }	{ <i>tpTrue, tpDeliverBook1</i> }	{ <i>tpTrue, tpDeliverBook1</i> }	{ <i>tpTrue, tpDeliverBook1, tpNotDeliverCD1</i> }
<i>KIsFalse</i>	<i>nothing</i>	<i>nothing</i>	<i>nothing</i>	{ <i>tpPayBook1</i> }

**Table 1.** Dynamic evolution of the state of the system

## 6 Conclusions and Related Works

The main contributions of this paper, with respect to our previous works and with respect to other approaches, are as follows. We show how conditional obligations and prohibitions with stating points and deadlines may be specified and

<sup>11</sup> <http://protege.stanford.edu/>

monitored using OWL and SWRL with significant advantages with respect to other approach that use other formal languages. Moreover we propose a hybrid solution, based on an OWL ontology, SWRL rules, and a Java program, of the problem of monitoring the time evolution of obligations and prohibitions.

In particular if we compare this specification with another one that we presented elsewhere based on Event Calculus [8] we observe significant improvement in performance (even if a complete comparison will be possible only when the complete *OCeAN* meta-model will be formalized with Semantic Web Technology). Moreover semantic web technologies are becoming an international standard for web applications and numerous tools, reasoners, and libraries are available to support the development and usage of ontologies. This is a crucial advantage with respect to other languages used in the multiagent community for the specification of norms and organizations, like as we already mentioned the Event Calculus [15, 1], or other specific formal languages like the one required by the rule engine Jess [11, 4].

In literature there are few approaches that use semantic web languages for the specification of multiagent systems. For example in [12] prohibited, obliged and permitted actions are represented as object properties from agents to actions. But without the reification of the notion of obligation and prohibition that we propose here, it is very difficult to find a feasible solution to express conditional commitments with deadlines. Moreover the approach proposed for detecting violations is based on the external performance of SPARQL queries and on the update of the ABox to register that an obligation/prohibition resulted violated; however SPARQL queries do not exploit the semantics specified by the ontology. In [2] a hybrid approach is presented: they define a communication acts ontology using OWL and express the semantics of those acts through social commitments that are formalized in the Event Calculus. This work is complementary with respect to our approach, in fact we specify also the semantics of social commitments using semantic web technologies. Semantic web technologies in multiagent systems can be used also to specify domain specific ontologies used in the content of norms like in [6]. Another interesting contribution is due also to the exemplification of a solution to the problem to performing closed world reasoning on certain classes in OWL. Another work that tackles a similar problem in a different domain, the ontology of software models, is [3].

Indeed this model is still incomplete e we plan to investigate how it is possible to manage the creation of commitments to model norm activations, and to model active sanctions, moreover we plan to study how to formalize the notion of power to express the semantics of declarative communicative acts and of passive sanctions.

## References

1. A. Artikis, M. Sergot, and J. Pitt. Animated Specifications of Computational Societies. In C. Castelfranchi and W. L. Johnson, editor, *Proceedings of the 1st International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS 2002)*, pages 535–542. ACM Press, 2002.

2. I. Berges, J. Bermdez, A. Goi, and A. Illarramendi. Semantic web technology for agent communication protocols. In *The Semantic Web: Research and Applications 5th European Semantic Web Conference, ESWC 2008, Tenerife, Canary Islands, Spain, June 1-5, 2008 Proceedings*, pages 5–18, 2008.
3. M. Bräuer and H. Lochmann. An ontology for software models and its practical implications for semantic web reasoning. In S. Bechhofer, M. Hauswirth, J. Hoffmann, and M. Koubarakis, editors, *ESWC*, volume 5021 of *LNCS*, pages 34–48. Springer, 2008.
4. V. T. da Silva. From the specification to the implementation of norms: an automatic approach to generate rules from norms to govern the behavior of agents. *Autonomous Agents and Multi-Agent Systems*, 17(1):113–155, August 2008.
5. F. M. Donini, M. Lenzerini, D. Nardi, A. Schaerf, and W. Nutt. An epistemic operator for description logics. *Artificial Intelligence*, 200(1-2):225274, 1998.
6. C. Felicissimo, J.-P. Briot, C. Chopinaud, and C. Lucena. How to concretize norms in NMAS? An operational normative approach presented with a case study from the television domain. In *International Workshop on Coordination, Organization, Institutions and Norms in Agent Systems (COIN@AAAI'08)*, 23rd AAAI Conference on Artificial Intelligence, Chicago, IL, Etats-Unis, 2008. AAAI, AAAI Press.
7. N. Fornara and M. Colombetti. Specifying and enforcing norms in artificial institutions. In M. Baldoni, T. Son, B. van Riemsdijk, and M. Winikoff, editors, *Declarative Agent Languages and Technologies VI 6th International Workshop, DALT 2008, Estoril, Portugal, May 12, 2008, Revised Selected and Invited Papers*, volume 5397 of *LNCS*, pages 1–17. Springer Berlin / Heidelberg, 2009.
8. N. Fornara and M. Colombetti. *Specifying Artificial Institutions in the Event Calculus*, chapter XIV, page to appear. Information science reference. IGI Global, 2009.
9. N. Fornara, F. Viganò, and M. Colombetti. Agent communication and artificial institutions. *Autonomous Agents and Multi-Agent Systems*, 14(2):121–142, April 2007.
10. N. Fornara, F. Viganò, M. Verdicchio, and M. Colombetti. Artificial institutions: A model of institutional reality for open multiagent systems. *Artificial Intelligence and Law*, 16(1):89–105, March 2008.
11. A. García-Camino, J. A. Rodríguez-Aguilar, C. Sierra, and W. Vasconcelos. Constraint rule-based programming of norms for electronic institutions. *Autonomous Agents and Multi-Agent Systems*, 18(1):186–217, 2009.
12. J. S.-C. Lam, F. Guerin, W. Vasconcelos, and T. J. Norman. Representing and reasoning about norm-governed organisations with semantic web languages. In *Sixth European Workshop on Multi-Agent Systems Bath, UK, 18-19 December 2008*, 2008.
13. F. López y López, M. Luck, and M. d’Inverno. A Normative Framework for Agent-Based Systems. In *Proceedings of the First International Symposium on Normative Multi-Agent Systems, Hatfield*, 2005.
14. E. Sirin, B. Parsia, B. C. Grau, A. Kalyanpur, and Y. Katz. Pellet: A practical owl-dl reasoner. *Web Semantics: Science, Services and Agents on the World Wide Web*, 5(2):51–53, 2007.
15. P. Yolum and M. Singh. Reasoning about commitment in the event calculus: An approach for specifying and executing protocols. *Annals of Mathematics and Artificial Intelligence*, 42:227–253, 2004.



# Prioritized Goals and Subgoals in a Logical Account of Goal Change – A Preliminary Report

Shakil M. Khan and Yves Lespérance

Department of Computer Science and Engineering  
York University, Toronto, ON, Canada M3J 1P3  
{skhan, lesperan}@cse.yorku.ca

**Abstract.** Most previous logical accounts of goal change do not deal with prioritized goals and do not handle subgoals and their dynamics properly. Many are restricted to achievement goals. In this paper, we develop a logical account of goal change that addresses these deficiencies. In our account, we do not drop lower priority goals permanently when they become inconsistent with other goals and the agent’s knowledge; rather, we make such goals inactive. We ensure that the agent’s chosen goals/intentions are consistent with each other and the agent’s knowledge. When the world changes, the agent recomputes her chosen goals and some inactive goals may become active again. This ensures that our agent maximizes her utility. We also propose an approach for handling subgoals and their dynamics. We prove that the proposed account has some intuitively desirable properties.

## 1 Introduction

There has been much work on modeling agent’s mental states, beliefs, goals, and intentions, and how they interact and lead to rational decisions about action. As well, there has been a lot of work on modeling belief change. But motivational attitudes and their dynamics have received much less attention. Most formal models of goals and goal change [1–7] assume that all goals are equally important and many only deal with achievement goals (one exception to this is the model of prioritized goals in [7]). Moreover, most of these frameworks do not guarantee that an agent’s goals will properly evolve when an action is performed or an event occurs, e.g. when the agent’s beliefs/knowledge changes or a goal is adopted or dropped. Also, they do not model the dependencies between goals and the subgoals and plans adopted to achieve these goals. For instance, subgoals and plans adopted to bring about a goal should be dropped when the parent goal becomes impossible, is achieved, or is dropped. Dealing with these issues is important for developing effective models of rational agency. It is also important for work on BDI agent programming languages, where handling declarative goals is an active research topic [8, 9].

In this paper, we present a formal model of prioritized goals and their dynamics that addresses some of these issues. Specifically, we propose a framework, where an agent can have multiple goals at different priority levels, possibly inconsistent with each other. We define intentions as the maximal set of highest priority goals that is consistent given the agent’s knowledge. Our formalization of goals and goal dynamics ensures that the

agent strives to maximize her utility. Our model supports the specification of general temporally extended goals, not just achievement goals, and handles subgoals and their dynamics.

We start with a (possibly inconsistent) initial set of *prioritized goals* or desires that are totally ordered according to priority, and specify how these goals evolve when actions/events occur and the agent’s knowledge changes. We define the agent’s *chosen goals* or intentions in terms of this goal hierarchy. Our agents maximize their utility; they will abandon a chosen goal  $\phi$  if an opportunity to commit to a higher priority but inconsistent with  $\phi$  goal arises. To this end, we keep all prioritized goals in the goal base unless they are explicitly dropped. At every step, we compute an optimal set of chosen goals given the hierarchy of prioritized goals, preferring higher priority goals, such that chosen goals are consistent with each other and with the agent’s knowledge. Thus at any given time, some goals in the hierarchy are active, i.e. chosen, while others are inactive. Some of these inactive goals may later become active, e.g. if a higher priority active goal that is inconsistent with the inactive goal becomes impossible. We also show how the dependencies between goals and subgoals can be modeled. Finally, we prove some interesting properties about the dynamics of chosen goals.

As mentioned above, our formalization of prioritized goals ensures that the agent always tries to maximize her utility, and as such a limitation of our framework is that it displays an idealized form of rationality. In Section 5, we discuss how this relates to Bratman’s theory of practical reasoning [10]. We use an action theory based on the situation calculus [11] along with our formalization of paths in the situation calculus as our base formalism.

The paper is organized as follows: in the next section, we outline our base framework. In Section 3, we formalize *paths* in the situation calculus to support modeling goals. In Section 4, we present our model of prioritized goals. In section 5, we present our formalization of goal dynamics and discuss some of its properties. In Section 6, we discuss what it means for an agent to have a subgoal and how subgoals change as a result of changes to their parent goals. Then in the last section, we summarize our results, discuss previous work in this area, and point to possible future work.

## 2 Action and Knowledge

Our base framework for modeling goal change is the situation calculus [11] as formalized in [12]. In the situation calculus, a possible state of the domain is represented by a situation. There is a set of initial situations corresponding to the ways the agents believe the domain might be initially, i.e. situations in which no actions have yet occurred.  $\text{Init}(s)$  means that  $s$  is an initial situation. The actual initial state is represented by a special constant  $S_0$ . There is a distinguished binary function symbol *do* where  $\text{do}(a, s)$  denotes the successor situation to  $s$  resulting from performing the action  $a$ . Thus the situations can be viewed as a set of trees, where the root of each tree is an initial situation and the arcs represent actions. Relations (and functions) whose truth values vary from situation to situation, are called relational (functional, resp.) fluents, and are denoted by predicate (function, resp.) symbols taking a situation term as their last argument. There is a special predicate  $\text{Poss}(a, s)$  used to state that action  $a$  is executable in situation  $s$ .

Our framework uses a theory  $D_{basic}$  that includes the following set of axioms:<sup>1</sup> (1) action precondition axioms, one per action  $a$  characterizing  $Poss(a, s)$ , (2) successor state axioms (SSA), one per fluent, that succinctly encode both effect and frame axioms and specify exactly when the fluent changes [12], (3) initial state axioms describing what is true initially including the mental states of the agents, (4) unique name axioms for actions, and (5) domain-independent foundational axioms describing the structure of situations [14].

Following [15, 16], we model knowledge using a possible worlds account adapted to the situation calculus.  $K(s', s)$  is used to denote that in situation  $s$ , the agent thinks that she could be in situation  $s'$ . Using  $K$ , the knowledge of an agent is defined as:<sup>2</sup>  $Know(\Phi, s) \stackrel{\text{def}}{=} \forall s'. K(s', s) \supset \Phi(s')$ , i.e. the agent knows  $\Phi$  in  $s$  if  $\Phi$  holds in all of her  $K$ -accessible situations in  $s$ .  $K$  is constrained to be reflexive, transitive, and Euclidean in the initial situation to capture the fact that agents' knowledge is true, and that agents have positive and negative introspection. As shown in [16], these constraints then continue to hold after any sequence of actions since they are preserved by the successor state axiom for  $K$ . We also assume that all actions are public, i.e. whenever an action (including exogenous events) occurs, the agent learns that it has happened. Note that, we work with knowledge rather than belief. Although much of our formalization should extend to the latter, we leave this for future work.

### 3 Paths in the Situation Calculus

To support modeling temporally extended goals, we introduce a new sort of *paths*, with (possibly sub/super-scripted) variables  $p$  ranging over paths. A path is essentially an infinite sequence of situations, where each successor situation along the path can be reached by performing some *executable* action in the preceding situation. We introduce a predicate  $OnPath(p, s)$ , meaning that the situation  $s$  is on the path  $p$ . Also,  $Starts(p, s)$  means that  $s$  is the starting situation of path  $p$ . A path  $p$  starts with the situation  $s$  iff  $s$  is the earliest situation on  $p$ .<sup>3</sup>

**Axiom 1**  $Starts(p, s) \equiv OnPath(p, s) \wedge \forall s'. OnPath(p, s') \supset s \leq s'$ .

In the standard situation calculus, paths are implicitly there, and a path can be viewed as a pair  $(s, F)$  consisting of a situation  $s$  representing the starting situation of the path, and a function  $F$  from situations to actions (called *Action Selection Functions* (ASF) or strategies in [5]), such that from starting situation  $s$ ,  $F$  defines an infinite

<sup>1</sup> We will be quantifying over formulae, and thus we assume  $D_{basic}$  includes axioms for encoding of formulae as first order terms, as in [13]. We will also be using lists of integers, and assume that  $D_{basic}$  includes axiomatizations of integers and lists.

<sup>2</sup> A state formula  $\Phi(s)$  takes a single situation as argument and is evaluated with respect to that situation.  $\Phi$  may contain a placeholder constant *now* that stands for the situation in which  $\Phi$  must hold.  $\Phi(s)$  is the formula that results from replacing *now* by  $s$ . Where the intended meaning is clear, we sometimes suppress the placeholder.

<sup>3</sup> In the following,  $s < s'$  means that  $s'$  can be reached from  $s$  by performing a sequence of executable actions.  $s \leq s'$  is an abbreviation for  $s < s' \vee s = s'$ .

sequence of situations by specifying an action for every situation starting from  $s$ . Thus, one way of axiomatizing paths is by making them correspond to such pairs  $(s, F)$ :

**Axiom 2**

$$\begin{aligned} \forall p. \text{Starts}(p, s) \supset (\exists F. \text{Executable}(F, s) \wedge \forall s'. \text{OnPath}(p, s') \equiv \text{OnPathASF}(F, s, s')), \\ \forall F, s. \text{Executable}(F, s) \supset \exists p. \text{Starts}(p, s) \wedge \forall s'. \text{OnPathASF}(F, s, s') \equiv \text{OnPath}(p, s'). \end{aligned}$$

This says that for every path there is an executable ASF that produces exactly the sequence of situations on the path from its starting situation. Also, for every executable ASF and situation, there is a path that corresponds to the sequence of situations produced by the ASF starting from that situation.

$$\begin{aligned} \text{OnPathASF}(F, s, s') &\stackrel{\text{def}}{=} s \leq s' \wedge \forall a, s^*. s < do(a, s^*) \leq s' \supset F(s^*) = a, \\ \text{Executable}(F, s) &\stackrel{\text{def}}{=} \forall s'. \text{OnPathASF}(F, s, s') \supset \text{Poss}(F(s'), s'). \end{aligned}$$

Here,  $\text{OnPathASF}(F, s, s')$  [6] means that the situation sequence defined by  $(s, F)$  includes the situation  $s'$ . Also, the situation sequence encoded by a strategy  $F$  and a starting situation  $s$  is executable iff for all situations  $s'$  in this sequence, the action selected by  $F$  in  $s'$  is executable in  $s'$ .

In our framework, we will use both state and path formulae. A state formula is a formula that has a free situation variable in it, whereas a path formula is one that has a free path variable. State formulae are used in the context of knowledge while path formulae are used in that of goals. We use  $\Phi(s), \Psi(s), \dots$  and  $\phi(p), \psi(p), \dots$  possibly with decorations to represent state and path formulae, respectively. Note that, by incorporating infinite paths in our framework, we can evaluate goals over these and handle arbitrary temporally extended goals; thus, unlike some other situation calculus based accounts where goal formulae are evaluated w.r.t. finite paths (e.g. [7]), we can handle for example unbounded maintenance goals.

We next define some useful constructs. A state formula  $\Phi$  *eventually holds* over the path  $p$  if  $\Phi$  holds in some situation that is on  $p$ , i.e.  $\diamond\Phi(p) \stackrel{\text{def}}{=} \exists s'. \text{OnPath}(p, s') \wedge \Phi(s')$ . Other Temporal Logic operators can be defined similarly, e.g. always  $\Phi$ :  $\Box\Phi(p)$ .

An agent *knows* in  $s$  that  $\phi$  has become *inevitable* if  $\phi$  holds over all paths that start with a  $K$ -accessible situation in  $s$ , i.e.  $\text{KIinevitable}(\phi, s) \stackrel{\text{def}}{=} \forall p. \text{Starts}(p, s') \wedge K(s', s) \supset \phi(p)$ . An agent knows in  $s$  that  $\phi$  is impossible if she knows that  $\neg\phi$  is inevitable in  $s$ , i.e.  $\text{KImpossible}(\phi, s) \stackrel{\text{def}}{=} \text{KIinevitable}(\neg\phi, s)$ .

Thirdly, we define what it means for a path  $p'$  to be a suffix of another path  $p$  w.r.t. a situation  $s$ :

$$\begin{aligned} \text{Suffix}(p', p, s) &\stackrel{\text{def}}{=} \text{OnPath}(p, s) \wedge \text{Starts}(p', s) \\ &\wedge \forall s'. s' \geq s \supset \text{OnPath}(p, s') \equiv \text{OnPath}(p', s'). \end{aligned}$$

Finally,  $\text{SameHistory}(s_1, s_2)$  means that the situations  $s_1$  and  $s_2$  share the same history of actions, but perhaps starting from different initial situations:

**Axiom 3**  $\text{SameHistory}(s_1, s_2) \equiv \text{Init}(s_1) \wedge \text{Init}(s_2)$

$$\vee (\exists a, s'_1, s'_2. s_1 = do(a, s'_1) \wedge s_2 = do(a, s'_2) \wedge \text{SameHistory}(s'_1, s'_2)).$$

## 4 Prioritized Goals

Most work on formalizing goals only deals with static goal semantics and not their dynamics. There are two main categories of motivational attitudes, namely goal [1, 17] (AKA choice [2], wish [10] or preference), and intention. While goals are sometimes allowed to be inconsistent [10], intentions are mostly required to be consistent. Another difference is that agents are committed to their intentions, but not necessarily to their goals [10]. Intention is sometimes primitive [17, 3] and sometimes a defined concept, specified in terms of goals [1, 2, 4]. In this section, we formalize goals or desires with different priorities, which we call *prioritized goals* (p-goals, henceforth). These p-goals are not required to be mutually consistent and need not be actively pursued by the agent. Using this, we define the consistent set of *chosen goals* or intentions (c-goals, henceforth) that the agent is committed to. In the next section, we formalize goal dynamics by providing a SSA for p-goals. The agent's c-goals are automatically updated when her p-goals change. We deal with subgoals and their dynamics in Section 6.

Not all of the agent's goals are equally important to her. Thus, it is useful to support a priority ordering over goals. This information can be used to decide which of the agent's c-goals should no longer be actively pursued in case they become mutually inconsistent. Following [6], we specify each p-goal by its own accessibility relation/fluents  $G$ . A path  $p$  is  $G$ -accessible at priority level  $n$  in situation  $s$  (denoted by  $G(p, n, s)$ ) if all the goals of the agent at level  $n$  are satisfied over this path and if it starts with a situation that has the same history (in terms of the actions performed so far) as  $s$ . The latter requirement ensures that the agent's p-goal-accessible paths reflect the actions that have been performed so far. A smaller  $n$  represents higher priority, and the highest priority level is 0. Thus in this framework, we assume that the set of p-goals are totally ordered according to priority. We say that an agent has the p-goal that  $\phi$  at level  $n$  in situation  $s$  iff  $\phi$  holds over all paths that are  $G$ -accessible at  $n$  in  $s$ :

$$\text{PGoal}(\phi, n, s) \stackrel{\text{def}}{=} \forall p. G(p, n, s) \supset \phi(p).$$

To be able to refer to all the p-goals of the agent at some given priority level, we also define *only p-goals*.

$$\begin{aligned} \text{OPGoal}(\phi, n, s) &\stackrel{\text{def}}{=} \text{PGoal}(\phi, n, s) \\ &\wedge \forall p, s'. \text{Starts}(p, s') \wedge \text{SameHist}(s', s) \wedge \phi(p) \supset G(p, n, s). \end{aligned}$$

An agent has the *only p-goal* that  $\phi$  at level  $n$  in situation  $s$  iff  $\phi$  is a p-goal at  $n$  in  $s$ , and any path over which  $\phi$  holds and that starts with a situation that has the same action history as in  $s$  is  $G$ -accessible at  $n$  in  $s$ .

A domain theory for our framework  $D$  includes the axioms of a theory  $D_{\text{basic}}$  as in the previous section, the axiomatization of paths i.e. axioms 1-3, domain dependent initial goal axioms (see below), the domain independent axioms 4-7 and the definitions that appear in this section and the next. The modeler must provide initial goal axioms of the following form:

$$\begin{aligned} \text{INITIAL GOAL AXIOMS (a)} \quad &\text{Init}(s) \supset ((G(p, 0, s) \equiv \text{Starts}(p, s') \wedge \text{Init}(s') \wedge \phi_0(p)) \\ &\wedge (G(p, 1, s) \equiv \text{Starts}(p, s') \wedge \text{Init}(s') \wedge \phi_1(p)) \wedge \dots \end{aligned}$$

- $\wedge (G(p, k - 1, s) \equiv \text{Starts}(p, s') \wedge \text{Init}(s') \wedge \phi_{k-1}(p))$ ,
- (b)  $\forall n, p, s. \text{Init}(s) \wedge n \geq k \supset (G(p, n, s) \equiv \text{Starts}(p, s') \wedge \text{Init}(s'))$ ,
- (c)  $\text{Init}(s) \supset \text{NPGoals}(s) = k$ .

The p-goals  $\phi_0, \phi_1, \dots, \phi_{k-1}$  (from highest to lowest priority) of the agent in the initial situations are specified by the Initial Goal Axiom (a); each of them defines a set of initial goal paths for a given priority level, and must be consistent. We assume that the agent has a finite number  $k$  of initial p-goals. For  $n \geq k$ , we make  $G(p, n, s)$  true for every path  $p$  that starts with an initial situation in (b). Thus at levels  $n \geq k$ , the agent has the trivial p-goal that she be in an initial situation. We also have a distinguished functional fluent  $\text{NPGoals}(s)$  that represents the number of prioritized goals that the agent has (i.e. the location of the first empty slot after the last p-goal). Initially  $\text{NPGoals}$  is set to  $k$  in (c). Later, we will specify the dynamics of p-goals by giving SSAs for  $G$  and  $\text{NPGoals}$ .

We use the following as a running example. We have an agent who initially has the following three p-goals:  $\phi_0 = \Box\text{BeRich}$ ,  $\phi_1 = \Diamond\text{GetPhD}$ , and  $\phi_2 = \Box\text{BeHappy}$  at level 0, 1, and 2, respectively (see second column of Table 1). Assume that while all

G-Level	$S_0, S'_1$	$S_1$	$S_2$	$S_3$
4	TRUE	TRUE	$\Box\text{BeRich} \wedge \Box\text{WorkHard} \wedge \Box\text{BeEnergetic}$	TRUE
3	TRUE	$\Box\text{BeRich} \wedge \Box\text{WorkHard}$	$\Box\text{BeRich} \wedge \Box\text{WorkHard}$	TRUE
2	$\Box\text{BeHappy}$	$\Box\text{BeHappy}$	$\Box\text{BeHappy}$	$\Box\text{BeHappy}$
1	$\Diamond\text{GetPhD}$	$\Diamond\text{GetPhD}$	$\Diamond\text{GetPhD}$	$\Diamond\text{GetPhD}$
0	$\Box\text{BeRich}$	$\Box\text{BeRich}$	$\Box\text{BeRich}$	$\Box\text{BeRich}$

**Table 1.** Example of an Agent's PGoals and their Dynamics

of her p-goals are individually achievable initially, her p-goal  $\Diamond\text{GetPhD}$  is inconsistent with her highest priority p-goal  $\Box\text{BeRich}$  as well as with her p-goal  $\Box\text{BeHappy}$ , while the latter are consistent with each other. It is straightforward to specify a domain action theory such that it entails this. Also assume that, after the action *goBankrupt* happens in  $S_0$ , the p-goal  $\Box\text{BeRich}$  becomes impossible. Thus in our example, we have  $\text{OPGoal}(\phi_i(p) \wedge \text{Starts}(p, s) \wedge \text{Init}(s), i, S_0)$ , for  $i = 0, 1, 2$ . Also, for any  $n \geq 3$ , we have  $\text{OPGoal}(\text{Starts}(p, s) \wedge \text{Init}(s), n, S_0)$ .

Using p-goals, we next define c-goals. While p-goals or desires are allowed to be known to be impossible to achieve, an agent's c-goals or intentions must be realistic. Not all of the  $G$ -accessible paths are realistic in the sense that they start with a knowledge accessible situation. To filter these out, we define *realistic* p-goal accessible paths:

$$G_R(p, n, s) \stackrel{\text{def}}{=} G(p, n, s) \wedge \text{Starts}(p, s') \wedge K(s', s),$$

i.e., a path  $p$  is  $G_R$ -accessible at level  $n$  in situation  $s$  if it is  $G$ -accessible at  $n$  in  $s$ , and if  $p$  starts with a situation that is  $K$ -accessible in  $s$ . Thus  $G_R$  prunes out the paths from  $G$  that are known to be impossible. We will define c-goals in terms of realistic p-goals, so this ensures that agents' c-goals are realistic.

The idea of how we compute c-goal-accessible paths is as follows: the set of  $G_R$ -accessibility relations represents a set of prioritized temporal propositions that are can-

didates for the agent's c-goals. Given  $G_R$ , in each situation we want to compute the agent's c-goals such that it is the *maximal consistent* set of higher priority realistic p-goals. We do this iteratively starting with the set of all possible paths (i.e. paths that starts with a  $K$ -accessible situation). At each iteration we compute the intersection of this set with the next highest priority set of  $G_R$ -accessible paths. If the intersection is not empty, we thus obtain a new chosen set of paths at level  $i$ . We call a p-goal chosen by this process an *active* p-goal. If on the other hand, the intersection is empty, then it must be the case that the p-goal represented by this level is either in conflict with another active higher priority p-goal/a combination of two or more active higher priority p-goals, or known to be impossible. In that case, that p-goal is ignored (i.e. marked as inactive), and the chosen set of paths at level  $i$  is the same as at level  $i - 1$ . We repeat this until we reach  $i = NPGoals$ . Axiom 4 computes this intersection:<sup>4</sup>

**Axiom 4**  $G_{\cap}(p, n, s) \equiv$  **if**  $(n = 0)$  **then**  $Starts(p, s') \wedge K(s', s)$   
**else if**  $\exists p'.(G_R(p', n - 1, s) \wedge G_{\cap}(p', n - 1, s))$   
**then**  $(G_R(p, n - 1, s) \wedge G_{\cap}(p, n - 1, s))$   
**else**  $G_{\cap}(p, n - 1, s)$ .

C-goal accessible paths are the result of this intersection after all priority levels have been considered:

$$G_C(p, s) \stackrel{\text{def}}{=} G_{\cap}(p, NPGoals(s), s).$$

We define an agent's c-goals in terms of the  $G_C$ -accessible paths:

$$CGoal(\phi, s) \stackrel{\text{def}}{=} \forall p. G_C(p, s) \supset \phi(p),$$

i.e., the agent has the c-goal that  $\phi$  if  $\phi$  holds over all of her  $G_C$ -accessible paths.

We also define what it means for an agent to have a c-goal at some level  $n$ :

$$CGoal(\phi, n, s) \stackrel{\text{def}}{=} \forall p. G_{\cap}(p, n + 1, s) \supset \phi(p),$$

i.e. an agent has the c-goal at level  $n$  that  $\phi$  if  $\phi$  holds over all paths that are in the prioritized intersection of the set of  $G_R$ -accessible paths up to level  $n$ .

In our example, the agent's p-goals are  $\Box BeRich$ ,  $\Diamond GetPhD$ , and  $\Box BeHappy$  in order of priority. The  $G_{\cap}$ -accessible paths at level 1 in  $S_0$  are the ones that start with a  $K$ -accessible situation and where  $\Box BeRich$  holds. The  $G_{\cap}$ -accessible paths at level 2 in  $S_0$  are the same as at level 1, since there are no  $K$ -accessible paths over which both  $\Diamond GetPhD$  and  $\Box BeRich$  hold. Finally, the  $G_{\cap}$ -accessible paths at level 3 in  $S_0$  and hence the  $G_C$ -accessible paths are those that start with a  $K$ -accessible situation and over which  $\Box BeRich \wedge \Box BeHappy$  holds. Also, it can be shown that initially our example agent has the c-goals that  $\Box BeRich$  and  $\Box BeHappy$ , but not  $\Diamond GetPhD$ .

Note that by our definition of c-goals, the agent can have a c-goal that  $\phi$  in situation  $s$  for various reasons: 1)  $\phi$  is known to be inevitable in  $s$ ; 2)  $\phi$  is an active p-goal at some level  $n$  in  $s$ ; 3)  $\phi$  is a consequence of two or more active p-goals at different levels in  $s$ . To be able to refer to c-goals for which the agent has a primitive motivation, i.e. c-goals that result from a single active p-goal at some priority level  $n$ , in contrast to those

<sup>4</sup> **if**  $\phi$  **then**  $\delta_1$  **else**  $\delta_2$  is an abbreviation for  $(\phi \supset \delta_1) \wedge (\neg\phi \supset \delta_2)$ .

that hold as a consequence of two or more active p-goals at different priority levels, we define *primary* c-goals:

$$\text{PrimCGoal}(\phi, s) \stackrel{\text{def}}{=} \exists n. \text{PGoal}(\phi, n, s) \wedge \exists p. G_C(p, s) \wedge G(p, n, s).$$

That is, an agent has the primary c-goal that  $\phi$  in situation  $s$ , if  $\phi$  is a p-goal at some level  $n$  in  $s$ , and if there is a  $G_C$ -accessible path  $p$  in  $s$  that is also  $G$ -accessible at  $n$  in  $s$ . Thus if an agent has a primary c-goal that  $\phi$ , then she also has the c-goal that  $\phi$ , but not necessarily vice-versa. It can be shown that initially our example agent has the primary c-goals that  $\Box\text{BeRich}$  and  $\Box\text{BeHappy}$ , but not their conjunction.

## 5 Goal Dynamics

An agent's goals change when her knowledge changes as a result of the occurrence of an action (including exogenous events), or when she adopts or drops a goal. We formalize this by specifying how p-goals change. C-goals are then calculated using (realistic) p-goals in every new situation as explained above.

We introduce two actions for adopting and dropping a p-goal,  $\text{adopt}(\phi)$  and  $\text{drop}(\phi)$ , and a third for adopting a subgoal  $\psi$  w.r.t. a supergoal  $\phi$ ,  $\text{adopt}(\psi, \phi)$ . The action precondition axioms for these are as follows:

**Axiom 5**  $\text{Poss}(\text{adopt}(\phi), s) \equiv \neg \exists n. \text{PGoal}(\phi, n, s),$

$$\text{Poss}(\text{adopt}(\psi, \phi), s) \equiv \neg \exists n. \text{PGoal}(\psi, n, s) \wedge \exists n'. \text{PGoal}(\phi, n', s),$$

$$\text{Poss}(\text{drop}(\phi), s) \equiv \exists n. \text{PGoal}(\phi, n, s).$$

That is, an agent can adopt the p-goal that  $\phi$ , if she does not already have  $\phi$  as her p-goal at some level. An agent can adopt a subgoal  $\psi$  w.r.t. the parent goal that  $\phi$  if she does not already have the p-goal that  $\psi$  at some level, and if at some level she currently has the parent goal that  $\phi$ . The  $\text{drop}(\phi)$  action is possible in  $s$  if  $\phi$  is a p-goal at some level  $n$  in  $s$ .

In the following, we specify the dynamics of p-goals by giving the SSA for  $G$  and then discuss each case one at a time:

**Axiom 6 (SSA for G)**  $G(p, n, \text{do}(a, s)) \equiv$

$$\forall \phi, \psi. (a \neq \text{adopt}(\phi) \wedge a \neq \text{adopt}(\psi, \phi) \wedge a \neq \text{drop}(\phi) \wedge \text{Progressed}(p, n, a, s))$$

$$\vee \exists \phi. (a = \text{adopt}(\phi) \wedge \text{Adopted}(p, n, a, s, \phi))$$

$$\vee \exists \phi, \psi. (a = \text{adopt}(\psi, \phi) \wedge \text{SubGoalAdopted}(p, n, a, s, \psi, \phi))$$

$$\vee \exists \phi. (a = \text{drop}(\phi) \wedge \text{Dropped}(p, n, a, s, \phi)).$$

The overall idea of the SSA for  $G$  is as follows. First of all, to handle the occurrence of a non-adopt/drop (i.e. regular) action  $a$ , we progress all  $G$ -accessible paths to reflect the fact that this action has just happened; this is done using the  $\text{Progressed}(p, n, a, s)$  construct, which replaces each  $G$ -accessible path  $p'$  with starting situation  $s'$ , by its suffix  $p$  provided that it starts with  $\text{do}(a, s')$ :

$$\text{Progressed}(p, n, a, s) \stackrel{\text{def}}{=} \exists p'. G(p', n, s) \wedge \text{Starts}(p', s') \wedge \text{Suffix}(p, p', \text{do}(a, s')).$$



Any path over which the next action performed is not  $a$  is eliminated from the respective  $G$  accessibility level.

Secondly, to handle adoption of a p-goal  $\phi$ , we add a new proposition containing the p-goal to the agent's goal hierarchy. We assume that the newly adopted p-goal  $\phi$  has the lowest priority. Thus in addition to progressing the  $G$ -accessible paths at all levels as above, we eliminate the paths over which  $\phi$  does not hold from the  $NPGoals(s)$ -th  $G$ -accessibility level, and the agent acquires the p-goal that  $\phi$  at level  $NPGoals(s)$ .

$$\text{Adopted}(p, n, a, s, \phi) \stackrel{\text{def}}{=} \text{if } (n = NPGoals(s)) \text{ then } (\text{Progressed}(p, n, a, s) \wedge \phi(p)) \\ \text{else } \text{Progressed}(p, n, a, s).$$

The third case of subgoal adoption is discussed in the next section.

Finally, to handle dropping of a p-goal  $\phi$ , we replace the propositions that imply the dropped goal in the agent's goal hierarchy by the "trivial" proposition that the history of actions in the current situation has occurred. Thus in addition to progressing all  $G$ -accessible paths as above, we add back all paths that share the same history with  $do(a, s)$  to the existing  $G$ -accessibility levels where the agent has the p-goal that  $\phi$ , and thus these  $G$ -accessibility levels now amounts to the "trivial" p-goal that  $\text{CorrectHist}(s, path)$ .<sup>5</sup>

$$\text{Dropped}(p, n, a, s, \phi) \stackrel{\text{def}}{=} \text{if } \text{PGoal}(\phi, n, s) \text{ then } \text{Starts}(p, s') \wedge \text{SameHistory}(s', do(a, s)) \\ \text{else } \text{Progressed}(p, n, a, s).$$

The SSA for  $NPGoals(s)$  is as follows:

**Axiom 7 (SSA for NPGoals)**  $NPGoals(do(a, s)) = k \equiv$   
 $a \neq \text{adopt}(\phi) \wedge a \neq \text{adopt}(\psi, \phi) \wedge NPGoals(s) = k \vee$   
 $a = \text{adopt}(\phi) \wedge NPGoals(s) + 1 = k \vee$   
 $a = \text{adopt}(\psi, \phi) \wedge \text{AdjustSubGoalAdopt}(\phi, s) = k.$

That is, when the agent adopts a p-goal, her current  $NPGoals$  is incremented by one. We discuss the adjustment of  $NPGoals$  required for subgoal adoption in the next section. Finally,  $NPGoals$  is not affected by any other action.

Returning to our example, recall that our agent has the c-goals/active p-goals in  $S_0$  that  $\Box\text{BeRich}$  and  $\Box\text{BeHappy}$ , but not  $\Diamond\text{GetPhD}$ , since the latter is inconsistent with her higher priority p-goal  $\Box\text{BeRich}$ . On the other hand, in  $S'_1 = do(goBankrupt, S_0)$ , the agent has the c-goal that  $\Diamond\text{GetPhD}$ , but not  $\Box\text{BeRich}$  nor  $\Box\text{BeHappy}$ ;  $\Box\text{BeRich}$  is excluded from the set of c-goals since it has become impossible to achieve (i.e. unrealistic). Also, since her higher priority p-goal  $\Diamond\text{GetPhD}$  is inconsistent with her p-goal  $\Box\text{BeHappy}$ , the agent will make  $\Box\text{BeHappy}$  inactive.

Note that, while it might be reasonable to drop a p-goal (e.g.  $\Diamond\text{GetPhD}$ ) that is in conflict with another higher priority active p-goal (e.g.  $\Box\text{BeRich}$ ), in our framework we keep such p-goals around. The reason for this is that although  $\Box\text{BeRich}$  is currently inconsistent with  $\Diamond\text{GetPhD}$ , the agent might later learn that  $\Box\text{BeRich}$  has become impossible to bring about (e.g. after  $goBankrupt$  occurs), and then might want to pursue

<sup>5</sup>  $\text{CorrectHist}(s, path)$  is defined as  $\text{Starts}(path, s') \wedge \text{SameHistory}(s', s)$ ; here  $path$  is a placeholder that stands for a path and  $s$  represents the current situation.

◇GetPhD. Thus, it is useful to keep these inactive p-goals since this allows the agent to maximize her utility (that of her chosen goals) by taking advantage of such opportunities. As mentioned earlier, c-goals are our analogue to intentions. Recall that Bratman’s model of intentions limits the agent’s practical reasoning – agents do not always optimize their utility and don’t always reconsider all available options in order to allocate their reasoning effort wisely. In contrast to this, our c-goals are defined in terms of the p-goals, and at every step, we ensure that the agent’s c-goals maximizes her utility so that it is the set of highest priority goals that is consistent given the agent’s knowledge. Thus, our notion of c-goals is not as persistent as Bratman’s notion of intention [10]. For instance as mentioned above, after the action *goBankrupt* happens in  $S_0$ , the agent will lose the c-goal that  $\Box\text{BeHappy}$ , although she did not drop it and it did not become impossible or achieved. In this sense, our model is that of an idealized agent. There is a tradeoff between optimizing the agent’s chosen set of prioritized goals and being committed to chosen goals. In our framework, chosen goals behave like intentions with an automatic filter-override mechanism [10] that forces the agent to drop her chosen goals when opportunities to commit to other higher priority goals arise. In the future, it would be interesting to develop a logical model that captures the pragmatics of intention reconsideration by supporting control over it.

We now show that our formalization of prioritized goals has some desirable properties. Some of these (e.g. Proposition 3a) are analogues of the AGM postulates; others (e.g. adopting logically equivalent goals has the same result, etc.) were left out for space reasons. First we show that c-goals are consistent:

**Proposition 1 (Consistency)**  $D \models \forall s. \neg \text{CGoal}(\text{False}, s)$ .

Thus, the agent cannot have both  $\phi$  and  $\neg\phi$  c-goals in a situation  $s$  and there is a path that is  $G_C$ -accessible in  $s$ . Even if all of the agent’s p-goals become known to be impossible, the set of  $G_C$ -accessible paths will be precisely the ones that starts with a  $K$ -accessible situation, and thus the agent will only choose the propositions that are known to be inevitable.

We also have the property of realism [1], i.e. if an agent knows that something has become inevitable, then she has this as a c-goal:

**Proposition 2 (Realism)**  $D \models \text{KInevitable}(\phi, s) \supset \text{CGoal}(\phi, s)$ .

Note that this is not necessarily true for p-goals and primary c-goals – an agent may know that something has become inevitable and not have it as her p-goal/primary c-goal, which is intuitive. In fact, this is the reason why we define p-goals in terms of  $G$ -accessible paths rather than  $G_R$ . A consequence of Proposition 1 and 2 is that an agent does not have a c-goal that is known to be impossible, i.e.  $D \models \text{CGoal}(\phi, s) \supset \neg \text{KImpossible}(\phi, s)$ .

We next discuss some properties of the framework w.r.t. goal change. Proposition 3 says that (a) an agent acquires the p-goal at some level  $n$  that  $\phi$  after she adopts it, and (b) an agent acquires the primary c-goal (and c-goal) that  $\phi$  after she adopts it in  $s$ , provided that she does not have the c-goal that  $\neg\phi$  in  $s$ .

**Proposition 3 (Adoption)** (a)  $D \models \exists n. \text{PGoal}(\phi, n, \text{do}(\text{adopt}(\phi), s))$ ,  
 (b)  $D \models \neg \text{CGoal}(\neg\phi, s) \supset \text{PrimCGoal}(\phi, \text{do}(\text{adopt}(\phi), s))$ .

We can also show that after dropping the p-goal that  $\phi$  in  $s$ , an agent does not have the p-goal (and thus the primary c-goal) that  $\phi$ .

**Proposition 4 (Drop)**  $D \models \neg \exists n. \text{PGoal}(\phi, n, \text{do}(\text{drop}(\phi), s))$ .

Note that, this does not hold for  $\text{CGoal}$ , as  $\phi$  could still be a consequence of two or more of her remaining primary c-goals.

The next few properties concern the persistence of these motivational attitudes. In the following, we prove a persistence property for achievement p-goals:

**Proposition 5 (Persistence of Achievement PGoals)**

$D \models \text{PGoal}(\diamond\Phi, n, s) \wedge \text{Know}(\neg\Phi(\text{now}), s) \wedge \forall\psi. a \neq \text{drop}(\psi) \supset \text{PGoal}(\diamond\Phi, n, \text{do}(a, s))$ .

This says that if an agent has a p-goal that  $\diamond\Phi$  in  $s$ , then she will retain this p-goal after some action  $a$  has been performed in  $s$ , provided that she knows that  $\Phi$  has not yet been achieved, and  $a$  is not the action of dropping a p-goal. Note that, we do not need to ensure that  $\diamond\Phi$  is still known to be possible or consistent with higher priority active p-goals, since the SSA for  $G$  does not automatically drop such incompatible p-goals from the goal hierarchy.

For achievement chosen goals we have the following:

**Proposition 6 (Persistence of Achievement Chosen Goals)**

$$D \models \text{OPGoal}(\diamond\Phi \wedge \text{CorrectHist}(s), n, s) \wedge \text{CGoal}(\diamond\Phi, s) \\ \wedge \text{Know}(\neg\Phi(\text{now}), s) \wedge \forall\psi. a \neq \text{drop}(\psi) \wedge \neg\text{CGoal}(\neg\diamond\Phi, n-1, \text{do}(a, s)) \\ \supset \text{CGoal}(\diamond\Phi, n, \text{do}(a, s)).$$

Thus, in situation  $s$ , if an agent has the only p-goal at level  $n$  that  $\diamond\Phi$  and the correct history of actions in  $s$  has been performed, and if  $\diamond\Phi$  is also a c-goal in  $s$  (and thus she has the primary c-goal that  $\diamond\Phi$ ), then she will retain the c-goal that  $\diamond\Phi$  at level  $n$  after some action  $a$  has been performed in  $s$ , provided that:

- she knows that  $\Phi$  has not yet been achieved,
- that  $a$  is not the action of dropping a p-goal,
- and that at level  $n-1$  the agent does not have the c-goal that  $\neg\diamond\Phi$ , i.e.  $\diamond\Phi$  is consistent with higher priority c-goals.

Note that this property also follows if we replace the consequent with  $\text{CGoal}(\diamond\Phi, \text{do}(a, s))$  or  $\text{PrimCGoal}(\diamond\Phi, \text{do}(a, s))$ , and thus it deals with the persistence of (primary) c-goals. Note however that, it does not hold if we replace the  $\text{OPGoal}$  in the antecedent with  $\text{PGoal}$ ; the reason for this is that the agent might have a p-goal at level  $n$  in  $s$  that  $\phi$  and the c-goal in  $s$  that  $\phi$ , but not have  $\phi$  as a primary c-goal in  $s$ , e.g.  $n$  might be an inactive level because another p-goal at  $n$  has become impossible, and  $\phi$  could be a c-goal in  $s$  because it is a consequence of two other primary c-goals. Thus even if  $\neg\phi$  is not a c-goal after  $a$  has been performed in  $s$ , there is no guarantee that the level  $n$  will be active in  $\text{do}(a, s)$  or that all the active p-goals that contributed to  $\phi$  in  $s$  are still active.

We believe that the dropping of an unrelated p-goal will not affect persistence, and hence it should be possible to strengthen Proposition 5 and 6. Also, in the future we would like to generalize these two propositions to deal with arbitrary temporally extended goals.

## 6 Handling Subgoals

In this section, we deal with the dynamics of subgoals. As mentioned earlier, a subgoal must be dropped when the parent goal is dropped or becomes impossible. When adopting a subgoal  $\psi$  with respect to a supergoal  $\phi$ , in addition to recording the newly adopted goal  $\psi$ , we need to model the fact that  $\psi$  is a subgoal of  $\phi$ . This information can later be used to drop the subgoal when the parent goal is dropped. One way of modeling this is to ensure that the adoption of a subgoal  $\psi$  w.r.t. a parent goal  $\phi$  adds new p-goals that contain *both this subgoal and this parent goal i.e.  $\psi \wedge \phi$  at a lower priority than the parent goal  $\phi$* . This ensures that when the parent goal is dropped, the subgoal is also dropped. To see this, recall that to handle the dropping of a goal  $\phi$ , we drop the p-goals at all  $G$ -accessibility levels that imply  $\phi$ . Thus, if we drop the parent goal  $\phi$ , it will also drop all of its subgoals including  $\psi$ , since the  $G$ -accessibility levels where the parent goal  $\phi$  holds include the  $G$ -accessibility levels where the subgoal  $\psi$  holds. Note that, if there are more than one level where the supergoal  $\phi$  is a p-goal, then we copy all these levels, i.e. for each level  $n$  where  $\phi$  is a p-goal, we add a (lower priority) level to the goal hierarchy. As we will see, this ensures that the sub-subgoals and sub-sub-subgoals etc. are also properly dropped when the supergoal is dropped. Also, this means that dropping a subgoal does not necessarily drop the supergoal.

Before going over the formal details, let us mention some useful bookkeeping tools that we will use:  $\text{Length}(l)$  returns the number of elements in list  $l$ ;  $\text{Nth}(l, i)$  returns the  $i$ -th element in list  $l$ , and -1 if  $i > \text{Length}(l)$ ;  $\text{Sort}(l)$  returns a sorted version of list  $l$ . The part of the SSA for  $G$  that handles subgoal adoption is defined as follows:

$$\begin{aligned} \text{SubGoalAdopted}(p, n, a, s, \psi, \phi) \stackrel{\text{def}}{=} & (n < \text{NPGoals}(s) \wedge \text{Progressed}(p, n, a, s)) \vee \\ & (\text{NPGoals}(s) \leq n < \text{NPGoals}(s) + \text{Length}(\text{AddList}(\phi, s)) \\ & \wedge \exists i, m. (n = \text{NPGoals}(s) + i \wedge m = \text{Nth}(\text{AddList}(\phi, s), i) \\ & \wedge \text{Progressed}(p, m, a, s) \wedge \psi(p))) \vee \\ & (n \geq \text{NPGoals}(s) + \text{Length}(\text{AddList}(\phi, s)) \wedge \text{Progressed}(p, n, a, s)). \end{aligned}$$

That is, if the action involves the adoption of a subgoal  $\psi$  w.r.t. a supergoal  $\phi$ , we adjust  $G$  to incorporate (possibly several) new p-goals. We will discuss each case in turn. First, note that the existing p-goals are just carried over by progressing them; this is handled by the first disjunct.

Secondly, we adjust  $G$  starting at level  $\text{NPGoals}(s)$ . We add a number of new levels that include the conjunction of the only p-goal and the subgoal at a lower priority for all the current only p-goals that imply the parent goal  $\phi$ . For example, say at level  $i$  we have an OPGoal that  $\phi_i$  and it implies the parent goal that  $\phi$ ; then we add at a lower priority the conjoined goal of the progressed version of  $\phi_i$  and the subgoal  $\psi$ . Our formalization of this uses the abbreviation  $\text{AddList}(\phi, s)$  which is a sorted list of levels such that the parent goal is implied by the only p-goal at this level.  $\text{AddList}$  is defined as:  $\text{AddList}(\phi, s) \stackrel{\text{def}}{=} \text{Sort}([n \mid \text{PGoal}(\phi, n, s)])$ . The length of this list indicates the number of lower priority goals that needs to be added. As discussed above, this ensures that the agent will drop the subgoal when the parent goal is dropped (but not necessarily vice-versa). Note that if this process adds two or more new p-goals to the

agent's goal hierarchy, we maintain the original ordering; e.g. suppose that the agent adopted  $\psi$  w.r.t.  $\phi$ , that there are two  $G$ -accessibility levels  $m$  and  $n$  such that the agent has the only p-goal that  $\phi_m$  at  $m$  and  $\phi_n$  at  $n$ , that  $\phi_m$  implies  $\phi$  and  $\phi_n$  implies  $\phi$ , and that  $n > m$ . In that case, the SSA for  $G$  will add the p-goal  $\phi_m \wedge \phi$  at level  $NPGoals(s)$  and the p-goal  $\phi_n \wedge \phi$  at level  $NPGoals(s) + 1$ .

Finally, all the remaining levels involving trivially true goals are just carried over by progressing them.

The part of the SSA for  $NPGoals$  that handles subgoal adoption is defined as follows:

$$\text{AdjustSubGoalAdopt}(\phi, s) \stackrel{\text{def}}{=} NPGoals(s) + \text{Length}(\text{AddList}(\phi, s)).$$

That is, when the agent adopts a subgoal w.r.t. a parent goal, her current  $NPGoals$  is incremented by the number of new p-goals adopted in this process.

Let us go back to our example. Suppose that the agent knows that one way of always being rich is to always work hard, which in turns can be fulfilled by always being energetic. Assume that with this in mind, our agent adopts the subgoal that  $\Box\text{WorkHard}$  w.r.t. the p-goal that  $\Box\text{BeRich}$ , and then adopts the sub-subgoal that  $\Box\text{BeEnergetic}$  w.r.t. the subgoal that  $\Box\text{WorkHard}$  starting in  $S_0$ . Then the agent's goal hierarchy in  $S_1 = do(adopt(\Box\text{WorkHard}, \Box\text{BeRich}), S_0)$  should include the p-goal that  $\Box\text{WorkHard}$  and in  $S_2 = do(adopt(\Box\text{BeEnergetic}, \Box\text{WorkHard}), S_1)$  should also include the p-goal that  $\Box\text{BeEnergetic}$ . According to the SSA for  $G$ , our agent's goal hierarchy in  $S_1$  and in  $S_2$  will be as in Table 1.<sup>6</sup> In  $S_0$ , the supergoal  $\Box\text{BeRich}$  holds at level 0 and thus  $\text{AddList}(\Box\text{BeRich}, S_0) = [0]$ . Similarly in  $S_1$ , the supergoal  $\Box\text{WorkHard}$  holds at level 3 and thus  $\text{AddList}(\Box\text{WorkHard}, S_1) = [3]$ . Now, suppose that in  $S_2$  the agent wants to drop the p-goal that  $\Box\text{WorkHard}$ . Then in  $S_3 = do(drop(\Box\text{WorkHard}), S_2)$ , she should no longer have  $\Box\text{BeEnergetic}$  as a p-goal, but should retain the supergoal that  $\Box\text{BeRich}$ . After the agent drops the p-goal that  $\Box\text{WorkHard}$ , by the SSA for  $G$  we can see that all the  $G$ -accessible levels where  $\Box\text{WorkHard}$  holds will be replaced by the only p-goal that  $\text{CorrectHist}(S_2, path)$  (see  $S_3$  in Table 1). This shows that dropping  $\Box\text{WorkHard}$  results in the dropping of all of its subgoals (in this case  $\Box\text{BeEnergetic}$ ), but that its parent goal  $\Box\text{BeRich}$  is retained.

We define the SubGoal relation as follows:

$$\begin{aligned} \text{SubGoal}(\psi, \phi, s) \stackrel{\text{def}}{=} & \exists n. \text{PGoal}(\phi, n, s) \wedge \neg \text{PGoal}(\psi, n, s) \\ & \wedge \forall n. \text{PGoal}(\psi, n, s) \supset \text{PGoal}(\phi, n, s). \end{aligned}$$

This says that  $\psi$  is a subgoal of  $\phi$  in situation  $s$  iff there exists an  $G$ -accessibility level  $n$  in  $s$  such that  $\phi$  is a p-goal at  $n$  while  $\psi$  is not, and for all  $G$ -accessibility levels in  $s$  where  $\psi$  is a p-goal,  $\phi$  is also a p-goal. Note that, while our formalization of subgoal dynamics allows a subgoal to have multiple parents, in this definition we assume that a subgoal can't have more than one parent. In the future, we will work on relaxing this constraint.

We now discuss some properties concerning the dynamics of subgoals and the dependencies between a subgoal and its parent goal. Proposition 7 states that (a) an agent

<sup>6</sup> For simplicity in Table 1, we only show the agent's relevant p-goals rather than its only p-goals (which in addition reflect the actions that have been performed so far, i.e.  $\text{CorrectHist}(s)$ ).

acquires the p-goal that  $\psi$  after she adopts it as a subgoal of another goal  $\phi$  in  $s$ , provided that she has the p-goal at some level in  $s$  that  $\phi$ , and (b) she also acquires the primary c-goal that  $\psi$  after she adopts it as a subgoal of  $\phi$  in  $s$ , provided that she has the primary c-goal in  $s$  that  $\phi$ , and that she does not have the c-goal that  $\neg\psi$  in  $s$ .

**Proposition 7 (Subgoal Adoption)**

- (a)  $D \models \exists m. \text{PGoal}(\phi, m, s) \supset \exists n. \text{PGoal}(\psi, n, do(adopt(\psi, \phi), s))$ ,
- (b)  $D \models \text{PrimCGoal}(\phi, s) \wedge \neg \text{CGoal}(\neg\psi, s) \supset \text{PrimCGoal}(\psi, do(adopt(\psi, \phi), s))$ .

The next property says that after dropping the p-goal that  $\phi$  in  $s$ , an agent does not have the p-goal (and thus the primary c-goal) that  $\psi$ , provided that  $\psi$  is a subgoal of  $\phi$  in  $s$ .

**Proposition 8 (Supergoal Drop)**

$$D \models \text{SubGoal}(\psi, \phi, s) \supset \neg \exists n. \text{PGoal}(\psi, n, do(drop(\phi), s)).$$

As with Proposition 4, this does not hold if we replace PGoal in the consequence with CGoal since  $\psi$  could be a consequence of a combination of other active p-goals.

The next two properties say that dropping a subgoal does not effect the parent goal.

**Proposition 9 (Subgoal Drop)**

- (a)  $D \models \text{SubGoal}(\psi, \phi, s) \supset \exists n. \text{PGoal}(\phi, n, do(drop(\psi), s))$ ,
- (b)  $D \models \text{SubGoal}(\psi, \phi, s) \wedge \text{PrimCGoal}(\phi, s) \supset \text{PrimCGoal}(\phi, do(drop(\psi), s))$ .

That is, (a) an agent retains the p-goal that  $\phi$  after she drops a subgoal  $\psi$  of  $\phi$ , and (b) she also retains the primary c-goal that  $\phi$  after she drops a subgoal  $\psi$  of  $\phi$  in  $s$ , provided that she has the primary c-goal that  $\phi$  in  $s$ .

Finally, it can be shown that the SubGoal relation is transitive, i.e. if  $\psi_1$  is a subgoal of  $\phi$  in  $s$ , and if  $\psi_2$  is a subgoal of  $\psi_1$  in  $s$ , then  $\psi_2$  must also be a subgoal of  $\phi$  in  $s$ .

## 7 Discussion and Future Work

In this paper, we presented a formalization of prioritized goals, subgoals, and their dynamics. Our formalization ensures that an agent's chosen goals are always consistent and that goals and subgoals properly evolve as a result of actions and as a result of adopting and dropping goals. Although we made some simplifying assumptions, in this paper we have focused on developing an expressive framework that captures an idealized form of rationality without worrying about tractability. It would be desirable to study restricted fragments of the logic where reasoning is tractable. Also, before defining more limited forms of rationality, one should have a clear specification of what ideal rationality really is so that one understands what compromises are being done.

While in our account chosen goals are closed under logical consequence, primary c-goals are not. To this end, our formalization of primary c-goals is related to the non-normal modal formalizations of intentions found in the literature [3], and as such it does not suffer from the side-effect problem [1]. For instance, in our framework an agent can have the primary c-goal to get her teeth fixed and know that this always involves pain, but not have the primary c-goal to have pain.

Also, since we are using the situation calculus, we can easily represent procedural goals/plans, e.g. the goal to do  $a_1$  and then  $a_2$  can be written as:  $\text{PGoal}(\text{Starts}(p, s_1) \wedge \text{OnPath}(p, s) \wedge s = \text{do}(a_2, \text{do}(a_1, s_1)), 0, S_0)$ . Golog [12] can be used to represent complex plans/programs. So we can model the adoption of plans as subgoals.

Recently, there have been a few proposals that deal with goal change. Shapiro *et al.* [18] present a situation calculus based framework where an agent adopts a goal when she is requested to do so, and remains committed to this goal unless the requester cancels this request; a goal is retained even if the agent learns that it has become impossible, and in this case the agent's goals become inconsistent. Shapiro and Brewka [7] modify this framework to ensure that goals are dropped when they are believed to be impossible or when they are achieved. Their account is similar to ours in the sense that they also assume a priority ordering over the set of (in their case, requested) goals, and in every situation they compute chosen goals by computing a maximal consistent goal set that is also compatible with the agent's beliefs. However, their model has some unintuitive properties: the agent's chosen set of goals in  $\text{do}(a, s)$  may be quite different from her goals in  $s$ , although  $a$  did not make any of her goals in  $s$  impossible or inconsistent with higher priority goals, because inconsistencies between goals at the same priority level are resolved differently (this can happen because goals are only partially ordered). Note that, while one might argue that a partial order over goals might be more general, allowing this means that additional control information is required to obtain a single goal state after the agent's goals change. Also, we provide a more expressive formalization of prioritized goals – we model goals using infinite paths, and thus can model many types of goals that they cannot. Finally they model prioritized goals by treating the agent's p-goals as an arbitrary set of temporal formulas, and then defining the set of c-goals as a subset of the p-goals. But our possible world semantics has some advantages over this: it clearly defines when goals are consistent with each other and with what is known. One can easily specify how goals change when an action  $a$  occurs, e.g. the goal to do  $a$  next and then do  $b$  becomes the goal to do  $b$  next, the goal that  $\diamond\Phi \vee \diamond\Psi$  becomes the goal that  $\diamond\Psi$  if  $a$  makes achieving  $\Phi$  impossible, etc.

There has been much work on agent programming languages with declarative goals where the dynamics of goals and intentions and the dependencies between goals and sub-goals are modeled (e.g. [19], [9] and the references therein). However, most of these are not based on a formal theory of agency, and to the best of our knowledge, none maintains the consistency of (chosen) goals (e.g. when adopting a plan to achieve a goal, these frameworks do not ensure that this plan is consistent with the agent's other concurrent goals/plans). Also, most of these do not deal with temporally extended goals, and as a result they often need to accommodate inconsistent goal-bases to allow the agent to achieve conflicting states at different time points (e.g. the default logic based framework in [20]); chosen goals are required to be consistent. In [6], the authors present a situation calculus based agent programming language where the agent executes a program while maximizing the achievement of a set of prioritized goals. However, they do not formalize goal dynamics.

One limitation of our account is that we assume that the agent's p-goals are totally ordered in terms of priority. Also, newly adopted p-goals are assigned the lowest priority. A consequence of this is that an agent's c-goals depend on the adoption order of

her p-goals. For instance, given a fixed starting situation, an agent can end up with two different sets of c-goals by adopting  $\phi$  followed by  $\psi$ , and by adopting  $\psi$  followed by  $\phi$ . This has very different results when  $\phi$  and  $\psi$  conflict with each other. We would like to address this by incorporating the priority of the p-goal as an argument to the *adopt* action, and handling this in the framework. Finally, one can argue that our agent spends too much resources trying to optimize her c-goals at every step. It would be interesting to develop an account where the agent is strongly committed to her chosen goals, and where the filter override mechanism is only triggered under specific conditions.

## References

1. Cohen, P.R., Levesque, H.J.: Intention is Choice with Commitment. *Artificial Intelligence* **42**(2–3) (1990) 213–361
2. Sadek, M.D.: A Study in the Logic of Intention. In: *Third Intl. Conf. on Principles of Knowledge Representation and Reasoning (KR&R-92)*, Cambridge, MA (1992) 462–473
3. Konolige, K., Pollack, M.E.: A Representationalist Theory of Intention. In: *Thirteenth Intl. J. Conf. on Artificial Intelligence (IJCAI-93)*, Chambéry, France (1993) 390–395
4. Singh, M.P.: *Multiagent Systems: A Theoretical Framework for Intentions, Know-How, and Communications*. Volume 799 of LNAI. Springer-Verlag, Germany (1994)
5. Shapiro, S., Lespérance, Y., Levesque, H.J.: Goals and Rational Action in the Situation Calculus - A Preliminary Report. In: *Working Notes of the AAAI Fall Symposium on Rational Agency: Concepts, Theories, Models, and Applications*, Cambridge, MA (November 1995) 117–122
6. Sardina, S., Shapiro, S.: Rational Action in Agent Programs with Prioritized Goals. In: *Second Intl. J. Conf. on Autonomous Agents and Multi-Agent Sys. (AAMAS-03)*, Melbourne, Australia (2003) 417–424
7. Shapiro, S., Brewka, G.: Dynamic Interactions Between Goals and Beliefs. In: *Twentieth Intl. J. Conf. on Artificial Intelligence (IJCAI-07)*, India (2007) 2625–2630
8. Winikoff, M., Padgham, L., Harland, J., Thangarajah, J.: Declarative and Procedural Goals in Intelligent Agent Systems. In: *Eighth Intl. Conf. on Principles and Knowledge Representation and Reasoning (KR&R-02)*, Toulouse, France (2002) 470–481
9. Bordini, R.H., Dastani, M., Dix, J., Fallah-Seghrouchni, A.E., eds.: *Multi-Agent Programming: Languages, Platforms and Applications*. Springer (2005)
10. Bratman, M.E.: *Intentions, Plans, and Practical Reason*. Harvard University Press, Cambridge, MA (1987)
11. McCarthy, J., Hayes, P.J.: Some Philosophical Problems from the Standpoint of Artificial Intelligence. *Machine Intelligence* **4** (1969) 463–502
12. Reiter, R.: *Knowledge in Action. Logical Foundations for Specifying and Implementing Dynamical Systems*. MIT Press, Cambridge, MA (2001)
13. DeGiacomo, G., Lespérance, Y., Levesque, H.J.: ConGolog, a Concurrent Programming Language Based on the Situation Calculus. *Artificial Intelligence* **121** (2000) 109–169
14. Levesque, H.J., Pirri, F., Reiter, R.: Foundations for a Calculus of Situations. *Electronic Transactions of AI (ETAI)* **2**(3–4) (1998) 159–178
15. Moore, R.C.: A Formal Theory of Knowledge and Action. In Hobbs, J.R., Moore, R.C., eds.: *Formal Theories of the Commonsense World*. Ablex (1985) 319–358
16. Scherl, R., Levesque, H.: Knowledge, Action, and the Frame Problem. *Artificial Intelligence* **144**(1–2) (2003) 1–39



17. Rao, A.S., Georgeff, M.P.: Modeling Rational Agents with a BDI-Architecture. In Fikes, R., Sandewall, E., eds.: *Second Intl. Conf. on Principles of Knowledge Representation and Reasoning (KR&R-91)*, San Mateo, CA, Morgan Kaufmann Publishers (1991) 473–484
18. Shapiro, S., Lespérance, Y., Levesque, H.J.: Goal Change in the Situation Calculus. *J. of Logic and Computation* **17**(5) (2007) 983–1018
19. Sardina, S., deSilva, L., Padgham, L.: Hierarchical Planning in BDI Agent Programming Languages: A Formal Approach. In: *Fifth Intl. J. Conf. on Autonomous Agents and Multi-Agent Systems (AAMAS-06)*, Hakodate, Japan (2006) 1001–1008
20. van Riemsdijk, M.B., Dastani, M., Meyer, J.J.C.: Semantics of Declarative Goals in Agent Programming. In: *Fourth Int'l J. Conf. on Autonomous Agents and Multiagent Sys. (AAMAS-05)*. (2005) 133–140

# Microscopic Analysis of Edge Creation Process in Trust-based Social Networks

Babak Khosravifar<sup>1</sup>, Maziar Gomrokchi<sup>2</sup>, and Jamal Bentahar<sup>3</sup>  
*b\_khosr@encs.concordia.ca, m\_gomrok@encs.concordia.ca,*  
*bentahar@ciise.concordia.ca*

<sup>1</sup>Department of Electrical and Computer Engineering, <sup>2</sup>Computer Science,  
<sup>3</sup>Concordia Institute for Information System Engineering, Concordia University,  
Montreal, Canada

**Abstract.** In last years, online social networks are enjoying drastic increase in their population and connectivity. Meanwhile, trust is known as essential factor in quality of the connections between diverse nodes in the network. To address the efficiency in the interactions of nodes, we propose in this paper a trust-based architecture applicable to maintain interactions in multi-agent-based social network. We provide a detailed discussion over the network formation by taking into account the edge creation factors classified as homophily, confounding and influence. We systematically inspire different involving factors to observe evolution of the network of trust-based interconnections in a microscopic manner. We also provide a theoretical analysis of the proposed model assessment and discuss the system implementation, along with simulations obtained from a number of executions compared with the broadly known frameworks.

**General Terms.** Human Factors, Measurement, Experimentation.

**Keywords.** Trust establishment , Edge Creation, Multi-Agent Systems, Agent Communication, Social Networks.

## 1 Introduction

Online social networks are drastically being enlarged. Facebook, Flickr, Yahoo Answers are among very popular social networks that are gaining a very high traffic in terms of the users and their conductivities. In general, the impact of the features of these networks and analysis on how they form the behavior of the users have been a hot interest during the very recent years. A number of theoretical and empirical works have been proposed analyzing the users behavior in forming the connection between them. For example, the analysis on the edge creation process [10, 6] related to the sociality of a node led to hypothesis to observe the distribution of a heavy-traffic degree of popular nodes in the network. In [1], the authors address the source of the correlation between agents that led them to extend their activity and create edges. In [2], the correlation between agents are analyzed in an online large scale community. In fact, the relation between the agents that just joined the community and the agents that are already in the community is discussed.

Here we take a different perspective over the social correlation of agents. In edge creation process, we analyze diverse impacts upon the trust that is already established between two nodes and consequently, we generalize to analyze the socializing of the agents that use different trust evaluation systems. To this end, we first discuss the trust evaluating method and upon that, we analyze the activity area extension of agents in a microscopic approach. To maintain a trust-based network, different computational frameworks have been proposed in the literature. Some models consider the direct interaction of two parties [15, 14, 8]. Some models rely, to some extent, on the suggested rating provided by other agents [12, 9, 13]; and some others also consider the suggested rating of the agent being evaluated [4, 8]. Since agents are self-interested, it is hard to analyze an agent's likely behavior based on previous direct interactions given the fact that the collected information from other agents may be non-reliable and could lead to a non-accurate trust assessment. So far, these frameworks do not act properly if selfish agents tend to change their behaviors. Therefore, agents do not properly initiate a social activity in the sense that they cannot maintain a strong control on the adjacent environment.

In this paper, we use the proposed model that is fully presented in [11] and discuss the social network-related affects (classified in [1]) such as influence, confounding, and homophily on the edge creation process of the distributed agents. In the proposed model, we provide an efficient assessment process in a twofold contribution. In the first contribution, agents mutually interact and rate each other based on the interaction done (either satisfactory or dissatisfactory). The obtained ratings are accumulated to assess the direct interaction rating of a particular agent. In the evaluation process, we call the evaluator agent as *the trustor* and refer to him as  $Ag_a$ , and the evaluated agent as *the trustee* and refer to him as  $Ag_b$ . Also we refer to the involved agents in the trust assessment process as *consulting agents*. In the proposed framework,  $Ag_a$  evaluates the credibility of  $Ag_b$  by combining his own direct trust rating with the ratings provided by the consulting agents. The computed trust value is used as a doubt to extend the connection between the trustor and the trustee agent. In the second contribution of the proposed model, the trustor agent after a period of direct interaction with the trustee agent performs a retrospect trust adjustment (so-called maintenance) in order to update his belief set about the credibility of the consulting agents that provided information regarding to the trust level of trustee agents. Depending on the accuracy of the consulting agents, the trustor agent would either increase or decrease his trust ratings about consulting agents. By updating the trust values, agents basically would recognize the adjacent agents that are worth to extend the connection with. On the other hand, agents would stay apart from bad agents in terms of social activities. Doing so, gradually agents recognize more reliable consulting agents around in the network, which would form a more efficient social correlation between adjacent agents. This assessment is used in microscopic analysis of edge creation between the interacting agents.

The remainder of this paper is organized as follows. In Section 2, we briefly define our proposed framework as comprehensive trust assessment process, which

is composed of evaluation and maintenance process. In Section 3, we define the social network parameters, and environment that the interactions are initiated. In Section 4, we elaborate the affects of the discussed network properties in the experimental environment. Representing the testbed, we compare our model results with two well-known trust models in terms of efficiency in trust assessment and finally Section 5 concludes the paper.

## 2 Trust Evaluation Method

In this section, we formalize the trust assessment between interacting parties (customer and provider agents) in the social network. Since the proposed trust assessment model is fully described in [11], in this paper we only refer to the assessment formulation. In general, each customer agent  $c_i$  is linked to a set of customers it knows and a set of provider agents it has interacted within the past. Without loss of generality, a customer  $c_i$  has to evaluate correlating agents in the environment (could be consumers or providers) in any case of interaction. For instance, a direct evaluation of a provider  $p_j$  is possible if the customer had enough transactions with that provider. In the trust model that we use for evaluation [3], three elements are used to characterize the relationship between the trustor  $Ag_a$  and trustee agents  $Ag_b$ : 1) how much the trustor agent trusts the trustee:  $Tr_{Ag_a}^{Ag_b}$ ; 2) the number of past transactions:  $NT_{Ag_a}^{Ag_b}$ ; and 3) the time recency of the last transactions:  $TiR_{Ag_a}^{Ag_b}$ . The transactions could taken place caused from homophily, confounding or influence of any other agents in the network. Formally, we define a social network for service selection as follows:

**Definition 1.** A social network for service selection is a tuple  $\langle C, P, \xrightarrow{\text{tc}}, \xrightarrow{\text{tp}} \rangle$  where  $C$  is a set of customer services,  $P$  a set of provider services,  $\xrightarrow{\text{tc}} \subseteq C \times \mathbb{R}^3 \times C$  is a ternary relation (for labelled edges linking customers) and  $\xrightarrow{\text{tp}} \subseteq C \times \mathbb{R}^3 \times P$  is a ternary relation (for labelled edges linking customers to providers).

We use the usual notation for the labelled edges: if  $c_i, c_k \in C$  and  $v \in \mathbb{R}^3$ , then  $(c_i, v, c_k) \in \xrightarrow{\text{tc}}$  is written as  $c_i \xrightarrow{v}_{\text{tc}} c_k$ . Likewise, we write  $c_i \xrightarrow{v}_{\text{tp}} p_j$  instead of  $(c_i, v, p_j) \in \xrightarrow{\text{tp}}$ . Our social network for service selection has two types of nodes: type 1 for customers and type 2 for providers and two types of edges: type 1 for edges between customers and type 2 for edges linking customers to providers. The edges of type 1 represent friendship relations in the network, while edges of type 2 capture business relationships. The existence of an edge of type 1  $c_i \xrightarrow{v}_{\text{tc}} c_k$  means that  $c_i$  knows (is friend of)  $c_k$  such that:  $v = (Tr_{c_i}^{c_k}, NI_{c_i}^{c_k}, TiR_{c_i}^{c_k})$ . The existence of an edge of type 2  $c_i \xrightarrow{v}_{\text{tp}} p_j$  means that  $c_i$  had transactions with  $p_j$  such that:  $v = (Tr_{c_i}^{p_j}, NT_{c_i}^{p_j}, TiR_{c_i}^{p_j})$ . We note that there is no edges in this social network between providers. This does not mean that there is no social link between providers, but only the existing links (which could be collaborations or competitions) are not used in our framework.

The direct evaluation of a provider  $p_j$  by a customer  $c_i$  is based on the ratings  $c_i$  gave to  $p_j$  for each past interaction ( $r_l$ ) combined with the importance of that interaction ( $\lambda_l$ ) and its time recency. Let  $n$  be the number of total transactions between  $c_i$  and  $p_j$  ( $n = NT_{c_i}^{p_j}$ ), equation 1 gives the formula to compute this evaluation. To perform the indirect evaluation, the customer  $c_i$  solicits information about the provider  $p_j$  from other customers  $c_k$  such that there is an edge  $c_i \xrightarrow{v}_{c_c} c_k$  in the social network. The set of these customers  $c_k$  is denoted  $\mathcal{T}_{c_i}$ . The equation computing the indirect estimation is given by equation 2, where  $\alpha Tr_{c_i}^{c_k} = Tr_{c_i}^{c_k} \cdot NI_{c_i}^{c_k} \cdot TiR_{c_i}^{c_k}$ .

$$DTr_{c_i}^{p_j} = \frac{\sum_{l=1}^n (\lambda_l \cdot TiR_{c_i}^{p_j} \cdot r_l)}{\sum_{l=1}^n (\lambda_l \cdot TiR_{c_i}^{p_j})} \quad (1)$$

$$ITr_{c_i}^{p_j} = \frac{\sum_{c_k \in \mathcal{T}_{c_i}} \alpha Tr_{c_i}^{c_k} \cdot Tr_{c_k}^{p_j} \cdot TiR_{c_k}^{p_j} \cdot NT_{c_k}^{p_j}}{\sum_{c_k \in \mathcal{T}_{c_i}} \alpha Tr_{c_i}^{c_k} \cdot TiR_{c_k}^{p_j} \cdot NT_{c_k}^{p_j}} \quad (2)$$

To compute  $Tr_{c_i}^{p_j}$ , the direct and indirect evaluations are combined according to their proportional importance. The idea is that the customer relies, to some extent, on its own history (direct trust evaluation) and on consulting with its network (indirect trust evaluation). This merging method considers the proportional relevance of each trust assessment, rather than treating them separately. To this end,  $c_i$  assigns a contribution value for the trust assessment method ( $\omega$  for direct trust evaluation and  $1 - \omega$  for indirect trust evaluation when  $\omega < 1$ ). The value  $\omega$  is obtained from equation 3. Basically, the contribution of each approach in the evaluation of  $p_j$  is defined regarding to: (1) how informative the history is in terms of the number of direct transactions between  $c_i$  and  $p_j$  ( $NT_{c_i}^{p_j}$ ) and their time recency ( $TiR_{c_i}^{p_j}$ ); and (2) how informative and reliable the consulting customers are from  $c_i$ 's point of view ( $DTr_{c_i}^{p_j}$ ). Therefore, consultation with other agents is less considered if the history represents a comparatively higher entropy value  $\omega$ , which reflects lower uncertainty. Respecting the contribution percentage of the trust assessments,  $c_i$  computes the trust value for  $p_j$  using equation 4.

$$\omega = \frac{\ln(DTr_{c_i}^{p_j} \cdot NT_{c_i}^{p_j} \cdot TiR_{c_i}^{p_j})}{\sum_{c_k \in \mathcal{T}_{c_i}} \ln(DTr_{c_i}^{c_k} \cdot NI_{c_i}^{c_k} \cdot TiR_{c_i}^{c_k})} \quad (3)$$

$$Tr_{c_i}^{p_j} = \begin{cases} \omega \cdot DTr_{c_i}^{p_j} + (1 - \omega) \cdot ITr_{c_i}^{p_j} & \text{if } \omega < 1 \\ DTr_{c_i}^{p_j} & \text{if } \omega \geq 1 \end{cases} \quad (4)$$

Generally, the merging method is used to obtain the most accurate trust assessment. However, after a number of transactions, customers should analyze the quality of the received services regarding to what is expected (represented here by  $Tr_{c_i}^{p_j}$ ) and what is actually performed (so-called observed trust value  $\widehat{Tr}_{c_i}^{p_j}$ ). To this end, an adjustment trust evaluation should be periodically performed. The idea is to learn from past experiences, so that witnesses providing bad trust values, which are far from the observed one, will be removed from the list of

potential witnesses in the future. In addition, over the recent interactions, high quality providers are recognized and thus distributed to the adjacent agents. In general, using the maintenance process (for full description of algorithms, see [11]), correlated agents could increase their rate of influence to one another, which eventually would approach to a more active social network. This can be represented by the following equation:

$$\min_{c_k \in \mathcal{T}_{c_i}} |Tr_{c_k}^{p_j} - \widehat{Tr}_{c_i}^{p_j}| \quad (5)$$

### 3 Social Network Representation

To analyze our social network for service selection, many parameters described in the literature about social networks could be considered. A detailed list of such parameters are presented in [5]. For space limit, we consider only the following parameters and provide equations to compute them in our context of trust for service selection. Without loss of generality, we would like to measure the probability (likelihood) of edge creation between a customer and a provider agent. The focus of this paper is on the study of edge-by-edge evaluation of the social network in microscopic manner. We compare the network formation of different types of agents that are using different trust establishment method and use different strategies. Hence, we effectively analyze the effect of different trust models in socializing a particular agent that joins a network and seeks to increase his overall outcome (so-called utility). We basically distinguish between different models based on their strategy of network formation in agent arrival, edge arrival and interaction maintenance process (how after-interaction parameters affect the strategies that are used in the further actions of agents).

#### 3.1 Outdegree

Outdegree is a parameter for the extent to which an agent in the network conveys information regarding some other agents. Outdegree value from the customer point of view, is to what extent a customer agent knows the providers. The idea is to reflect the fact that a customer that is connected to more reliable providers has a higher outdegree than a customer linked to less reliable ones. In other words, the outdegree value reflects the extent to which an agent tries to set up and strengthen more edges connecting him to other agents. Equation 6 computes this parameter for a generalized agent  $Ag$ , that could be a customer or a provider agent, where  $\alpha Tr_{Ag}^{c_k} = Tr_{Ag}^{c_k} \cdot NI_{Ag}^{c_k} \cdot TiR_{Ag}^{c_k}$  and  $\alpha Tr_{Ag}^{p_j} = Tr_{Ag}^{p_j} \cdot NI_{Ag}^{p_j} \cdot TiR_{Ag}^{p_j}$ .

$$D_{out}(Ag) = \sum_{c_k \in \mathcal{T}_{Ag}} \alpha Tr_{Ag}^{c_k} + \sum_{p_j \in \mathcal{T}'_{Ag}} \alpha Tr_{Ag}^{p_j} \quad (6)$$

where  $\mathcal{T}'_{Ag} = \{p_j \in P \mid \exists c_i \xrightarrow{v} p_j \text{ in the social network}\}$

### 3.2 Indegree

Indegree is a parameter for the extent to which a customer in the network receives information regarding to a particular agent from some other agents. Indegree value from the customer point of view, is the extent that the agent is known by the close agents in the network. The idea is to reflect the fact that a customer that is connected to more reliable providers has a higher indegree than a customer linked to less reliable ones. Indegree value from a provide point of view, is the extent that a provider agent is popular in the social network that causes higher number of requests from the customer agents. In other words, the indegree value reflects the popularity of an agent in the sense that any agent would like to increase it and thus cares not to distract it. Although the agents that are known by many agents are supposed to be supported by them, however, regarding to their accuracy and quality of service, they may expect a portion of the adjacent agents for support. Equation 7 computes this parameter for a generalized agent  $Ag$ , that could be a customer or a provider agent.

$$D_{in}(Ag) = \sum_{c_k \in \mathcal{S}_{Ag}} \alpha Tr_{c_k}^{Ag} \quad (7)$$

where  $\mathcal{S}_{Ag} = \{c_k \in C \mid \exists c_k \xrightarrow{v}_{\bar{c}_c} Ag \text{ in the social network}\}$

### 3.3 Homophily

Homophily is a parameter for the extent to which a customer in the network chooses to interact with a provider that is known and is already evaluated (this is concept is derived from [1]). This basically raises to strengthen the correlation of adjacent agents. In the social network, agents that are known from previous interactions may tend to request for a service, which is expected to be satisfactory. This is the affect of being friend in a network. In general, it is likely that a customer agent re-selects a particular provider agent aiming to request for a new service. Thus, provider agents normally try to provide a quality service to keep their customers. The homophily of agents in the network is a factor that is not directly compared to other choices of the customer agent, that is seeking for a service. Basically it is the matter of how well-quality the provider agent would provide the new service. This means that, the customer agent's concern is to measure the probability of gaining the expected quality in the service given the fact that the provider agent is already provided a similar service to the same customer. This possibility measurement is mainly related to the indegree value of the provider agent in the sense that a provider with high indegree value is known to be popular, so there is less chance of disturbing its popularity by providing a not promised service quality. In Section 4, we analyze this effect in more details showing that the trust models with the after interaction policies could lead to a more accurate friendship evaluations.

Equation 8 computes the probability of selecting a provider with  $D_{in}(P_j)$  as indegree value. In this equation, we do not involve the trust measurement that

the customer agent  $C_i$  performs for evaluating the provider agent  $P_j$  ( $Tr_{C_i}^{P_j}$ ). The reason is that since the customer agent  $C_i$  is already in relation with the provider  $p_j$ , then based on the previous evaluation, could decide whether it worths to select this provider again. If by any chance, the previous history does not reflect the efficiency of the provider  $P_j$ , there is no point for investigating the probability of the provider's efficiency if being selected. In equation 8, the value  $\omega$  is set to be the entropy value (see equation 3) of the history between the customer agent  $C_i$  and the provider agent  $P_j$ . And the value  $\beta$  represents the coefficient set for the system inconsistency. In the trust models with after interaction strategies, this value is dynamically modified reflecting the system accuracy level, however without maintenance process, the value is set initially and remains fixed.

$$p(D_{in}(P_j)) = \frac{e^{\omega \ln(D_{in}(P_j)+1)+\beta}}{1 + e^{\omega \ln(D_{in}(P_j)+1)+\beta}} \quad (8)$$

### 3.4 Confounding

Confounding is a parameter for the extent to which a provider as an external agent influences a customer agent to request for a particular service (this is concept is derived from [1]). This influence affects some close agents in the network to set up an edge with an unknown provider under the promising conditions that the provider defines. In general, the providers that join the network, seek for the agents that are more likely to request for their providing service. In other words, when a provider agent is being activated, tries to socialize himself in the network. Thus, starting from very close customer agents, the provider agent encourages them to request for his service. To this end, the provider at the beginning acts generously in order to attract the customers and gain more popularity. Moreover, upon high quality service, the customer agents may influence their adjacent agents to request for the same service. So, the provider agent takes the outdegree value of the customer agents into account and based on the interaction span of the customer agents, provides high quality services.

In confounding factor, the probability of activating an agent with a provider agent is computed in equation 9. As it is assumed that the provider  $P_j$  is unknown to the customer  $C_i$ , so the customer agent would evaluate the social trustworthiness value of the provider. Given the fact that the trust measurement requires some information from the other adjacent agents, the customer agent takes the entropy value into account in order to partially consider the indirect trust value ( $ITr_{C_i}^{P_j}$ ) and the rest for the popularity of the provider agent. Thus, the customer  $C_i$  first evaluates the provider  $P_j$  and then considers the  $P_j$ 's indegree value together with the network inconsistency level. If the information obtained for evaluating  $P_j$  is not enough, the entropy value  $\omega$  would be high, so that mostly the trust evaluation part would be considered. This would normally cause to lower the overall probability of activation.



$$p(D_{in}(P_j)) = \omega \times ITr_{C_i}^{P_j} + (1 - \omega) \times \frac{e^{\ln(D_{in}(P_j)+1)+\beta}}{1 + e^{\ln(D_{in}(P_j)+1)+\beta}} \quad (9)$$

### 3.5 Influence

Influence is a parameter for the extent to which an agent is prompted to initiate a request caused by an adjacent agent (this is concept is derived from [1]). This could taken place in a friendship of agents that they distribute the idea of some services to be requested. When an agent in encouraged to think about a particular service from a provider, the agent may have already a set up edge with the provider, by which can evaluate the provider, or may need to set up a new edge upon which could obtain a service. This is the affect of getting encouraged by a friend in a network. In general, it is likely that a person does action because his friend is already done it. Thus it is the matter of activation of a new edge, which is set up between a customer agent and the provider agent, that is already been requested for a service by the customer agent’s adjacent agent (friend).

In the confounding factor, we mentioned that when a typical provider advertises his service to a couple of adjacent customer agents, he considers that some of the customers may propagate his quality of service to their adjacent agents, which could lead to more service requests. On the other hand, the customer agent that is being prompted to take a service produced by a particular provider, needs to evaluate both the advertising adjacent agent  $C_j$  ( $DTr_{C_i}^{C_j}$ ) and the provider itself  $P_j$  ( $ITr_{C_i}^{P_j}$ ). Equation 10 computes the influence-based probability of activation of a customer agent  $C_i$  regarding to taking the service produced by a provider agent  $P_j$ . In this equation,  $\omega_{C_j}$  is the entropy value regarding to the information  $C_i$  has and thus could rely on, and  $\omega_{P_j}$  is the entropy value that  $C_i$  has regarding to the provider  $P_j$  and would consider for interaction.

$$p(D_{in}(P_j)) = \omega_{C_j} \times DTr_{C_i}^{C_j} + (1 - \omega_{C_j}) \times \Theta \quad (10)$$

where

$$\Theta = \omega_{P_j} \times ITr_{C_i}^{P_j} + (1 - \omega_{P_j}) \times \frac{e^{\ln(D_{in}(P_j)+1)+\beta}}{1 + e^{\ln(D_{in}(P_j)+1)+\beta}}$$

## 4 Experimental Results and Related Work

In this section, we describe the implementation of proof of concept prototype. In the implemented prototype, agents are implemented as *Jadex*<sup>©TM</sup> agents. Like in [7], the testbed environment is populated with two agent types: (1) *service provider agents*; and (2) *service consumer agents*. The simulation consists of a number of consequent Runs in which agents are activated and build their private knowledge, keep interacting with one another, and enhance their overall knowledge about the environment. Depending on the agent interactions, agent may extend their connections hoping to be more socialized. However, there is

always the chance of investing on wrong agents that lead to no outcome. Here we distinguish agents by the service (or information) quality that they provide and agents do not know about that. Table 1 represents four types of the service providers we consider in our simulation: good, ordinary, bad and fickle. The first three provide the service regarding to the assigned mean value of quality with a small range of deviation. Fickle providers are more flexible as their range of service quality covers the whole possible outcomes. Upon interaction with service providers, service consumer agents obtain utilities and consequently rate the quality of the providers (for simplicity, we assume only the consumers are interconnected to the provider agents). In the simulation environment, agents are equipped with different trust models in the sense that their edge creation policies are different. In the proposed model, we try to establish a trust mechanism where an agent, firstly can maintain an effective trust assessment process and secondly, accurately updates his belief set, which reflects the other agents likely accuracy. In order to confirm the mentioned characteristics, we compare the proposed model with other trust models in two perspectives. In former comparison view, we use the agents that only perform a direct trust assessment process. We refer to this group of agents as *Direct Trust Group (DTG)*. In later overview, we use the agents that (in addition to the direct trust assessment mechanism), perform maintenance process for evaluating the consulting agents in order to increase their information accuracy. We refer to this group of agents as *Maintenance-based Trust Group (MTG)*. The reason of decomposing the proposed model to two groups is to focus on the efficiency of each model, which enables us to analyze the impact of each contribution on the accuracy of the agent in edge creation process. In order to discuss the proposed model's overall performance, we compare it with BRS <sup>1</sup> [9] and Travos <sup>2</sup> [13] trust models. These models are similar to the proposed model in the sense that they do consider other agents' suggestions while evaluating the trust of some specific agents and discard inaccurate suggestions aiming to perform best edge creation. The detailed description of these models is provided in [4]. Here we basically distinguish [10] between different models based on their strategy of network formation in agent arrival, edge arrival and interaction maintenance process (how after-interaction parameters affect the strategies that are used in the further actions of agents).

We start the discussion by the probability of selecting the providers over their different popularity values. As we discussed earlier, the indegree value of a node reflects their popularity in the social network. Thus we could conclude that the chance of selection for a particular service provider agent would be proportionally relevant to its indegree value (ordinary selection attitude). However, the trust evaluation method together with its distribution process would affect

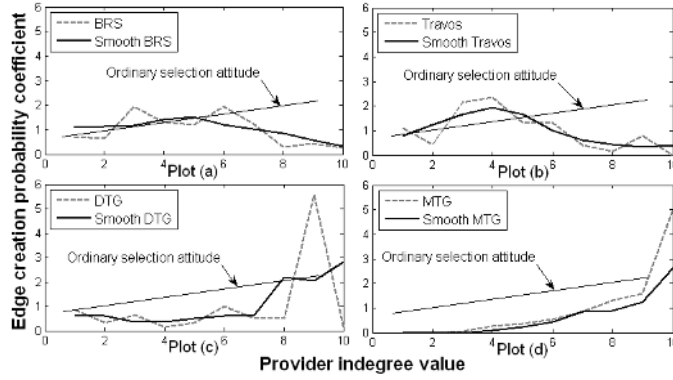
---

<sup>1</sup> BRS trust model collects the after-interaction ratings and estimates the trust using beta distribution method. This trust model ignores the ratings from such agents that deviate the most from the majority of the ratings.

<sup>2</sup> Travos trust model is similar to BRS in collecting the after-interaction ratings and estimating the trust using beta distribution method. But Travos ignores the ratings from agents that provide intermittent reports in the form of suggestions.

**Table 1.** Simulation summarization over the obtained measurements.

Service provider type	Density in the network	Utility range	Utility SD
Good	15.0%	] + 5, +10]	1.0
Ordinary	30.0%	] - 5, +5]	2.0
Bad	15.0%	] - 10, -5]	2.0
Fickle	40.0%	] - 10, +10]	-



**Fig. 1.** Probability of edge creation with provider agent vs. the provider’s indegree value.

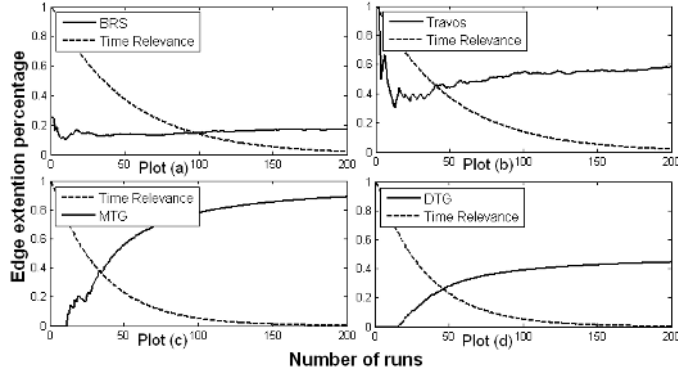
this probability of selection. Illustrated in figure 1, the BRS agents act independently of the mentioned probability as strategically, the BRS agents do not consider the popularity of the provider. Travos agents also do not consider such value, however, the probability of selection of the popular providers increase, as they take less risk of changing their behaviors and thus perform satisfactory services, which would lead to their selection. In general, because of inaccuracy detection feature of Travos agents, the percentage of selection of provider agents with high indegree value increase in a gentle manner. At some certain point, the selection of popular providers are coming down (see *plot b*). This is explained by the fact that a popular provider has large number of recommenders that provide diverse range of information to the agent, that is trying to evaluate the provider. So this diversity would cause to confusion state (the state that this system would generalize the majority of the information that is obtained and could be inaccurate), which in Travos would cause the drop of the suggestions and thus the selection would be less. The proposed model agents (*DTG* and *MTG*) follow the information propagation feature as the adjacent agents influence each other to select the high quality providers. There is a difference in the slope of selection graph in *MTG* and *DTG* models. This is explained by the fact that the *MTG* group are characterized by the maintenance process that enable them to recognize high quality provider agents and thus their accuracy in influencing adjacent agents are more than regular *DTG* agents. In general, since

the maintenance feature does not exist in *DTG* group, the customer agents loose the track of high quality provider agents, and thus the probability of selection would not increase so fast.

In general in the defined testbed, the agents that are obtaining a high quality service are encouraged to distribute their experience to other adjacent agents (influence others). This activity of agents would basically get increased over the time, or say over the age of the agent. In figure 2, we have compared the activity of different groups of agents by comparing edge extension of the agents (out-degree value). Without loss of generality, the edge extension is proportionally related to the accuracy of agent in detecting the high quality providers. In BRS model, the extension over the time is not increasing as the agent gets involved with high number of adjacent agents and would be difficult to effectively extend the social activity, so more or less would be independent of the age of the agent. Travos and *DTG* models are increasing, however relatively with small slope. In *MTG* group, because of the maintenance process the agents would be encouraged to initiate a request to high quality service providers and thus extend their activity. In this graph, the slope is relatively large as over the time, the agent could manage to categorize the providers that could possibly act beneficially for the agent, and thus would enlarge his activity area. In figure 2, the second line represents how fast the agents would drop the previous data and use the recent data for their analysis. This dropping factor is also relevant to how active an agent is and thus, to what extent there would be available resource that agents could drop obsolete data. *DTG* and *MTG* group use the same dropping feature ( $TiR(\Delta t_{Ag_a}^{Ag_b})$ ), which is derived in equation 11. Variable  $\lambda$  is an application-dependent coefficient. In some applications, recent interactions are more desirable to be considered ( $\lambda$  is set to relatively large number). In contrast, in some other applications, even the old interactions are still valuable source of information. In that case, a relatively smaller value to  $\lambda$  is used.

$$TiR(\Delta t_{Ag_a}^{Ag_b}) = e^{-\lambda \ln(\Delta t_{Ag_a}^{Ag_b})} \quad \lambda \geq 0 \quad (11)$$

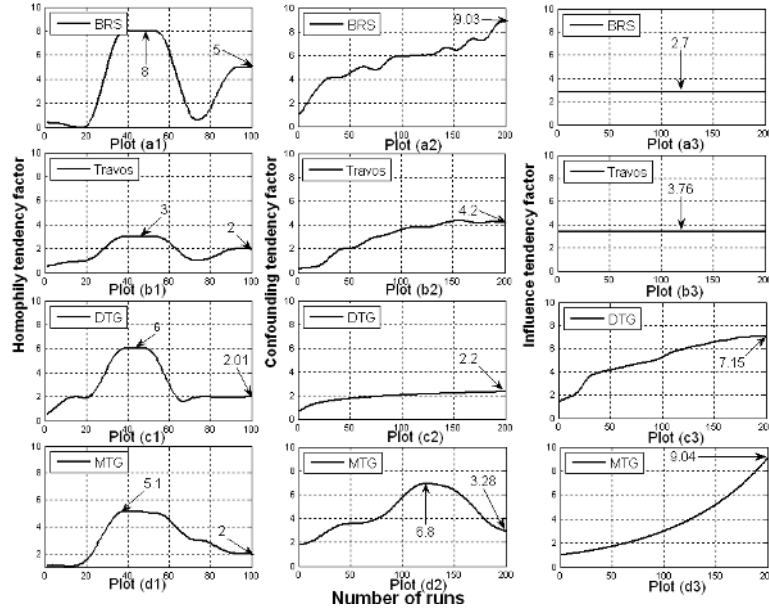
We would like to go further into the details of the selection history in terms of the microscopic social network affects (homophily, confounding, and influence) and illustrate them in figure 3. in this section, we observe the diverse impacts of homophily, confounding and influence features on each group in the sense that we would capture their edge creation reasons. Note that the edge creation is not the important issue, however, the concern is to extend to the agents that are known to be trustworthy. Therefore, we elaborate the overall outcome of different agents at the following. The homophily aspect would be caused by the friendship relation of the agents that have history interaction between them. This is a very general case in the sense that consumer agents over the time would get to know and select the provider agents. If the interacted service is satisfactory for the agent, then the consumer agent may re-select the same provider agent in some future. BRS agents are the ones that mostly rely on the homophily affect in the sense that they keep the history of the interaction in order to re-evaluate the provider agent. The providers that remain trustworthy would be selected over



**Fig. 2.** Agent edge extension vs. the agents age.

the time. As it is clear from plot *a1*, once the providers change their policies, the selection of them would be affected so fast, as the BRS agents recognize that they should start seeking for the appropriate friends. Travos agents also rely on the previous history and re-select the previously interacted service providers (see plot *b1*). However, over the time the reports regarding to the accuracy of the providers would be divergent, which would lead to refuse the selection. The same reason is the case for *DTG* and *MTG* group (shown in plots *c1* and *d1*). These agents to some extent rely on the previous history and select the providers. After some certain time, these agents also recognize the inconsistency in the evaluation process of the history interacted providers. Overall, *DTG* and *MTG* agents evaluate the providers in a very accurate manner. The accuracy that Travos, *DTG* and *MTG* agents have cause the decremented manner after some certain time.

Confounding factor reflects the extent to which the provider agents advertise their service to the consumers (could be new or previously serviced ones). This feature also affects BRS group, as they start evaluating the advertising provider, and thus extend their activation area. Plot *a2* indicates that the BRS group are easy to involve in interaction with the advertising provider agent. Travos agents act in the same way as the provider agents could induce them to take their service. However, Travos agents are considering this case less, because they investigate the previous reports related to the advertising provider and doubt on the inconsistent ones (see plot *b2*). In general, the BRS and Travos agents accept the confounding-related interactions over the time, and thus their graph has an increasing manner. But in *DTG* and specially *MTG*, the agents would not accept this service all the time, as over the time, once the network inconsistency level increases, these agents would have confusion in accepting the confounding-related affect caused by unknown service providers (see plots *c2* and *d2*). *MTG* agents would accept this option from the providers, but since they are equipped with a maintenance process, they would distribute the performance of

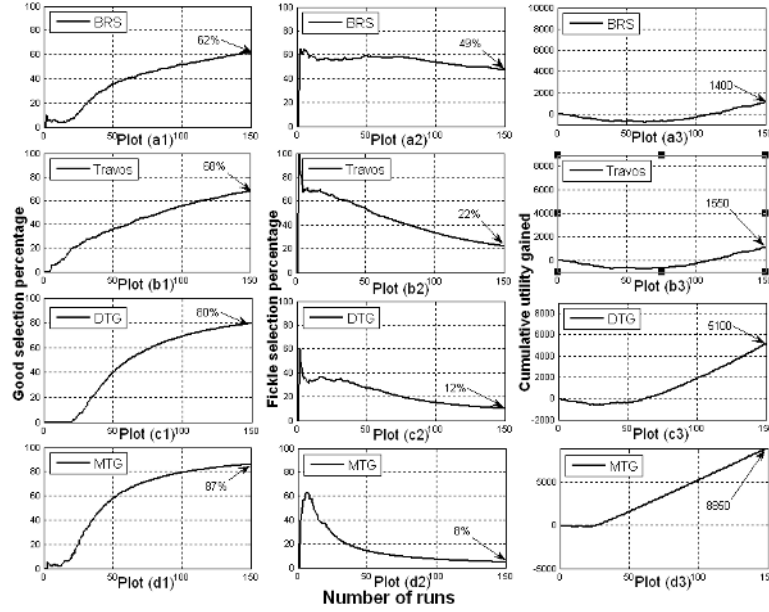


**Fig. 3.** Overall comparison of the proposed model with BRS and Travos in terms of (a) Homophily; (b) confounding; and (c) influence factors.

the providers to the adjacent agents, which would lead them to get to know the network faster than the other models. This would let the *MTG* agents to select the best providers, and thus would drop the request from most of the unknown agents while they are already in a good accuracy level.

Influence factor is mostly used by active agents, while they obtain service and tend to distribute the efficiency of the interaction to the adjacent agents. Since BRS agents independently select the providers, the influence is not a factor for these agents (plot *a3*). Travos agents would act almost independently, however the Travos agents are encouraged by the reports they obtain for the evaluation of a particular provider agent (plot *b3*). *DTG* group would be encouraged with the same factor as Travos agents. Upon evaluating provides, the *DTG* agents would consider the reports obtained from adjacent agents and recognize outstanding service provided by the provider that is just served an adjacent agent (see plot *c3*). The influence-related interactions are mostly initiated among *MTG* group, shown in plot *d3*. This is explained by the fact that the *MTG* group are equipped with maintenance feature, which enables them to reason about the accuracy and efficiency of the obtained services and propagate the information to the adjacent information.

Considering all the involved features, at the end we compare the models in general perspective, starting good provider selection efficiency. In such a biased environment, the number of good providers are comparatively low. Therefore,



**Fig. 4.** Overall comparison of the proposed model with BRS and in terms of (a) good selection percentage; (b) fickle selection percentage; and (c) cumulative utility gained.

the agents need to perform an accurate trust assessment to recognize the best providers. As it is clear from the Figures 4, plots *a1*, *b1*, and *c1*, *DTG* agents function better than other models (Travos and BRS). The reason is that in this model, agents are assessing the credibility of the providers using other agents suggestions depending on their credibility and to what extent they know the provider. Afterwards these agents rate the provider, which would be distributed to other agents upon their request (relatively in plots *a2*, *b2*, and *c2* the comparison of fickle selection percentage, and in *a2*, *b2*, and *c2*, the gained cumulative utility is shown). Not excluding the fact that *DTG* agents are considering partial ratings for consulting agents, we state that they weakly function when the environment contains agents that do not truthfully reveal their believes. *MTG* agents in addition to the direct trust assessment, provide incentives for consulting agents, which encourages them to effectively provide the information aiming to gain more utility. Plot *d1* shows that *MTG* agents outperform other models in best provider selection. This is expressed by the fact that *MTG* agents recognize the best providers ensuring that the best selected provider would provide the highest utility. Relatively plot *d2* shows an outperform in fickle selection and consequently higher cumulative utility in plot *d3*.

In BRS model, the trustor agent in the assessment process uses beta distribution method and discards the ratings that deviate the most from the majority of the ratings. Concerning this, BRS is comparatively a static trust method,

which causes a low-efficient performance in very dynamic environment. In general, if a BRS agent decides to evaluate an agent that he is not acquainted with, he considers the majority of ratings, which are supposed to be truthfully revealed about the trustee agent. In such a case that the trustee agent has just changed his strategy, the trustor agent would loose in trust assessment and does not verify the accuracy of the gained information. Therefore, as illustrated in figure 4, plots *a1*, the BRS agents would have less percentage of good providers selection, relatively higher percentage of fickle providers selection (plot *a2*), and consequently lower gained cumulative utility (plot *a3*).

Travos [13] trust model is similar to BRS in using beta distribution to estimate the trust based on the previous interactions. Travos model also does not have partial rating. Hence, the trustor agent merges his own experience with suggestions from other agents. However, unlike BRS model, Travos filters the surrounding agents that are fluctuating in their reports about a specific trustee agent. To some extent, this feature would cause a partial suggestion consideration and thus, Travos agents would adapt faster comparing to BRS agents. Rates concerning the good and fickle selection percentage shown in figures 4, plots *b1* and *b2* reflect higher efficiency of Travos compared to BRS. However, Travos model considers that agents do not change their behavior towards the elapsing time. These missing assumptions affect the accuracy of trust estimation in a very biased environment (lower gained cumulative utility in plot *b3*).

## 5 Conclusion

The contribution of this paper is the detailed investigation of a trust-based multi-agent architecture in edge creation and correlation formation in social network. The established trust is provided by the proposed framework, that is briefly explained here. The trust assessment procedure is based on integrating suggestion of consulting agents, objectively enhancing the accuracy of agents to make use of the information communicated to them. The surveillance over the surrounding environment, makes distributed agents eager to extend their activity area by interacting to high quality agents. In the proposed framework, maintenance process considers the communicated information to judge the accuracy of the consulting agents in the previous trust evaluation process. The ex-interacted analysis, makes the agents to propagate the recent and accurate information to their adjacent agents, which is considered as homophily and influence factors in edge creation analysis.

Our model has the advantage of being computationally efficient as it takes into account the important factors involved in extending the activity zone of agents. Moreover, we have done a detailed empirical analysis over the edge creation and behavior of agents over their age, while they are equipped with different trust mechanism protocols. The proposed mechanism efficiency is compared with other related models to prove the capabilities of the proposed model. Our plan for future work is to advance the assessment model to enhance the model efficiency. In the maintenance process we need to elaborate more on the optimization part,



trying to formulate it in the sense to be adaptable to diverse situations. Finally, we plan to maintain more detailed analysis in comparison with other models to capture more results reflecting the proposed model capabilities.

## References

1. A. Anagnostopoulos, R. Kumar, and M. Mahdian. Influence and correlation in social networks. Proceedings of the 14'th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Las Vegas, USA, KDD2008.
2. L. Backstrom, D. Huttenlocher, J. Kleinberg, and X. Lan. Group formation in large scale networks: membership, growth, and evolution. Proceedings of the 12'th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining , pp.44-54, KDD2006.
3. J. Bentahar, B. Khosravifar, and M. Gomrokchi. Social network-based trust for agent-based services. Proceedings of the 5'th IEEE International Symposium on Web and Mobile Information Services, Bradford, UK, WAMIS2009 (in press).
4. J. Bentahar, B. Khosravifar. Using trustworthy and referee agents to secure multi-agent systems. Proceedings of the 5'th IEEE International Conference on Information Technology: New Generations, pp. 477-482, Las Vegas, USA, ITNG 2008.
5. V. Buskens. Social networks and trust. Kluwer Academic Publishers. 2002.
6. N. A. Christakis, and J. H. Fowler. The spread of obesity in a large social network over 32 Years. The new England Journal of Medicine, 357(4):370-379, 2007.
7. T. Dong-Huynh, N.R. Jennings and N.R. Shadbolt. Fire: an integrated trust and reputation model for open multi-agent systems. Journal of Autonomous Agents and Multi-Agent Systems 13(2):119-154, 2006.
8. T. Dong-Huynh, N.R. Jennings and N.R. Shadbolt. Certified reputation: how an agent can trust a stranger. In Proceedings of the 5'th International Joint Conference on Autonomous Agents and Multiagent Systems, pp. 1217-1224, Japan, 2006.
9. A. Jesang and R. Ismail. The beta reputation system. 15'th Bled Electronic Commerce Conference e-Reality: Constructing the e-Economy, June 2002.
10. J. Leskovec , L. Backstrom , R. Kumar, and A. Tomkins. Microscopic evolution of social networks. In Proceedings of the 14'th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 462-470, 2008.
11. B. Khosravifar, M. Gomrokchi, J. Bentahar, and Ph. Thiran. A maintenance-based trust for Open multi-agent systems. Accepted in the 8'th International Conference on Autonomous Agents and Multiagent Systems, AAMAS, Budapest, Hungary 2009.
12. J. Sabater, M. Paolucci and R. Conte, RePage: REPutation and ImAGE among limited autonomous partners. Journal of Artificial Societies and Social Simulation 9(2), 2006.
13. W. T. Teacy, J. Patel, N.R. Jennings, and M. Luck. Travos: trust and reputation in the context of inaccurate information sources. Autonomous Agents and Multi-Agent Systems, 12(2):183-198, 2006.
14. Y. Wang, and M.P. Singh. Formal trust model for multiagent systems. Proceedings of the 20'th International Joint Conference on Artificial Intelligence (IJCAI), pp. 1551-1556, 2007.
15. L. Xiong and L. Liu. PeerTrust: supporting reputation-based trust for peer-to-peer electronic communities. Journal of IEEE Transactions on Knowledge and Data Engineering, 16(7):843-857, 2004.

# Computing Utility from Weighted Description Logic Preference Formulas

Azzurra Ragone<sup>1</sup>, Tommaso Di Noia<sup>1</sup>, Francesco M. Donini<sup>2</sup>, Eugenio Di Sciascio<sup>1</sup>,  
Michael P. Wellman<sup>3</sup>

<sup>1</sup> SisInfLab, Politecnico di Bari, Bari, Italy  
{a.ragone,t.dinoia,disciascio}@poliba.it

<sup>2</sup> Università della Tuscia, Viterbo, Italy  
donini@unitus.it

<sup>3</sup> Artificial Intelligence Laboratory–University of Michigan, Ann Arbor, USA  
wellman@umich.edu

**Abstract.** We propose a framework to compute the utility of a proposal w.r.t. a preference set in a negotiation process. In particular, we refer to preferences expressed as weighted formulas in a decidable fragment of First-order Logic. Specifically, we refer to expressive Description Logics (DL) endowed with disjunction.

DLs offer expressivity advantages over propositional representations, and allow us to relax the often unrealistic assumption of additive independence among attributes. We provide suitable definitions of the problem and present algorithms to compute utility in our setting. We also study complexity issues of our approach and demonstrate its usefulness with a running example in a multiattribute negotiation scenario.

## 1 Introduction

Effective and expressive specification of preferences is a particularly challenging research problem in knowledge representation. Preference representation is essential, for example, to instruct a software agent to act on behalf of the objectives of a human being. One common approach to preference representation appeals to *multiattribute utility theory* [1], which concerns the construction of *utility functions* mapping vectors of *attributes* to real values. Given that the size of a multiattribute domain is exponential in the number of attributes, applications typically exploit independence relations among the attributes, in the most extreme case to assume that all attributes are additively independent, so that the multiattribute utility function is a weighted sum of single-attribute utility functions.

However, most real-world domains pose significant preferential dependencies, ruled out by the fully additive model. For example, referring to the desktop computer realm, we could not with an additive value function capture the fact that the value of some combination of attributes it is not simple the sum of the single attribute values. Indeed, some operating systems performs very poorly without a minimum amount of memory, therefore the value (utility) given to a specific operating system will depend on the amount of memory available.

Some recent approaches support relaxation of the fully additive assumption, for example by providing generalized versions [2] or exploiting graphical models of dependence structure [3–5], while remaining within the multiattribute framework.

Logical languages likewise provide a means to express interdependencies, but unlike multiattribute formulations they do not necessarily require that we explicitly decompose the domain into an orthogonal set of attributes. Furthermore, logical languages support integration of preference knowledge with domain knowledge modeled through an ontology. Using an ontology, indeed, it is possible to model relations among attributes in the domain (*e.g.*, *a Centrino is an Intel processor with a 32-bit CPU*), as well as the fact that some combination of features may be infeasible (therefore of minimal or undefined preference) due to constraints in the ontology itself (*e.g.*, *a Centrino processor is not compatible with a processor with a 64-bit architecture*).

In decision making problems, preferences are expressed over a set of possible alternatives, in order to rank them. In many cases, such as *e.g.*, bilateral negotiation, auctions, resource allocation, it is important to compute a utility value for, respectively, an agreement, an offer, an allocation w.r.t. the set of preferences expressed by the agent. If preferences are expressed using Propositional Logic, then the utility can be computed considering a particular propositional model (agreement, offer, allocation), taking into account formulas satisfied by that model.

While for Propositional Logic it is possible to refer directly to models (interpretations) in order to compute utility, this computation for First-order Logic (FOL) is less straightforward, as the number of possible models is infinite.

The main contribution of this paper is an approach that, given a set of preferences, represented as weighted DL formulas w.r.t. a shared ontology, computes the utility of a formula (agreement, offer, allocation, etc.) based on its possible models (interpretations). To our knowledge, the only prior method proposed in the literature for this problem is subsumption, which has some limitations, as we show in Section 4.

We point out that even though the results we show in this paper can be easily applied to whatever decidable logic with a model-theoretic semantics, we ground our approach on DLs because of their importance in the development of the Semantic Web.

The remainder of the paper proceeds as follows. First, we introduce Description Logics, then we give a brief overview of the problem of preference representation in the field of logic languages. In Section 4, we first introduce the problem of computing utility of a concept w.r.t. a preference set, showing how, sometime, subsumption leads to counterintuitive results. Then we analyze some complexity issues. In Section 5 we illustrate our framework for the computation of utility for a set of weighted DL formulas with the help of a running example. Finally, we discuss some considerations about the computational properties of the framework. Conclusion closes the paper.

## 2 Description Logic Basics

Description logics (DLs) are a family of formalisms well-established in the field of knowledge representation. Readers familiar with DLs may safely skim this section, attending mainly to the notation and examples. Those interested in learning more may refer to the *Description Logic Handbook* [6] for a much more comprehensive treatment.

The basic syntax elements of Description Logics are *concept names*, *properties*, and *individuals*. Concept names stand for sets of objects in the domain<sup>4</sup> (Windows, Intel, LCDMonitor), and properties link (sets of) such objects (hasOS, hasCPU, hasMonitor). Individuals correspond to special named elements belonging to concepts (HP\_Pavilion, Apple\_iMac). When we do not use proper names, we denote concepts by symbols  $A, B, C, D, \dots, \top, \perp$ .

Description logics are usually endowed with a model-theoretic formal semantics. A semantic *interpretation* is a pair  $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ , where  $\Delta$  represents the *domain* and  $\cdot^{\mathcal{I}}$  is the *interpretation function*. This function maps every concept to a subset of  $\Delta$ , and every property to a subset of  $\Delta \times \Delta$ . Then, given a concept name  $CN$  and a property name  $R$  we have:  $CN^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}}$  and  $R^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$ . The symbols  $\top$  and  $\perp$  are used to represent the most generic concept and the most specific concept respectively. Hence their formal semantics correspond to  $\top^{\mathcal{I}} = \Delta^{\mathcal{I}}$  and  $\perp^{\mathcal{I}} = \emptyset$ .

Properties and concept names can be combined using *existential role quantification*. For example,  $\text{PC} \sqcap \exists \text{netSupport.WiFi}$  describes the set of PCs supporting a wireless connection. Similarly, we can use *universal role quantification*, as in  $\text{PC} \sqcap \forall \text{hasCPU.AMD}$ , to describe the set of PCs having only AMD processors on board. The formal semantics of universal and existential quantification is as follows:

$$\begin{aligned} \exists R.C &= \{x \in \Delta^{\mathcal{I}} \mid \exists y, (x, y) \in R^{\mathcal{I}} \wedge y \in C^{\mathcal{I}}\} \\ \forall R.C &= \{x \in \Delta^{\mathcal{I}} \mid \forall y, (x, y) \in R^{\mathcal{I}} \rightarrow y \in C^{\mathcal{I}}\} \end{aligned}$$

Concept expressions can be written using *constructors* to write concept and property *expressions*. Based on the set of allowed constructors we can distinguish different description logics. Essentially every DL allows one to form a *conjunction* of concepts, usually denoted as  $\sqcap$ ; some DLs include also disjunction  $\sqcup$  and complement  $\neg$  to close concept expressions under boolean operations.

$$\begin{aligned} (C \sqcap D)^{\mathcal{I}} &= C^{\mathcal{I}} \cap D^{\mathcal{I}} \\ (C \sqcup D)^{\mathcal{I}} &= C^{\mathcal{I}} \cup D^{\mathcal{I}} \\ (\neg C)^{\mathcal{I}} &= \Delta^{\mathcal{I}} \setminus C^{\mathcal{I}} \end{aligned}$$

Constructs involving number restriction enable us to define concepts in terms of the numbers of roles with specified properties. For example,  $\text{Mac} \sqcap (\geq 4 \text{hasUSBport})$  describes a Macintosh PC with at least four USB ports. An example of concrete domain restrictions appears in the expression  $\text{PC} \sqcap \exists \text{hasRam} . (\geq_2 \text{GB})$ , which describes a PC with at least 2 GB of memory. Notice that while properties, such as  $\text{hasUSBport}$ , are mapped to a subset of  $\Delta \times \Delta$ , concrete properties, such as  $\text{GB}$  are mapped to a subset  $\Delta \times D$  where  $D$  is a concrete domain.

$$\begin{aligned} (\geq n R)^{\mathcal{I}} &= \{x \in \Delta^{\mathcal{I}} \mid |\{b \in \Delta^{\mathcal{I}} \mid (a, b) \in R^{\mathcal{I}}\}| \geq n\} \\ (\leq m R)^{\mathcal{I}} &= \{x \in \Delta^{\mathcal{I}} \mid |\{b \in \Delta^{\mathcal{I}} \mid (a, b) \in R^{\mathcal{I}}\}| \leq m\} \\ (\leq_k R)^{\mathcal{I}} &= \{x \in \Delta^{\mathcal{I}} \mid R^{\mathcal{I}}(x) \leq k\} \\ (\geq_h R)^{\mathcal{I}} &= \{x \in \Delta^{\mathcal{I}} \mid R^{\mathcal{I}}(x) \geq h\} \end{aligned}$$

<sup>4</sup> We illustrate the main points here (and throughout the paper) using the domain of desktop computers.

In general, the expressiveness of a DL depends on the type of constructors allowed.

Given a generic concept  $C$ , we use the notation  $m \models C$  to say that  $m$  is a model of  $c$ , i.e.,  $m \in C^{\mathcal{I}}$ .

In order to formally represent domain knowledge and constraints operating among elements of the domain, we employ a set of background axioms, that is, an *ontology*. Formally, ontology  $\mathcal{T}$  (for Terminology) comprises axioms of the form  $D \sqsubseteq C$ , where  $D$  and  $C$  are well-formed formulas in the adopted DL, and  $R \sqsubseteq S$ , where both  $R$  and  $S$  are properties. The formal semantics of such axioms is:  $(C \sqsubseteq D)^{\mathcal{I}} = C^{\mathcal{I}} \subseteq D^{\mathcal{I}}$ ,  $(R \sqsubseteq S)^{\mathcal{I}} = R^{\mathcal{I}} \subseteq S^{\mathcal{I}}$ . We write  $C \equiv D$  to represent both  $C \sqsubseteq D$  and  $D \sqsubseteq C$ . In the rest of the paper we refer to the Ontology  $\mathcal{T}$  depicted in Figure 1.

```

DesktopComputer  $\sqsubseteq$   $\exists$ hasCPU  $\sqcap$   $\exists$ hasRam

CPUArchitecture  $\equiv$  64BitCPU  $\sqcup$  32BitCPU

32BitCPU  $\sqsubseteq$   $\neg$ 64BitCPU

Intel  $\sqcup$  AMD  $\sqsubseteq$  CPU

Intel  $\sqsubseteq$   $\neg$ AMD

hasUSBport  $\sqsubseteq$  hasPeripheralPort

Athlon64  $\sqsubseteq$  AMD  $\sqcap$   $\exists$ arch  $\sqcap$   $\forall$ arch.64BitCPU

Centrino  $\sqsubseteq$  Intel  $\sqcap$   $\exists$ arch  $\sqcap$   $\forall$ arch.32BitCPU

 $\exists$ hasCPU.Centrino  $\sqsubseteq$   $\exists$ netSupport.WiFi

IntelDuo  $\sqsubseteq$  Intel  $\sqcap$   $\exists$ arch  $\sqcap$   $\forall$ arch.32BitCPU

Sempron  $\sqsubseteq$  AMD  $\sqcap$   $\exists$ arch  $\sqcap$   $\forall$ arch.64BitCPU

```

**Fig. 1.** Reference ontology.

Description Logics are considered to be highly expressive representation languages, corresponding to decidable fragments of first-order logic. Reasoning systems based on DLs generally provide at least two basic inference services: *satisfiability* and *subsumption*.

**Satisfiability:** a concept expression  $C$  is satisfiable w.r.t. an ontology  $\mathcal{T}$  when  $\mathcal{T} \not\models C \sqsubseteq \perp$ , or equivalently  $C \not\sqsubseteq_{\mathcal{T}} \perp$ ;

**Subsumption:** a concept expression  $C$  is subsumed by a concept expression  $D$  w.r.t.  $\mathcal{T}$  when  $\mathcal{T} \models C \sqsubseteq D$ , or equivalently  $C \sqsubseteq_{\mathcal{T}} D$ .

In our setting satisfiability is useful, for example, to catch inconsistencies among a buyer’s expressed preferences or, analogously, among sellers’ offered configurations. On the other hand, subsumption can be employed to verify if a particular seller’s offer satisfies one or more buyer’s preferences.

### 3 Preference Representation Using Description Logics

The problem of preference representation deals with the expression and evaluation of preferences over a set of different alternatives (outcomes). This problem can be challenging even for a small set of alternatives, involving a moderate number of features, as the user has to evaluate all possible configurations of feature values in the domain.

In this work, we deal with this problem by a combination of expressive language, to facilitate preference specification, and preference structure exploitation, justified by multiattribute utility theory.

Several approaches to negotiation have exploited logic languages in order to express preferences, most of them using propositional logic [7–9], however only few of the approaches proposed in the literature have explored the possibility to use also an ontology to model relations among attributes [10] or the use of more expressive logics as DLs [11, 12]. Lukasiewicz and Schellhase [13] propose a framework to model conditional preferences in DLs for matchmaking. In their framework they refer to set of concepts and the formal semantics of implication is defined in terms of set membership. Such a formulation well suits their target matchmaking task. Indeed, they are not interested in computing a utility value for a concept, e.g. an agreement, but they focus on ranking a set of results w.r.t. a query.

We point out the importance to refer to a background knowledge, *i.e.*, having an ontology  $\mathcal{T}$ , in order to model not only interdependencies among attributes in preference statements, but also to model inner relations among attributes that cannot be disregarded *e.g.*, is-a, disjoint or equivalence relations. In order to lay out the importance of an ontology and why we cannot abstract from it, we use a simple example involving one perspective buyer and two sellers.

*Example 1.* Let us suppose the buyer has among her preferences:

$$P = \exists \text{hasCPU} . (\text{AMD} \sqcap \exists \text{arch} \sqcap \forall \text{arch} . 64\text{BitCPU}) \sqcap \exists \text{hasRam} . (\geq_2 \text{ GB})$$

(PC with a 64-bit AMD processor and at least 2 GB of memory)

and there are two sellers, A and B, that can provide the respective configurations:

$$A = \exists \text{hasCPU} . \text{Athlon64} \sqcap \text{hasRam} (=_2 \text{ GB})$$

$$B = \exists \text{hasCPU} . \text{Centrino} \sqcap \text{hasRam} (=_1 \text{ GB})$$

If we refer to the ontology  $\mathcal{T}$  in Figure 1 we can state that seller A can satisfy the preference expressed by the buyer —from  $\mathcal{T}$  we know that Athlon64 is a 64-bit AMD processor. Conversely, seller B cannot satisfy buyer’s preference, because Centrino is not a 64-bit AMD processor.  $\triangle$

We extend the well-known approach of weighted propositional formulas [7–9], representing preferences as DL formulas, where at each formula we associate a value  $v$  representing the relative importance of that formula.

**Definition 1.** Let  $\mathcal{T}$  be an ontology in a DL. A **Preference** is a pair  $\phi = \langle P, v \rangle$  where  $P$  is a concept such that  $P \not\sqsubseteq_{\mathcal{T}} \perp$  and  $v$  is a real number assigning a worth to  $P$ . We call a set  $\mathcal{P}$  of preferences a **Preference Set**.

*Example 2 (Desktop Computer Negotiation).* Imagine a negotiation setting where buyer and seller are negotiating on the characteristics of a *desktop computer*. The buyer will have some preferences, while the seller will have some different configurations to offer to the buyer in order to satisfy her preferences. Let us hence suppose the buyer is looking for a desktop PC endowed with an AMD CPU. Otherwise, if the desktop PC has an Intel CPU, it should only be a Centrino one. The buyer also wants a desktop PC supporting wireless connection. Following Definition 1 the buyer’s Preference set is  $\mathcal{P} = \{\langle P_1, v_1 \rangle, \langle P_2, v_2 \rangle, \langle P_3, v_3 \rangle\}$ , with:

$$\begin{aligned} P_1 &= \forall \text{hasCPU.Centrino} \sqcap \exists \text{hasCPU} \\ P_2 &= \exists \text{hasCPU.AMD} \\ P_3 &= \exists \text{netSupport.WiFi} \end{aligned}$$

On the other side, the seller could offer a *Desktop computer supporting either a wireless connection or an AMD CPU, specifically a Sempron one, and he does not have a desktop PC endowed with a Centrino CPU*.

$$A = \text{DesktopComputer} \sqcap \neg \exists \text{hasCPU.Centrino} \sqcap (\exists \text{netSupport.WiFi} \sqcup \forall \text{hasCPU.Sempron})$$

Therefore, given a preference set  $\mathcal{P}$  and a proposal  $A$ , how to evaluate the utility of this proposal w.r.t. buyer’s preferences? Intuitively, the utility value should be the sum of the value  $v_i$  of the preferences satisfied by the seller’s proposal. Next sections will address this problem, showing a computation method for weighted DL-formulas.

△

## 4 Utility

Weighted formulas have been introduced for propositional logic to assign a utility value to a propositional model representing *e.g.*, the final agreement. The computation of the model and its corresponding utility is quite easy since in propositional logic we deal with a finite set of models. Following Chevaleyre et al. [9] utility is computed as:

$$\sum \{v \mid \langle P, v \rangle \in \mathcal{P} \text{ and } m \models P\}$$

Where  $\mathcal{P}$  is a propositional Preference Set,  $m$  is a propositional interpretation (model) *e.g.*, representing the final agreement.<sup>5</sup> We call this approach *model-based*. Less straightforward is the case of more expressive logics. Some attempts in this direction have been

<sup>5</sup>  $\sum\{\cdot\}$  indicates the summation over all the elements in the set  $\{\cdot\}$ .

made by Ragone et al. [11, 12] adapting the weighted formulas framework to Description Logics. There, they do not consider models as final agreements but formulas, and the utility value is computed as:

$$\sum \{v \mid \langle P, v \rangle \in \mathcal{P} \text{ and } A \sqsubseteq_{\mathcal{T}} P\}$$

Where  $\mathcal{P}$  is a DL preference set,  $\mathcal{T}$  is a DL ontology and  $A$  is a concept *e.g.*, representing a proposal in a negotiation process. We call this approach *implication-based*. Although very simple to compute and immediate, this basic approach may lead to counter-intuitive examples when dealing with logic languages allowing disjunction in the final agreement.

*Example 3.* Consider the following preference set  $\mathcal{P}$  (here for the sake of clarity we do not consider the ontology  $\mathcal{T}$ )

$$\begin{aligned}\phi_1 &= \langle A_1, v_1 \rangle \\ \phi_2 &= \langle A_2, v_2 \rangle \\ \phi_3 &= \langle A_3, v_3 \rangle\end{aligned}$$

and a concept  $A$ .

$$A = (A_1 \sqcup A_3) \sqcap A_2$$

△

In Example 3, following an *implication-based* approach the final utility is

$$u^{impl}(A) = v_2$$

Indeed, only  $A \sqsubseteq P_2$  holds.

On the other hand, if we use a *model-based* approach we can say that the final utility value is:

$$u^{model}(A) \in \{v_1 + v_2, v_2 + v_3, v_1 + v_2 + v_3\}$$

If we consider models  $m$  of  $A$  we may have that only one of the conditions below holds:

$$\begin{aligned}m &\models A_1 \sqcap \neg A_3 \sqcap A_2 \\ m &\models \neg A_1 \sqcap A_3 \sqcap A_2 \\ m &\models A_1 \sqcap A_3 \sqcap A_2\end{aligned}$$

Using a model-based approach, if we want to be as conservative as possible we may consider:

$$u^{model}(A) = \min\{v_1 + v_2, v_2 + v_3, v_1 + v_2 + v_3\}$$

Conversely, in the most optimistic case we consider:

$$u^{model}(A) = \max\{v_1 + v_2, v_2 + v_3, v_1 + v_2 + v_3\}$$

In the rest of the paper we will refer to the conservative situation but all the results can be easily adapted to the optimistic case. Consequently, we give a definition of Minimal Model and of its corresponding Minimal Utility Value.



**Definition 2 (Minimal Models – Minimal Utility Value).** Given an ontology  $\mathcal{T}$ , a concept  $A$ , such that  $\mathcal{T} \not\models A \sqsubseteq \perp$ , and a set of preferences  $\mathcal{P}$ , a **Minimal Model** is a model  $m$  such that:

1.  $m \models \{A\} \cup \mathcal{T}$
2. the value  $u^c(A) = \sum\{v \mid \langle P, v \rangle \in \mathcal{P} \text{ and } m \models P\}$  is minimal

We call  $u^c(A)$  a **Minimal Utility Value** for  $A$  w.r.t. to  $\mathcal{P}$ .

#### 4.1 Complexity

In this section, we give some consequences on the complexity of computing the utility of a formula  $A$ , when utility is attached to models. In what follows, we abstract from the particular logic language  $\mathcal{L}$ , which gives our results the maximum attainable generality.

**Lemma 1.** Given two concepts  $A, P \in \mathcal{L}$ , and an ontology  $\mathcal{T}$ , we have  $A \sqsubseteq_{\mathcal{T}} P$  iff there exists a minimal model assigning value  $v$  to  $A$  when preferences are  $\mathcal{P} = \{\langle P, v \rangle\}$ .

*Proof.* Given a singleton set of preferences  $\mathcal{P} = \{\langle P, v \rangle\}$ , then  $u^c(A)$  can be equal to either 0, or  $v$ . Now if there exists a minimal model with value 0, then such a model is a model of  $\mathcal{T}$  and  $A$ , and it must be not a model of  $P$ ; hence, when the minimal model has utility 0,  $\mathcal{T} \not\models A \sqsubseteq P$ . On the other hand, there exists a minimal model assigning utility value  $v$  to  $A$ , then every model of  $\mathcal{T}$  and  $A$  is also a model of  $P$ . But this is just the condition expressing semantically that  $A \sqsubseteq_{\mathcal{T}} P$ .  $\square$

A consequence of the above lemma is that computing a minimal model is at least as hard as computing subsumption in a DL  $\mathcal{L}$ ; below we can even generalize this result to logical implication.

**Theorem 1.** Given a language  $\mathcal{L}$  in which deciding logical implication is  $\mathcal{C}$ -hard, deciding the existence of a minimal model with a given utility value is  $\mathcal{C}$ -hard, too.

We observe that the result is the same (by definition) when utilities are directly assigned to *formulas*, as done e.g., by Ragone et al. [11].

We now move to upper bounds on computing utilities over formulas by minimal models. We first assess the upper bound of the decision problem corresponding to computing the utility of a concept.

**Theorem 2.** Let  $\mathcal{L}$  be a language in which the satisfiability problem belongs to the complexity class  $\mathcal{C}$ , and such that  $\text{NP} \subseteq \mathcal{C}$ ; moreover, let  $v$  be a positive real number,  $\mathcal{P}$  be a set of preferences,  $\mathcal{T}$  be a Terminology and  $A$  a concept, all expressed in  $\mathcal{L}$ . Then, deciding whether  $u^c(A) < v$  is a problem in  $\mathcal{C}$ .

*Proof.* Let  $\mathcal{P} = \{\langle P_1, v_1 \rangle, \dots, \langle P_n, v_n \rangle\}$ . Then, for each integer  $m$  between 0 and  $2^n - 1$ , let  $(m)_2 = b_1 b_2 \dots b_n$  be the binary representation of  $m$ , and let  $D_m$  be the concept defined as  $A \sqcap B_1 \sqcap \dots \sqcap B_n$ , where for each  $i = 1, \dots, n$ , if  $b_i = 0$  then  $B_i = \neg P_i$ , else  $B_i = P_i$ . Intuitively, the  $i$ -th bit of  $(m)_2$  decides whether  $P_i$  appears positively or negatively in  $D_m$ . Now let  $\mathcal{S} = \{m \mid 0 \leq m \leq 2^n - 1 \text{ and } D_m \not\sqsubseteq_{\mathcal{T}} \perp\}$ ,

*i.e.*, the set of all integers  $m \in [0, 2^{n-1}]$  such that  $D_m$  be satisfiable in  $\mathcal{T}$ . Then, the utility of  $A$  can be expressed as

$$\min \left\{ \sum_{i=1}^n b_i * v_i \mid m \in \mathcal{S} \text{ and } (m)_2 = b_1 b_2 \dots b_n \right\} \quad (1)$$

Intuitively, one searches the minimum of the objective function (1) over a *subset*  $\mathcal{S}$  of the hypercube  $\{0, 1\}^n$ , where the vertices in  $\mathcal{S}$  are only the ones which define a satisfiable combination of  $A$  and (possibly negated) preferences in  $\mathcal{P}$ . Clearly, for every satisfiable conjunction at least one model  $M$  exists, and when the utility computed by (1) is minimum,  $M$  is a minimal model.

Finally, observe that a “right” number  $m$  between 0 and  $2^{n-1}$ —*i.e.*, a number individuating a minimal model—can be guessed nondeterministically in polynomial time. Hence, a nondeterministic Turing machine deciding  $u^c(A) < v$  guesses a number  $m$ , checks the satisfiability of  $D_m$  (a problem in  $\mathcal{C}$ ), computes  $u = \text{sum}_{i=1}^n b_i * v_i$ , and halts with “yes” if  $u < v$ , otherwise “no”. Therefore, when  $\mathcal{C} = \text{NP}$ , the overall decision problem is in NP; when  $\text{NP} \subset \mathcal{C}$ , the satisfiability check in  $\mathcal{C}$  dominates the overall complexity.  $\square$

For languages such that  $\text{PSPACE} \subseteq \mathcal{C}$ , the above theorem yields also an upper bound on the complexity of *computing* the utility of a concept. For instance, for  $\mathcal{L} = \mathcal{ALC}$ , and simple Terminologies, satisfiability is a problem PSPACE-complete [14]. Then computing (1) is a problem in PSPACE, too.

For languages such that  $\mathcal{C} = \text{NP}$ , the above theorem yields an NPO upper bound on the complexity of *computing* the utility of a concept. We discuss this case in more detail, for  $\mathcal{L} = \text{Propositional Logic (PL)}$ . For this case, Theorem 1 implies that the decision problem is NP-hard; yet it does not tell us whether the computation of  $u^c(A)$  admits some approximation schema, or not. We can show also NPO-hardness of computing (1), through a (not difficult) reduction from MAX-2-SAT, which is the problem of finding a model maximizing the number of clauses in a CNF of 2-literals (a.k.a. Krom) clauses. For convenience, we use the dual problem of finding a model that minimizes the number of *unsatisfied* conjunctions in a DNF of 2-literals conjunctions  $D = D_1 \vee \dots \vee D_n$ . Then such a model minimizes also the utility of any (unused) literal  $C$  w.r.t. the set of preferences  $\mathcal{P}_D = \{\langle D_1, 1 \rangle, \dots, \langle D_n, 1 \rangle\}$ . Since computing such a model is NPO-hard, our claim follows.

Observe that the above theorem does not hold for classes below NP, which need separate discussions. The case  $\mathcal{C} = \text{PTIME}$  covers the most simple logics, such as  $\mathcal{L} = \text{Conjunctions of Literals}$ , or the DL  $\mathcal{L} = \mathcal{FL}^-$  [6]. In fact, satisfiability of a conjunction amounts to check the absence of a literal and its negation in the conjunction. Yet, observe that if  $\mathcal{P} \supseteq \{\langle A, v_1 \rangle, \langle \neg A, v_2 \rangle\}$ , then for every formula  $C$ ,  $u^c(C) \geq \min(v_1, v_2)$ . In general, deciding if  $u^c(C) \leq k$  is NP-complete, based on a simple reduction from 3-TAUT: given a DNF  $D = D_1 \vee \dots \vee D_n$ , where each  $D_i$  is a conjunction,  $D$  is *not* a tautology iff  $u^c(A) \leq k$  ( $A$  being any literal) w.r.t. the preferences  $\mathcal{P} = \{\langle D_1, 2k \rangle, \dots, \langle D_n, 2k \rangle\}$ . Less obviously, the same reduction holds for  $\mathcal{L} = \mathcal{FL}^-$ , using a role  $R_i$  for every propositional atom  $A_i$ , and letting the encoding  $\gamma$  be:  $\gamma(\wedge) = \sqcap$ ,  $\gamma(A_i) = \exists R_i$ , and  $\gamma(\neg A_i) = \forall R_i.B$  (for some concept name  $B$ ). Then the previous DNF  $D$  is not a tautology iff  $u^c(B) \leq k$  w.r.t.  $\mathcal{P} = \{\langle \gamma(D_1), 2k \rangle, \dots, \langle \gamma(D_n), 2k \rangle\}$ . The fact that deciding the util-

ity of an  $\mathcal{FL}^-$  concept w.r.t. a set of  $\mathcal{FL}^-$  preferences is NP-hard is remarkable, since satisfiability in  $\mathcal{FL}^-$  is trivial (every  $\mathcal{FL}^-$  concept is satisfiable).

## 5 Computation of Minimal Utility Value

In this section we show how the computation of the *minimal utility value* for a set of preferences  $\mathcal{P}$  w.r.t. a concept  $A$  can be turned out in solving an optimization problem.

Given the set  $\mathcal{P} = \{\phi_1, \phi_2, \phi_3\}$  of preferences and the concept  $A$  as in Example 3, we note that:

- (a)  $A$  is more specific than the simple preference specification  $A_2$ ;
- (b)  $A$  is more specific than a disjunction of preference specification (that, in the most general case, may appear even negated).

$$\begin{aligned} \text{(a)} \quad & A \sqsubseteq_{\mathcal{T}} A_2 \longrightarrow \mathcal{T}^{\mathcal{I}} \cap A^{\mathcal{I}} \subseteq \mathcal{T}^{\mathcal{I}} \cap (A_2)^{\mathcal{I}} \\ \text{(b)} \quad & A \sqsubseteq_{\mathcal{T}} A_1 \sqcup A_3 \longrightarrow \mathcal{T}^{\mathcal{I}} \cap A^{\mathcal{I}} \subseteq \mathcal{T}^{\mathcal{I}} \cap (A_1)^{\mathcal{I}} \cap (A_3)^{\mathcal{I}} \end{aligned}$$

On the other hand, due to constraints modeled within the ontology we may have some interrelations among elements of  $\mathcal{P}$ . For instance, it might result that:

- (c) two preferences  $\phi_1$  and  $\phi_2$  cannot be satisfied at the same time;
- (d) the conjunction of the former with the complement of the latter could be unsatisfiable;
- (e) the combination of the complement of  $\phi_1$  and  $\phi_2$  is more specific than (*i.e.*, it implies) a third preference  $\phi_3$ .

In other words, (c) no model of  $A_1$  can be also a model of  $A_2$ , (d) no model of  $A_1$  can be also a model of  $\neg A_2$ , (e) all models of both  $\neg A_1$  and  $A_2$  are also models of  $A_3$ .

$$\begin{aligned} \text{(c)} \quad & A_1 \sqcap A_2 \sqsubseteq_{\mathcal{T}} \perp \longrightarrow \mathcal{T}^{\mathcal{I}} \cap (A_1)^{\mathcal{I}} \cap (A_2)^{\mathcal{I}} = \emptyset \\ \text{(d)} \quad & A_1 \sqcap \neg A_2 \sqsubseteq_{\mathcal{T}} \perp \longrightarrow \mathcal{T}^{\mathcal{I}} \cap (A_1)^{\mathcal{I}} \cap (\Delta \setminus (A_2)^{\mathcal{I}}) = \emptyset \\ \text{(e)} \quad & \neg A_1 \sqcap A_2 \sqsubseteq_{\mathcal{T}} A_3 \longrightarrow \mathcal{T}^{\mathcal{I}} \cap (\Delta \setminus (A_1)^{\mathcal{I}}) \cap (A_2)^{\mathcal{I}} \subseteq \mathcal{T}^{\mathcal{I}} \cap (A_3)^{\mathcal{I}} \end{aligned}$$

Actually, if we consider also a concept  $A$  it is easy to see that whenever (c), (d) or (e) hold, then also the corresponding relations represented below are true (while the vice versa is not true).

$$\begin{aligned} \text{(f)} \quad & A \sqcap A_1 \sqcap A_2 \sqsubseteq_{\mathcal{T}} \perp \longrightarrow \mathcal{T}^{\mathcal{I}} \cap A^{\mathcal{I}} \cap (A_1)^{\mathcal{I}} \cap (A_2)^{\mathcal{I}} = \emptyset \\ \text{(g)} \quad & A \sqcap A_1 \sqcap \neg A_2 \sqsubseteq_{\mathcal{T}} \perp \longrightarrow \mathcal{T}^{\mathcal{I}} \cap A^{\mathcal{I}} \cap (A_1)^{\mathcal{I}} \cap (\Delta \setminus (A_2)^{\mathcal{I}}) = \emptyset \\ \text{(h)} \quad & A \sqcap \neg A_1 \sqcap A_2 \sqsubseteq_{\mathcal{T}} A_3 \longrightarrow \mathcal{T}^{\mathcal{I}} \cap A^{\mathcal{I}} \cap (\Delta \setminus (A_1)^{\mathcal{I}}) \cap (A_2)^{\mathcal{I}} \subseteq \mathcal{T}^{\mathcal{I}} \cap (A_3)^{\mathcal{I}} \end{aligned}$$

In fact, in order to compute a *Minimal Utility Value*, if  $A$  represents *e.g.*, a final agreement, we are more interested in those models satisfying the latter equations rather than the ones satisfying (c), (d) and (e) because they are also models of  $A$  (as the *Minimal Model* is). Obviously, (a),(b),(f),(g),(h) can be generalized to whatever *Preference Set*.

Noteworthy is that, since we use an ontology  $\mathcal{T}$ , the above observations apply to preference specifications represented as general concept expressions  $C$  and not only as concept names. We illustrate the above ideas with the help of an example.

*Example 4 (Desktop Computer Negotiation cont'd).* We again refer to the ontology  $\mathcal{T}$  depicted in Figure 1.

We recall that the buyer's Preference Set is  $\mathcal{P} = \{\langle P_1, v_1 \rangle, \langle P_2, v_2 \rangle, \langle P_3, v_3 \rangle\}$ , with:

$$\begin{aligned} P_1 &= \forall \text{hasCPU.Centrino} \sqcap \exists \text{hasCPU} \\ P_2 &= \exists \text{hasCPU.AMD} \\ P_3 &= \exists \text{netSupport.WiFi} \end{aligned}$$

while the seller proposal is:

$$A \equiv \text{DesktopComputer} \sqcap \neg \exists \text{hasCPU.Centrino} \sqcap (\exists \text{netSupport.WiFi} \sqcup \forall \text{hasCPU.Sempron})$$

Notice that, because of the axioms in ontology  $\mathcal{T}$ :

$$\begin{aligned} \text{Intel} &\sqsubseteq \neg \text{AMD} \\ \text{Centrino} &\sqsubseteq \text{Intel} \sqcap \exists \text{arch} \sqcap \forall \text{arch.32BitCPU} \end{aligned}$$

preferences  $P_1$  and  $P_2$  cannot be satisfied at the same time, and, moreover, due to the axiom

$$\exists \text{hasCPU.Centrino} \sqsubseteq \exists \text{netSupport.WiFi}$$

in the ontology  $\mathcal{T}$ , preference  $P_1$  is more specific than preference  $P_3$ .  $\triangle$

**Definition 3 (Preference Clause).** Given a set of preferences  $\mathcal{P} = \{\langle P_i, v_i \rangle\}$ ,  $i = 1 \dots n$ , an ontology  $\mathcal{T}$  and a concept  $A$  such that  $A \not\sqsubseteq_{\mathcal{T}} \perp$ , we say that  $\mathcal{P}$  is constrained if the following condition holds:

$$A \sqcap \hat{P}_1 \sqcap \dots \sqcap \hat{P}_n \sqsubseteq_{\mathcal{T}} \perp \quad (2)$$

Where  $\hat{P}_i \in \{P_i, \neg P_i\}$ . We call (2) a **Preference Clause** if there is no strict subset  $\mathcal{Q} \subset \mathcal{P}$  such that  $\mathcal{Q}$  is constrained.

Note that with a *Preference Clause* one can represent not only relations (f) and (g) but also relations (a),(b) and (h) thanks to the well known equivalence:

$$C \sqsubseteq_{\mathcal{T}} D \iff C \sqcap \neg D \sqsubseteq_{\mathcal{T}} \perp$$

We may say that a *Preference Clause* contains the minimal set of preferences such that Equation (2) holds.

**Definition 4 (Preference Closure).** Given a Preference set  $\mathcal{P} = \{\phi_i\}$ ,  $i = 1 \dots n$ , an ontology  $\mathcal{T}$  and a concept  $A \not\sqsubseteq_{\mathcal{T}} \perp$ , we call **Preference Closure**, denoted as  $\mathcal{CL}$ , the set of Preference Clauses built, if any, for each set in  $2^{\mathcal{P}}$ .

In other words, a *Preference Closure* represents the set of all possible *Preference Clauses* over  $\mathcal{P}$ . It represents all possible (minimal) interrelations occurring between  $A$  and preference descriptions in  $\mathcal{P}$  w.r.t. an ontology  $\mathcal{T}$ .

**Proposition 1.** Given a concept  $A$ , a **Preference Closure**  $\mathcal{CL}$  and an ontology  $\mathcal{T}$ , if  $m$  is a Minimal Model of  $A$  then

$$m \models \mathcal{CL} \quad (3)$$

*Proof.* We denote with  $\overline{\mathcal{P}}^m$  the set of preferences satisfying point 2 of Definition 2 and with  $(\overline{\mathcal{P}}^m)^\mathcal{I}$  the interpretation  $\bigcap_{\langle P, v \rangle \in \overline{\mathcal{P}}^m} P^\mathcal{I}$ . By construction of  $\mathcal{CL}$  we have  $\mathcal{T}^\mathcal{I} \cap A^\mathcal{I} \subseteq \mathcal{CL}^\mathcal{I}$  while by Definition 2 it results  $m \in \{\mathcal{T}^\mathcal{I}\} \cap \{A^\mathcal{I}\} \cap \{(\overline{\mathcal{P}}^m)^\mathcal{I}\}$ . Hence,  $m \in \mathcal{CL}^\mathcal{I}$ .  $\square$

In order to compute *Minimal Utility Value*  $u^c(A)$  we reduce to an optimization problem (OP). Usually, in an OP we have a set of constrained numerical variables and a function to be maximized/minimized. In our case we will represent constraints as a set  $\chi$  of linear inequalities over binary variables, i.e., variables whose value is in  $\{0, 1\}$ , and the function to be minimized as a weighted combination of such variables. In order to represent  $\chi$  we need some pre-processing steps.

1. Compute the *Preference Closure*  $\mathcal{CL}$  for  $\mathcal{P}$ ;
2. For each *Preference Clause*  $A \sqcap \hat{P}_1 \sqcap \dots \hat{P}_n \sqsubseteq_{\mathcal{T}} \perp \in \mathcal{CL}$ , compute a corresponding preference constraint set  $\overline{\mathcal{CL}} = \{\neg \hat{P}_1, \dots, \neg \hat{P}_n\}$ . We denote with  $\overline{\mathcal{CL}}^c = \{\overline{\mathcal{CL}}\}$  the set of all preference constraint sets.

**Observation 1** *The reason why we do not consider  $\neg A$  when computing  $\overline{\mathcal{CL}}$  is that in order to compute a Minimal Utility Value we are looking for Minimal Models, i.e., models of  $A$  (and  $\mathcal{T}$ ) satisfying properties of Definition 2. Each Preference Clause can be rewritten as  $\mathcal{T} \models \neg A \sqcup \neg \hat{P}_1 \sqcup \dots \neg \hat{P}_n$ . If we rewrite the right hand side of the relation in terms of interpretation functions, from the semantics of  $\sqcup$  operator, we have*

$$(\neg A)^\mathcal{I} \cup (\neg \hat{P}_1)^\mathcal{I} \cup \dots (\neg \hat{P}_n)^\mathcal{I} \quad (4)$$

Since a *Minimal Model* is a model in  $A^\mathcal{I}$ , then all the models we are looking for are such that  $(\neg A)^\mathcal{I} = \emptyset$ . As a consequence, for our computation, the term  $(\neg A)^\mathcal{I}$  in Equation (4) is meaningless. We will clarify further this point in Observation 2 while discussing the OP we build to compute the *Minimal Utility Value* in the following.

*Example 5 (Desktop Computer Negotiation cont'd).* Consider again the Desktop Computer negotiation of Example 2. Given the set of preference  $\mathcal{P} = \{\langle P_1, v_1 \rangle, \langle P_2, v_2 \rangle, \langle P_3, v_3 \rangle\}$  and the proposal  $A$ , if we compute the **Preference Closure**  $\mathcal{CL}$  we find:

$$\mathcal{CL} = \left\{ \begin{array}{l} A \sqcap P_1 \sqsubseteq_{\mathcal{T}} \perp; \\ A \sqcap \neg P_1 \sqcap \neg P_2 \sqcap \neg P_3 \sqsubseteq_{\mathcal{T}} \perp \end{array} \right\}$$

hence, the two corresponding preference constraint sets in  $\overline{\mathcal{CL}}^c$  are:

$$\begin{aligned} \overline{\mathcal{CL}}_1 &= \{\neg P_1\} \\ \overline{\mathcal{CL}}_2 &= \{P_1, P_2, P_3\} \end{aligned}$$

△

Based on well-known encoding of clauses into linear inequalities (e.g., [15, p.314]) we transform each set  $\overline{\mathcal{CL}} \in \overline{\mathcal{CL}}^c$  in a set of linear inequalities  $\chi$  and then define a function to be minimized in order to solve an OP.

**Definition 5 (Minimal Utility Value OP).** Let  $\mathcal{P}$  be a set of preferences and  $\overline{\mathcal{CL}}^c$  be the set of all preference constraint sets. We define a Minimal Utility Value OP, represented as  $\langle \chi, u(\mathbf{p}) \rangle$ , the optimization problem built as follows:

1. **define numerical variables** – for each preference  $\langle P_i, v_i \rangle \in \mathcal{P}$ , with  $i = 1, \dots, n$  introduce a binary variable  $p_i$  and define the corresponding array  $\mathbf{p} = (p_1, \dots, p_n)$  (see Example 6);
2. **define  $\chi$  as a set of linear inequalities** – pick up each set  $\overline{\mathcal{CL}} \in \overline{\mathcal{CL}}^c$  and build the linear inequalities

$$\sum \{(1 - p) \mid \neg P \in \overline{\mathcal{CL}}\} + \sum \{p \mid P \in \overline{\mathcal{CL}}\} \geq 1$$

3. **define the function to be minimized** – given the array  $\mathbf{p}$  of binary variables

$$u(\mathbf{p}) = \sum \{v \cdot p \mid p \text{ is the variable mapping } \langle P, v \rangle\}$$

**Observation 2** If we considered also  $\neg A$  when computing the sets  $\overline{\mathcal{CL}} \in \overline{\mathcal{CL}}^c$  we would have had inequalities in the form:

$$(1 - a) + \sum \{(1 - p) \mid \neg P \in \overline{\mathcal{CL}}\} + \sum \{p \mid P \in \overline{\mathcal{CL}}\} \geq 1$$

Since we are interested in models where  $A$  is interpreted as true, then variable  $a$  has to be equal to 1. Hence the first element of the above summation is always equal to 0. In other words, we can omit  $\neg A$  when computing a preference constraint set  $\overline{\mathcal{CL}}$ .

The solution to a *Minimal Utility Value OP* will be an assignment  $\mathbf{p}_s$  for  $\mathbf{p}$ , i.e., an array of  $\{0, 1\}$ -values, minimizing  $u(\mathbf{p})$ .

*Example 6 (Desktop Computer Negotiation cont'd).* Back to the Desktop Computer negotiation of Example 2, after the computation of Preference Closures and set  $\overline{\mathcal{CL}}^c$ , we build the corresponding optimization problem in order to find the model with the minimal utility value:

$$\begin{aligned} \mathbf{p} &= (p_1, p_2, p_3) \\ \chi &= \begin{cases} 1 - p_1 & \geq 1 \\ p_1 + p_2 + p_3 & \geq 1 \end{cases} \\ u(\mathbf{p}) &= v_1 \cdot p_1 + v_2 \cdot p_2 + v_3 \cdot p_3 \end{aligned}$$

Possible solutions are:

$$\begin{aligned} \mathbf{p}'_s &= (0, 1, 0), u(\mathbf{p}'_s) = v_2 \\ \mathbf{p}''_s &= (0, 0, 1), u(\mathbf{p}''_s) = v_3 \\ \mathbf{p}'''_s &= (0, 1, 1), u(\mathbf{p}'''_s) = v_2 + v_3 \end{aligned}$$

The minimal solution will be either  $\mathbf{p}'_s$  or  $\mathbf{p}''_s$ , depending of the value of  $v_2$  and  $v_3$ . △

Given a solution  $\mathbf{p}_s$  to a *Minimal Utility Value OP*  $\langle \chi, u(\mathbf{p}) \rangle$ , we call *Minimal Preference Set*  $\overline{\mathcal{P}}^m$  and *Minimal Assignment*  $A^m$ , respectively, the set and the formula built as in the following<sup>6</sup>:

$$\begin{aligned}\overline{\mathcal{P}}^m &= \{\langle P_i, v_i \rangle \mid p_i = 1 \text{ in the solution } \mathbf{p}_s\} \\ A^m &= \prod \{P_i \mid p_i = 1 \text{ in the solution } \mathbf{p}_s\} \cap \prod \{\neg P_i \mid p_i = 0 \text{ in the solution } \mathbf{p}_s\}\end{aligned}$$

**Theorem 3.** *Given a solution  $\mathbf{p}_s$  to a Minimal Utility Value OP  $\langle \chi, u(\mathcal{P}) \rangle$  and a Minimal Assignment  $A^m$ :*

1. *if  $m$  is a Minimal Model then  $m \models A^m$ ;*
2.  *$u(\mathbf{p}_s)$  is a Minimal Utility Value.*

*Proof.* First we show that there exists at least one model  $\underline{m} \models_{\mathcal{T}} A^m$ . If  $\underline{m}$  does not exist, then  $\mathcal{T}^{\mathcal{I}} \cap A^{\mathcal{I}} \cap (A^m)^{\mathcal{I}} = \emptyset$ . We can easily rewrite the latter relation as  $A \cap A^m \sqsubseteq_{\mathcal{T}} \perp$  or, in other words,  $A$  and  $A^m$  are inconsistent with each other w.r.t.  $\mathcal{T}$ . But this is not possible. In fact, if  $A$  and  $A^m$  were inconsistent with each other w.r.t.  $\mathcal{T}$  then, by Proposition 1 we should have the corresponding Preference Clause in  $\mathcal{CL}$  and the related inequality in  $\chi$ :

$$\sum \{(1-p) \mid P \text{ appears in } A^m\} + \sum \{p \mid \neg P \text{ appears in } A^m\} \geq 1$$

*In order to be satisfied, the latter inequality must have either (a) at least one variable assigned to 0 in the first summation or (b) at least one variable assigned to 1 in the second one. Case (a) means that the corresponding preference is not satisfied by  $A^m$  while case (b) means that the corresponding preference is satisfied by  $A^m$ . Both cases are conflicting with the definition of  $A^m$ . Hence  $(A^m)^{\mathcal{I}} \cap \mathcal{T}^{\mathcal{I}} \neq \emptyset$ .*

*By construction of  $\chi$ , we have that all models of  $A^m$  w.r.t.  $\mathcal{T}$  are also models of  $A$  w.r.t. to  $\mathcal{T}$ . In formulas  $A^m \sqsubseteq_{\mathcal{T}} A$ . Since  $A^m$  comes from the minimization of  $u(\mathbf{p}_s)$  then  $(A^m)^{\mathcal{I}} \cap \mathcal{T}^{\mathcal{I}}$  represents all those models  $m \in A^{\mathcal{I}} \cap \mathcal{T}^{\mathcal{I}}$  such that*

$$\sum \{v \mid \langle P, v \rangle \in \mathcal{P} \text{ and } m \models P\}$$

*is minimal.*

*It is straightforward to show that  $u(\mathbf{p}_s)$  is a Minimal Utility Value. □*

## 5.1 Computational properties of the method

We now relate the computation method proposed in this section with the computational complexity results of the previous section. First of all, we observe that the size of the Preference Closure  $|\mathcal{CL}|$  can be—in the worst case—exponential in  $n$ , the size of the preference set. Since Linear Integer Programming is an NPO-complete problem[15], overall our problem can be solved nondeterministically in exponential time and space.

<sup>6</sup> with  $\prod\{\cdot\}$  we denote the conjunction of all the concepts in the set  $\{\cdot\}$

However, we observe that  $|\mathcal{CL}|$  *does not depend* on the size of the ontology  $\mathcal{T}$ , which typically is much larger than the size of  $\mathcal{P}$ . In some sense,  $\mathcal{CL}$  *compiles out* all the complexity due to the satisfiability problem in the chosen DL, leaving the OP of the combinatorics related to compatibility of preferences among each other, and with the formula  $C$  whose minimal utility value has to be computed. This is perfectly reasonable when Satisfiability in the chosen DL is a PSPACE-complete problem, or harder, since the best known procedures for solving PSPACE-complete problems use exponential time anyway, and the space used is exponential only in the number of preferences, not in the size of  $\mathcal{T}$ .

For the cases in which the language for preferences has a low-complexity satisfiability problem, say, NP, or PTIME, though, preprocessing into  $\mathcal{CL}$  the complete structure of preference compatibilities may be an overshoot. In such cases, it would seem more reasonable to devise specialized procedures that compute on demand the satisfiability of a conjunction of preferences.

An orthogonal analysis can be done on the *scalability* of the method when the utilities of several offers  $C_1, \dots, C_m$  must be compared. Here it seems that one has to solve  $m$  separate OPs of size exponential in  $|\mathcal{P}|$ . While this is the worst case, some optimization based on the logic for offers is possible. In fact, observe that  $C_i \sqsubseteq_{\mathcal{T}} C_j$  implies  $u^c(C_j) \leq u^c(C_i)$  (a model of  $C_i$  is also a model of  $C_j$ ). Hence, when searching for the offer with the maximum least utility,  $C_j$  can be safely disregarded. Intuitively, more specific offers are preferred over more generic ones, with the intuition that a generic offer  $C_j$  has a worst-case utility  $u^c(C_j)$  which is less than the worst-case utility  $u^c(C_i)$  of a more specific offer  $C_i$ .

## 6 Conclusion

Logic languages have been proposed here as a natural and powerful preference representation tool for automated negotiation purposes. We have shown how it is possible to compute a utility value for a concept (agreement, proposal, allocation), when preferences are expressed as weighted DL formulas w.r.t. a shared ontology  $\mathcal{T}$ . Although we ground our framework in the DLs realm, we point out that the framework itself is completely general and suitable for whatever decidable fragment of FOL. We also reported complexity results and showed the applicability and benefits of our approach with the help of a meaningful example. Currently, we are studying how to combine this approach with graphical models, and in particular GAI (Generalized Additive Independence) [16, 3], in order to model multiattribute auctions.

## References

1. Keeney, R.L., Raiffa, H.: Decisions with Multiple Objectives: Preferences and Value Trade-offs. John Wiley & Sons, New York (1976)
2. Gonzales, C., Perny, P.: GAI networks for utility elicitation. In: Ninth Intl. Conf. on Principles of Knowledge Representation and Reasoning, Whistler, BC, Canada (2004) 224–234
3. Bacchus, F., Grove, A.: Graphical models for preference and utility. In: Eleventh Conf. on Uncertainty in Artificial Intelligence, Montreal (1995) 3–10



4. Boutilier, C., Brafman, R.I., Domshlak, C., Hoos, H.H., Poole, D.: CP-nets: A tool for representing and reasoning about conditional ceteris paribus preference statements. *Journal of Artificial Intelligence Research* **21** (2004) 135–191
5. Engel, Y., Wellman, M.P.: CUI networks: A graphical representation for conditional utility independence. *Journal of Artificial Intelligence Research* **31** (2008) 83–112
6. Baader, F., Calvanese, D., McGuinness, D., Nardi, D., Patel-Schneider, P., eds.: *The Description Logic Handbook*. Cambridge Univ. Press (2002)
7. Pinkas, G.: Propositional non-monotonic reasoning and inconsistency in symmetric neural networks. In: *Twelfth Intl. Joint Conf. on Artificial Intelligence*. (1991) 525–531
8. Lafage, C., Lang, J.: Logical representation of preferences for group decision making. In: *Seventh Intl. Conf. on Principles of Knowledge Representation and Reasoning*. (2000) 457–468
9. Chevaleyre, Y., Endriss, U., Lang, J.: Expressive power of weighted propositional formulas for cardinal preference modelling. In: *Tenth Intl. Conf. on Principles of Knowledge Representation and Reasoning*. (2006) 145–152
10. Ragone, A., Di Noia, T., Di Sciascio, E., Donini, F.: A logic-based framework to compute Pareto agreements in one-shot bilateral negotiation. In: *Seventeenth European Conference on Artificial Intelligence*. (2006) 230–234
11. Ragone, A., Di Noia, T., Di Sciascio, E., Donini, F.M.: Description logics for multi-issue bilateral negotiation with incomplete information. In: *Twenty-Second AAAI Conference on Artificial Intelligence*. (2007) 477–482
12. Ragone, A., Di Noia, T., Di Sciascio, E., Donini, F.M.: Alternating-offers protocol for multi-issue bilateral negotiation in semantic-enabled marketplaces. In: *6th International Semantic Web Conference (ISWC'07)*. Volume 4825 of *Lecture Notes in Computer Science*., Springer-Verlag (2007) 395–408
13. Lukasiewicz, T., Schellhase, J.: Variable-strength conditional preferences for matchmaking in description logics. In: *Tenth Intl. Conf. on Principles of Knowledge Representation and Reasoning*. (2006) 164–174
14. Baader, F., Hollunder, B.: *KRIS: Knowledge Representation and Inference System*. *SIGART Bulletin* **2**(3) (1991) 8–14
15. Papadimitriou, C.H., Steiglitz, K.: *Combinatorial Optimization: Algorithms and Complexity*. Prentice-Hall (1982)
16. Fishburn, P.C.: Interdependence and additivity in multivariate, unidimensional expected utility theory. *International Economic Review* **8** (1967) 335–342

# Explaining and Predicting the Behavior of BDI-Based Agents in Role-Playing Games<sup>\*</sup>

M.P. Sindlar, M.M. Dastani, F. Dignum, and J.-J.Ch. Meyer  
{michal,mehdi,dignum,jj}@cs.uu.nl

University of Utrecht  
P.O. Box 80.089, 3508 TB Utrecht, The Netherlands

**Abstract.** Virtual characters in games operate in a social context involving other characters and human players. If such socially situated virtual characters are to be considered believable, they should be able to adjust their behavior based on beliefs about the mental states of other characters and human players. Autonomous BDI-based agents are suitable for modeling characters that exhibit such intentional behavior. In this paper, it is illustrated how agent-based characters can infer the mental state of other players by observing those other players' actions in the context of a declarative game specification. The game specification can be utilized in explanation and prediction of agents' behavior, and as such can form the basis for developing socially aware characters.

## 1 Introduction

For games and simulations with interactive virtual characters to provide users with a satisfying experience, it is of vital importance that those characters are believable to the user. Appearing to pursue goals and to be responsive to social context are determining factors for believability [1], and interaction with virtual characters is richer and more enjoyable if these anticipate the behavior of other characters [2]. Believable characters that operate in a social context should exhibit social awareness and not only pursue their own interests, but also be able to take the mental states of other characters into account if they believe these to conflict or coincide with their own goals and beliefs.

Statistical approaches to game-based plan recognition exist [3] but require large amounts of gameplay data to be processed, which might not always be available. Recent work in the agent programming community has focused on recognizing an agent's plan on grounds of its program and observed actions, and inferring the mental state of the agent that plausibly explains observed behavior [4,5]. This offers promising directions for developing socially aware virtual characters, as characters which can infer other characters' mental states have the possibility to incorporate attributed mental states into their own decision-making

---

<sup>\*</sup> This research has been supported by the GATE project, funded by the Netherlands Organization for Scientific Research (NWO) and the Netherlands ICT Research and Innovation Authority (ICT Regie).

process. This allows for plausibly misguided behavior and thus contributes to believability [6]. However, current approaches ignore the setting in which agents operate and inferred explanations are independent of social context. If characters in games are designed to behave as autonomous agents, then a game can be regarded as an agent society [7]. The behavior of characters in such a game takes place in the context of the society, which can be captured and specified by means of an organizational model [8].

This paper gives a declarative solution to abducing the mental state of virtual characters implemented as BDI-based agents, which takes into account the context in which these characters operate. In Sect. 2 mental state abduction is reviewed. Sect. 3 introduces the game-related context in terms of scenes and roles, and in Sect. 4 it is shown how this context can be utilized in explaining and predicting agent behavior. Sect. 5 ties things together in an example, and Sect. 6 concludes with a brief discussion and ideas for future research.

## 2 Mental State Abduction

In this section it is described how the observed behavior of agents can be related to an explanation in terms of a description of their mental state, recapitulating our work in [5]. In Defs. 1–2 the behavior of agents is described, and in the remainder of this section it is shown how behavior is observed and explained.

Let  $\mathcal{L}$  be a propositional domain language with negation and conjunction over propositions.  $\text{Lit} \in \mathcal{L}$  is the set of literals in this language, and  $\mathcal{L}_\Gamma \subseteq \mathcal{L}$  a simple language allowing only (conjoined) literals. Let  $\text{Act}$  be a set of action names. The behavior of an agent can be described as an expression consisting of actions, which are either atomic observable actions  $\alpha \in \text{Act}$ , or tests  $\phi?$  on propositions  $\phi \in \mathcal{L}$ . Actions can be composed by means of operators for sequential composition (;) and choice (+).

**Definition 1 (behavioral description)** *Let  $\alpha \in \text{Act}$  be an atomic observable action, and  $\phi?$  the test action on proposition  $\phi \in \mathcal{L}$ . The set of behavioral descriptions  $\mathcal{L}_\Pi$  with typical element  $\pi$  is then defined as follows.*

$$\pi ::= \alpha \mid \phi? \mid \pi_1; \pi_2 \mid \pi_1 + \pi_2$$

Note that there is no notion of iteration in  $\mathcal{L}_\Pi$ . It is assumed, though, that behavior can be iteratively performed if an agent reapplies a behavioral rule of the type defined in Def. 2. Such a rule states that the behavior described by  $\pi$  is appropriate for achieving the goal  $\gamma$  if the condition  $\beta$  is believed to hold, and is taken to be interpreted in the operational context of an agent program [9,10].

**Definition 2 (behavioral rules)** *Behavioral rules specify a relation between behavior  $\pi \in \mathcal{L}_\Pi$  and a mental state consisting of achievement goals  $\gamma \in \mathcal{L}_\Gamma$  and beliefs  $\beta \in \mathcal{L}$ . The set of behavioral rules  $\mathcal{L}_{\mathcal{BR}}$  has  $\text{br}$  as its typical element.*

$$\text{br} ::= \gamma \leftarrow \beta \uparrow \pi$$

The atomic actions  $\alpha \in \text{Act}$  of an agent are taken to be publicly observable, such that a perceived action can be directly related to the performed action. Because behavioral descriptions do not allow for concurrent actions and interpretation of the agent program is assumed to do so neither, the perception of multiple actions performed by a single agent is taken to be sequential. To distinguish sequential perception of actions — which is an incremental process — from actions in sequence as part of a plan, percepts are presented as a list.

**Definition 3 (percepts)** *Let  $\alpha \in \text{Act}$  be an observable action and  $\epsilon$  a special empty (null) action, such that  $(\epsilon\delta) = (\delta) = (\delta\epsilon)$ . The set of percept expressions  $\mathcal{L}_\Delta$ , with typical element  $\delta$ , is then defined as follows.*

$$\delta ::= \alpha \mid \epsilon \mid \delta_1\delta_2$$

In order to explain perceived actions in intentional terms, a relation has to be established between perceived actions, descriptions of behavior, and the behavioral rules that connect a mental state (goals and beliefs) to behavior. This relation should be defeasible, because on grounds of observed actions alone it is not necessarily possible to analytically infer the mental state that caused the agent to perform those actions. For this reason the behavioral rules of Def. 2 are described as the logical implications defined in Def. 4, which state that a precondition consisting of a goal and belief description implies a certain behavior.

**Definition 4 (rule description)** *Let  $(\gamma \leftarrow \beta \uparrow \pi) \in \mathcal{L}_{\mathcal{BR}}$  be a behavioral rule. The set of rule descriptions  $\mathcal{L}_{\mathcal{RD}}$ , with typical element  $\text{rd}$ , is defined as follows.*

$$\text{rd} ::= \text{goal}(\gamma) \wedge \text{belief}(\beta) \Rightarrow \text{behavior}(\pi)$$

*The function  $\text{desc} : \mathcal{L}_{\mathcal{BR}} \longrightarrow \mathcal{L}_{\mathcal{RD}}$  maps rules to their description, such that  $\text{desc}(\text{br}) = (\text{goal}(\gamma) \wedge \text{belief}(\beta) \Rightarrow \text{behavior}(\pi))$  for any  $\text{br} = (\gamma \leftarrow \beta \uparrow \pi) \in \mathcal{L}_{\mathcal{BR}}$ .*

Because behavioral rules, as defined in Def. 2, are interpreted in an operational context, the implications in the rule descriptions defined in Def. 4 do not hold in a classical logical way with respect to describing agent operation. However, if perceived actions  $\delta \in \mathcal{L}_\Delta$  can be related to the behavioral descriptions  $\pi \in \mathcal{L}_\Pi$ , then logical abduction — which states that on grounds of an observation  $\psi$  and a logical implication  $\phi \Rightarrow \psi$ , the defeasible explanation  $\phi$  can be abduced — can be used to infer the preconditions of logical rule descriptions. To relate action sequences to descriptions of behavior, a function is defined that maps behavioral descriptions to sets of *observable traces* of behavior by filtering out internal tests, which are taken to be unobservable, and branching observable traces at points where choice occurs. The operator  $\cup$  is standard set union, and  $\circ : \wp(\mathcal{L}_\Delta) \times \wp(\mathcal{L}_\Delta) \longrightarrow \wp(\mathcal{L}_\Delta)$  is a non-commutative composition operator defined as  $\Delta_1 \circ \Delta_2 = \{ \delta_1\delta_2 \mid \delta_1 \in \Delta_1 \text{ and } \delta_2 \in \Delta_2 \}$ , where  $\Delta_1, \Delta_2 \subseteq \mathcal{L}_\Delta$ .

**Definition 5 (observable trace function)** *Let  $\alpha \in \text{Act}$ ,  $\phi \in \mathcal{L}$  and  $\pi \in \mathcal{L}_\Pi$ . The function  $\tau : \mathcal{L}_\Pi \longrightarrow \wp(\mathcal{L}_\Delta)$  is then defined as follows.*

$$\begin{aligned} \tau(\alpha) &= \{\alpha\} & \tau(\pi_1 + \pi_2) &= \tau(\pi_1) \cup \tau(\pi_2) \\ \tau(\phi?) &= \{\epsilon\} & \tau(\pi_1; \pi_2) &= \tau(\pi_1) \circ \tau(\pi_2) \end{aligned}$$

In order to abduce mental state preconditions of the logical rule descriptions in Def. 4 on grounds of observed actions, a relation must be established between an observed action sequence and the traces that represent the observable aspect of the behavior described in the behavioral description part of a rule. If every action of the agent is observed and the rules completely describe an agent's possible behavior, then an observed sequence of actions can be related to the (non-strict) prefix of an observable trace of behavior described by some  $\pi \in \mathcal{L}_\Pi$ .

**Definition 6 (structural relations)** *Let  $\preceq \subseteq \mathcal{L}_\Delta \times \mathcal{L}_\Delta$  be the prefix relation on sequences  $\delta, \delta' \in \mathcal{L}_\Delta$  and  $\succcurlyeq \subseteq \mathcal{L}_\Delta \times \mathcal{L}_\Delta$  the suffix relation, defined as follows.*

$$\delta \preceq \delta' \text{ iff } \exists \delta'' \in \mathcal{L}_\Delta : (\delta' = \delta\delta'') \quad \delta \succcurlyeq \delta' \text{ iff } \exists \delta'' \in \mathcal{L}_\Delta : (\delta' = \delta''\delta)$$

*Note that every  $\delta$  is a prefix and suffix of itself iff  $\delta''$  is the empty action  $\epsilon$ .*

Defining different structural relations, as shown in [5], may allow for relating observed actions to observable traces also in the case that not every action is observed. It was proven that this leads to an increase in abduced explanations, and that the set of explanations inferred on grounds of complete observation is a subset of explanations in case of partial observation. In order to focus on the way contextual information can be used to facilitate the process of mental state abduction, it is assumed here that observation is complete.

An agent's behavior is to be explained in terms of a description of its mental state. In order to refer to these preconditions, let the set  $\mathcal{L}_\Omega$  with typical element  $\omega$  be defined as  $\mathcal{L}_\Omega = \{ \text{goal}(\gamma) \wedge \text{belief}(\beta) \mid \gamma \in \mathcal{L}_\Gamma, \beta \in \mathcal{L} \}$ . An explicit 'lifting' notation for functions is used, such that for any  $f : D \rightarrow D'$ , the *lifted* version of  $f$  is  ${}^\wp f : \wp(D) \rightarrow \wp(D')$ , such that for  $\Phi \subseteq D$ ,  ${}^\wp f(\Phi) = \{ f(\phi) \mid \phi \in \Phi \}$ .

An explanatory function is now defined that maps observed action sequences  $\delta \in \mathcal{L}_\Delta$  to preconditions of rule descriptions of the type  $\text{rd} \in \mathcal{L}_{\mathcal{RD}}$ , such that  $\delta$  is a (partial) trace of the behavior 'implied' by the rule description.

**Definition 7 (explanatory function)** *The function  $\chi : \mathcal{L}_{\mathcal{RD}} \rightarrow \wp(\mathcal{L}_\Omega \times \mathcal{L}_\Delta)$  maps a rule description to a set of tuples of precondition and trace.*

$$\chi(\omega \Rightarrow \text{behavior}(\pi)) = \{ (\omega, \delta) \mid \delta \in \tau(\pi) \}$$

*Let  $\delta \in \mathcal{L}_\Delta$  be a percept and  $\mathcal{RD} \subseteq \mathcal{L}_{\mathcal{RD}}$  a set of rule descriptions. The explanatory function  $\text{explain} : \mathcal{L}_\Delta \times \wp(\mathcal{L}_{\mathcal{RD}}) \rightarrow \wp(\mathcal{L}_\Omega)$  is then defined as follows.*

$$\text{explain}(\delta, \mathcal{RD}) = \{ \omega \mid \exists (\omega, \delta') \in {}^\wp \chi(\mathcal{RD}) : [\delta \preceq \delta'] \}$$

Somewhat less formally, the explanatory function defined in Def. 7 states that the precondition of a rule description is in the set of explanations for a certain observed sequence of actions, if the observed sequence is a (non-strict) prefix of any trace of the behavioral description which is described in the postcondition of the rule description. The function as defined here is not intended to be computationally efficient. It can be proven, though, that  $\text{explain}(\delta\delta', \mathcal{RD}) \subseteq \text{explain}(\delta, \mathcal{RD})$  for any  $\delta, \delta' \in \mathcal{L}_\Delta$ , allowing for an efficient implementation.<sup>1</sup>

<sup>1</sup> The authors express their thanks to Henry Prakken for pointing out this stronger and more concise version of their proofs in [5].

### 3 Agents Playing Games

Agents in a multi-agent system each have their mental state and (inter)act in pursuit of their private goals, taking into account their beliefs and the means provided by their environment(s) [9]. An *agent-based game*, must be more than a regular multi-agent system, as the latter lacks particular qualities that a game might be required to have. When implementing game characters as autonomous agents, a designer gives away part of the behavioral control that a scripting-based approach to character design provides [11]. In return, the daring move of the designer is rewarded with emergent stories that take unexpected turns because of decisions made by the autonomous characters. However, there are certain aspects of the game’s ‘flow of events’ that the game designer wants to ensure, without having to rely on providence or agents’ good insight.

In this section declarative game-related concepts are defined, inspired by organizational principles, that are used to illustrate how an agent-based game can be designed that respects some storyline marked out by the designer. Moreover, the same concepts can be used as a guideline with respect to agents’ expected behavior, as will be shown in Sect. 4.

#### 3.1 A Declarative Game Specification

The general concept ‘game’ is hard to define, so that in the present approach a particular kind of game is considered, namely the *agent-based role-playing game*. Such a game is considered to be populated by virtual characters implemented as autonomous BDI-based agents, which play roles similar to the way actors do in a movie. Autonomous agents, however, may be allowed more freedom in the way they enact their role than movie actors are. Such role-enacting agents can be allowed to have private goals that supercede or conflict with those specified by their role, which will show in the behavior they exhibit [12,7].

Most definitions of the concept ‘role’ recognize that a role comes with obligations, permissions, authority, or the right to perform certain actions. As such, a role describes behavior which can be expected of the role-enacting agent [13]. Roles are therefore defined in a way that encompasses both the descriptive and prescriptive aspect by providing an agent with goals to achieve, information made available to the role-enacting agent, and behavioral rules. These concepts correspond to the B(eliefs), D(esires), and I(intentions) of the BDI-paradigm, and as such can form the basis for design of the role-enacting agent. Note that the relation of role-derived goals and goal-directed rules is taken to be not necessarily one-to-one; multiple rules for a single goal can exist, or a conjunctive goal may be provided by the role where only rules for literal goals exist. To be able to refer to roles and other entities uniquely, a set of constants ID is introduced.

**Definition 8 (role)** Let  $\Gamma_{\leq} = (\Gamma, \leq_{\Gamma})$  be an ordered set of goals  $\Gamma \subseteq \mathcal{L}_{\Gamma}$ , with  $\leq_{\Gamma} \subseteq \mathcal{L}_{\Gamma} \times \mathcal{L}_{\Gamma}$  a partial order on  $\Gamma$ . Let  $\mathcal{I} \subseteq \mathcal{L}$  be role-accessible information, and  $\mathcal{BR} \subseteq \mathcal{L}_{\mathcal{BR}}$  a set of behavioral rules. A role R, identified by a unique identifier  $r \in \text{ID}$ , is then defined as  $R = \langle r, \Gamma_{\leq}, \mathcal{I}, \mathcal{BR} \rangle$ .

A typical role, featured in many games of the role-playing game (RPG) genre, is that of the *thief*. Unsurprisingly, the thief-role may provide the goal to take possession of a particular item by stealing it. Moreover, a thief could have the goal to steal the item whilst double-checking that nobody is near. If the thief assesses a particular situation to be risky, the goal to steal the item but also ensure that nobody is around might supercede the goal to just steal the item.

Roles do not per definition remain unchanged throughout a game. The context in which a role is enacted influences the way it should be enacted, and this context may change as things happen in the game. Autonomous agents can be given the liberty to enact their role as they see fit, resulting in different types of behavior given the same role specification. Nevertheless, agents are restricted in their actions by the opportunities provided by their environment, and by the norms of the agent society in which they operate. To formalize the norms that regulate behavior, a language of normative expressions  $\mathcal{L}_{\mathcal{N}}$  is defined which captures *prima facie norms*, with typical element  $N$ , such that  $N ::= F(\alpha) \mid O(\alpha)$ . The expression  $F(\alpha)$  states that the action  $\alpha \in \text{Act}$  is forbidden,  $O(\alpha)$  states that the action is obligatory.

In [14], *prima facie norms* are defined to be *norms [which] usually do not arise from actions, but arise in certain situations [and remain] valid as long as the situation in which they arise stays valid*. Scenes are taken to constitute the norm-governed context in which roles remain unchanged. The scene definition includes the roles figuring in the scene, a set of norms pertaining to the scene, and a set of literals denoting an initial environment state.

**Definition 9 (scene)** *Let  $\mathcal{R}$  be a set of roles as defined in Def. 8,  $\mathcal{N} \subseteq \mathcal{L}_{\mathcal{N}}$  a set of norms, and  $E \subseteq \text{Lit}$  the initial environment state. Scene  $S$ , with unique identifier  $s$ , is then defined as  $S = \langle s, \mathcal{R}, \mathcal{N}, E \rangle$ .*

Take a scene in an RPG that features the thief and a store owner in some store where goods can be purchased. In this scene it is most likely forbidden to take items without paying for them, or to damage the merchandise.<sup>2</sup> Now take a scene in which the thief and the store owner are joined in the store by a city guard. The same norms may apply as in the previous scene, but the thief now gives priority to ensuring nobody is around before stealing anything because of the presence of the guard, whereas the store owner might be more at ease in knowing that the eyes of the law are keeping watch over her belongings.

### 3.2 The Multi-Agent Game

A game is taken to be a composition of scenes. The way scenes are composed (the ‘storyboard’ of the game) is defined in a game specification, which identifies the scenes occurring in a game, and specifies when a specific scene makes a transition to another scene. Such transitions might depend on conditions being

<sup>2</sup> Note that the *prima facie norms* of  $\mathcal{L}_{\mathcal{N}}$  do not allow for conditional statements, and it is therefore not possible to express statements such as the fact that it is obligatory to pay for an item after taking it.

fulfilled with respect to the environment of the scene, on specific actions being (jointly) executed by agents, or even some condition becoming true with respect to agents' mental states. To have a system of agents obey this specification, scene transition has to be operationalized in the semantics of the agent system. Because a detailed presentation of how the scene transition is realized does not contribute to the scope of the present approach, this is left unspecified and the game is taken to simply be a set of scenes, of the type defined in Def. 9.

Agents in 2APL [9] are defined by a configuration, which specifies their mental state in terms of goals, belief, plans, and rules. In this paper, we do not commit ourselves to an assumption about the specific language in which the agents are implemented, but do require that the behavior of agents is in accordance with the declarative specification of the role they enact, which contains elements that can be directly related to elements of agents' mental states, such as goals, information (beliefs) and behavioral rules. Specifically, the following assumptions and restrictions are enforced.

- *Every agent in the multi-agent game behaves in accordance with a role, such that the behavior of agents is completely described by the behavioral description which is part of the rules accompanying their role.*
- *The role of agents prescribes specific partially ordered goals, and the rules accompanying the role are taken to enable achievement of all these goals. However, it is not necessarily the case that every goal which the agent may have on grounds of its rules is part of the goals that the agent's role prescribes.*
- *It is assumed that agents do not interleave plans, even if they have multiple goals. If an agent has adopted multiple goals and has selected a plan based on the application of a behavioral rule for one of its goals, it will not apply a new rule until its selected plan is completed.*

The first scene of the multi-agent game is determined by the game's initial state, and consecutive scenes are determined as the game evolves; ie. as agents act in pursuit of their goals and 'things happen' in the game. Because it is not in the interest of the topic at hand, which is explanation and prediction of agents' behavior in the context of a multi-agent role-playing game, the operational transition of configurations of the multi-agent game will not be presented formally. Instead, it is assumed that the game takes place in some (known) scene, which provides a guideline with respect to behavior that can be expected of the agents populating the scene, as the behavior of each of them is based on some role.

## 4 Explaining and Predicting Agent Behavior

Mental state abduction can be used to abduce the mental state of BDI-based agents whose behavior can be observed. If this behavior is performed in the context of a multi-agent game, then information about the scene of the game and the role which agents enact helps improve the abduction process. If an agent's role is known, the set of rules the agent is taken to have at its disposition is reduced to the set of rules provided by the role, as the behavior descriptions in the rules of the roles completely describes behavior of the agents.



#### 4.1 Explaining Agent Behavior

In the approach to mental state abduction as described in Sect. 2 (and in [5] in more detail), the behavior of an agent is explained on grounds of all rules this agent can be assumed to have if context is not considered. In the present setting, only the rules which are ascribed to the agent on account of its role in a particular scene are considered in the explanatory process. This ensures that the explanations provided for its behavior are *contextually grounded*, and that the set of rules which need to be considered is restricted in size.

The role of the agent contains behavioral rules and a partially ordered set of goals. There might exist agents which dutifully pursue the goals their role prescribes, and others which don't care about their role in the least. To capture these aspects of *role conformance*, two refined versions of the explanatory function are defined. Because the role-prescribed goals do not necessarily have a one-to-one correspondence with the goals that form the head of behavioral rules, a relation between the two has to be established. The functions  $g$  and  $\bar{b}$  are defined, such that for  $\omega = (\text{goal}(\gamma) \wedge \text{belief}(\beta))$ , it holds that  $g(\omega) = \{\gamma\}$  and  $\bar{b}(\omega) = \{\beta\}$ .  $\text{Cn}(\Phi)$  denotes the closure of the set  $\Phi$  under the consequence operator  $\text{Cn}$ , defined as  $\text{Cn}(\Phi) = \{\phi \mid \Phi \models \phi\}$ .

**Definition 10 (loosely role-conformant explanation)** *Let  $\delta \in \mathcal{L}_\Delta$  be a percept and  $\langle r, (\Gamma, \leq_\Gamma), \mathcal{I}, \mathcal{BR} \rangle$  a role, as defined in Def. 8. The function  $\text{explain}_{lrc}$  for loosely role-conformant explanation is then defined as follows.*

$$\begin{aligned} \text{explain}_{lrc}(\delta, \langle r, (\Gamma, \leq_\Gamma), \mathcal{I}, \mathcal{BR} \rangle) &= (\Omega, \leq_\Omega) \\ \text{where } \Omega &= \text{explain}(\delta, {}^p\text{desc}(\mathcal{BR})), \text{ and for any } \omega, \omega' \in \Omega \\ \leq_\Omega &= \{ (\omega, \omega') \mid [\text{Cn}(g(\omega)) \not\subseteq \text{Cn}(\Gamma)] \wedge [\text{Cn}(g(\omega')) \subseteq \text{Cn}(\Gamma)] \} \\ &\cup \{ (\omega, \omega') \mid [\text{Cn}(g(\omega)) \cap \text{Cn}(\Gamma) = \emptyset] \wedge [\text{Cn}(g(\omega')) \not\subseteq \text{Cn}(\Gamma)] \} \\ &\cup \{ (\omega, \omega) \} \end{aligned}$$

A rule can be said to be relevant to a role, if the goal for which this rule applies is in the closure of the role-derived goals. Thus, the rules for goals  $\phi$  and  $\psi$  are both relevant to a role that prescribes the goal  $\phi \wedge \psi$ , just as the rule for  $\phi \wedge \psi$  is relevant to a role that prescribes  $\phi$  and  $\psi$  independently. The function  $\text{explain}_{lrc}$  maps to a poset of explanations, where the explanations are ordered on grounds of an ordering that ranks explanations containing role-derived goals over those with goals that derive from behavioral rules only. Explanations which contain exclusively role-derived goals rank over those with some role-derived goals, which in turn rank over explanations without role-derived goals.

**Definition 11 (strictly role-conformant explanation)** *The definition of the function  $\text{explain}_{src}$  is based on  $\text{explain}_{lrc}$ , but takes into account the order on  $\Gamma$ .*

$$\begin{aligned} \text{explain}_{src}(\delta, \langle r, (\Gamma, \leq_\Gamma), \mathcal{I}, \mathcal{BR} \rangle) &= (\Omega, \leq_\Omega) \\ \text{where } \text{explain}_{lrc}(\delta, \langle r, (\Gamma, \leq_\Gamma), \mathcal{I}, \mathcal{BR} \rangle) &= (\Omega, \leq'_\Omega), \text{ and for any } \omega, \omega' \in \Omega \\ \leq_\Omega &= \{ (\omega, \omega') \mid \exists \gamma \in \text{Cn}(g(\omega)), \exists \gamma' \in \text{Cn}(g(\omega')) : [\gamma <_\Gamma \gamma'] \wedge \\ &\quad \neg \exists \gamma \in \text{Cn}(g(\omega)), \exists \gamma' \in \text{Cn}(g(\omega')) : [\gamma' <_\Gamma \gamma] \} \cup \leq'_\Omega \end{aligned}$$

Strictly role-conformant agents are taken to also obey the priority ordering on goals specified by their role, and therefore  $\text{explain}_{src}$  takes this ordering into account this as well. Because not all goals need to be explicitly ordered, it is defined that some explanation  $\omega$  is preferred to  $\omega'$  on grounds of  $\text{explain}_{src}$  if and only if some goal  $\gamma$ , derived from  $\omega$ , has explicit priority over some goal  $\gamma'$ , derived from  $\omega'$ , and no goal derived from  $\omega'$  has explicit priority over any goal derived from  $\omega$ .

Instead of conforming to their role, agents might *rebel* against their role. Also, as explained in [12], agents which are allowed to have private objectives along with role-derived objectives can enact their roles in a selfish or social manner. This could imply an ordering which is the reverse of that seen in loose role conformance, or even of strict role conformance. Although it is not further dealt with, the fact that our approach allows for modeling explicit rebellion and different types of role enactment deserves pointing out.

## 4.2 Predicting Agent Behavior

An observed and explained sequence of actions can be regarded as the performed part of a trace. Given that the goal for which this plan was selected is still active, the agent can be expected to perform the remaining actions, which are the suffix of the trace of which the observed actions are the prefix. In explaining an agent's behavior, it was defined that a description of the agent's mental state can be regarded as an explanation. When predicting the behavior of the agent with respect to actions it has been observed to perform, multiple (distinct) action sequences may be predicted based on different assumed mental states. A predictive function is defined, taking these aspects into account.

**Definition 12 (predictive function)** *Let  $\delta, \delta' \in \mathcal{L}_\Delta$  be percepts,  $\omega \in \mathcal{L}_\Omega$  a mental state description and  $\mathcal{RD} \in \mathcal{LRD}$  a set of rule descriptions. The function  $\text{predict} : \mathcal{L}_\Delta \times \mathcal{LRD} \rightarrow \wp(\mathcal{L}_\Omega \times \mathcal{L}_\Delta)$  is then defined as follows.*

$$\text{predict}(\delta, \mathcal{RD}) = \{ (\omega, \delta') \mid (\omega, \delta\delta') \in \wp\chi(\mathcal{RD}) \}$$

Agents in a norm-governed society can be assumed to take norms into account in choosing their actions, either by design or by deliberation [15]. Similar to the explanatory functions taking into account role conformance of the agent (Defs. 10 & 11), one can consider *norm obedience* when predicting agent behavior. The norms of  $\mathcal{LN}$  were defined to state about actions whether these are either forbidden (*F*) or obligatory (*O*). Informally,  $F(\alpha)$  is taken to mean that the action  $\alpha \in \text{Act}$  is forbidden and that agents may be punished if they perform the action, whereas  $O(\alpha)$  states that agents are obliged to perform action  $\alpha$  and that may be punished if they do not perform it. Note that it is not defined what it means that the agent “*may* be punished”, but the explanation that the behavior of the agent is somehow monitored (possibly by law-enforcing agents in the game), and that this monitoring is not infallible, should suffice.

Thus, it may occur that the agent performs a forbidden action, but gets away with it. The predicates `forb` and `obl` are defined on  $\delta \in \mathcal{L}_\Delta$ , such that

$$\begin{aligned} \mathcal{N} \models \text{forb}(\delta) & \text{ iff } \exists \alpha, \delta', \delta'' \in \mathcal{L}_\Delta : [(\delta = \delta' \alpha \delta'') \wedge (F(\alpha) \in \mathcal{N})] \\ \mathcal{N} \models \text{obl}(\delta) & \text{ iff } \exists \alpha, \delta', \delta'' \in \mathcal{L}_\Delta : [(\delta = \delta' \alpha \delta'') \wedge (O(\alpha) \in \mathcal{N})] \end{aligned}$$

Based on the above, a predictive function is defined which takes norm obedience into account. This function predicts a sequence of actions on grounds of an observed sequence of actions and behavioral rules, and relates it to the presumed mental state which would account for observed behavior if it were the agent's actual mental state. Moreover, this predictive function takes into account that norms may exist which forbid or oblige the agent to perform specific actions, as expressed in the ordering that ranks pairs with an action sequence containing some obliged but no forbidden actions above all others, and pairs with sequences that contain some forbidden but no obliged actions below all others.<sup>3</sup>

**Definition 13 (norm-obedient prediction)** *Let  $\delta \in \mathcal{L}_\Delta$  be a percept, the tuple  $\langle r, \Gamma_{\leq}, \mathcal{I}, \mathcal{BR} \rangle$  a role as defined in Def. 8 and  $\langle s, \mathcal{R}, \mathcal{N}, \mathbf{E} \rangle$  a scene as defined in Def. 9. The predictive function  $\text{predict}_{no}$  is then defined as follows.*

$$\begin{aligned} \text{predict}_{no}(\delta, \langle r, \Gamma_{\leq}, \mathcal{I}, \mathcal{BR} \rangle, \langle s, \mathcal{R}, \mathcal{N}, \mathbf{E} \rangle) &= (\Theta, \leq_\Theta) \\ \text{where } \Theta &= \text{predict}(\delta, {}^{\circ}\text{desc}(\mathcal{BR})), \text{ and for any } (\omega, \delta), (\omega', \delta') \in \Theta \\ \leq_\Theta &= \{ ((\omega, \delta), (\omega', \delta')) \mid \mathcal{N} \models [\text{obl}(\delta') \wedge \neg \text{forb}(\delta')] \} \\ &\cup \{ ((\omega, \delta), (\omega', \delta')) \mid \mathcal{N} \models [\text{forb}(\delta) \wedge \neg \text{obl}(\delta)] \} \\ &\cup \{ ((\omega, \delta), (\omega, \delta)) \} \end{aligned}$$

In Sect. 4.1 the remark was made that agents can explicitly rebel against their role. Similarly, agents might rebel against ‘society’, which can be modeled by means of explicitly presumed norm disobedience, such that traces with forbidden actions are considered to be preferred by the agent.

### 4.3 The Observer

To explain and predict behavior, an abstract external Observer is proposed (in line with our approach in [5]) which perceives the atomic observable actions performed by agents, attempting to explain those actions in context of the game and making predictions about actions it expects agents to perform next. The Observer maintains a model of each of the agents it observes, which contains the role the Observer attributes to the agent and a sequence of actions the agent has been observed to perform, along with explanations and predictions based on observed behavior in context of the attributed role.

**Definition 14 (agent model)** *Let  $R$  be a role of the type in Def. 8,  $\delta \in \mathcal{L}_\Delta$  a list of perceived actions,  $\Omega \subseteq \mathcal{L}_\Omega$  a set of explanations and  $\Theta \subseteq \mathcal{L}_\Theta$ . An agent model, with a unique identifier  $i \in \text{ID}$ , is then defined as  $A = \langle i, R, \delta, \Omega, \Theta \rangle$ .*

<sup>3</sup> Note that a sequence with forbidden as well as obligatory actions is treated no differently than one that has only ‘neutral’ actions.

The Observer is assumed to have perfect observation of the environment and the actions agents perform. In many games, the roles of characters are evident from their external characteristics. The role might be indicated by the color of a suit, or simply by a label hovering over the character. In the following, it is assumed that agent  $i$ 's role  $r$  can be deduced from the state of the environment, such that  $E \models \text{enacts}(i, r)$ . As scene transitions are taken to depend only on changes in the environment, the Observer always knows the scene in which the game takes place if it is made aware of the initial scene when the game starts, and is always correct about roles it attributes to agents.

Relaxing these assumptions — either by introducing more uncertainty on part of the Observer by design or because the game does not allow for perfect observation of the environment, scene transitions, or agents' roles — leads to interesting scenarios. Instead of just performing mental state abduction, the Observer is forced to perform role abduction and/or scene abduction as well. If observation of the environment or agents' actions is imperfect as well, yet more defeasibility is introduced. Given that our goal is to allow for designing agent-based game characters which have uncertainty about other characters' mental states, it is not in our interest to introduce any more uncertainty than necessary. Partial observation was discussed in [5], but here perfect observation is assumed.

**Definition 15 (Observer)** *Let  $\mathcal{G}$  be a set of scenes as defined in Def. 9, and  $s \in \text{ID}$  a scene identifier such that  $\langle s, \dots \rangle \in \mathcal{G}$ . The set of literals  $E \subseteq \text{Lit}$  is the environment state and for every (perceived) agent  $i$ ,  $A_i$  is an agent model. The Observer is then defined as  $\langle \mathcal{G}, s, E, \{A_i, \dots, A_j\} \rangle$ .*

The Observer as defined in this approach is an abstract entity, which serves to illustrate the explanatory and predictive process ultimately to be used by agents that observe other agents' behavior in some environment. For this reason the details of how the Observer configuration evolves with successive action observations and scene transitions are left to the imagination of the reader, and instead the focus is on the procedures defined in Sect. 4.1 and 4.2. Given a single sequence of observed actions for some agent  $i$ , the Observer can explain as well as predict this sequence of actions. Prop. 1 shows that each explanation — in terms of an agent's mental state — is accompanied by a matching prediction.

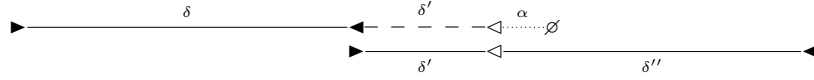
**Proposition 1 (explanation matches prediction).** *Given an agent model  $\langle i, \langle r, \Gamma_{\leq}, \mathcal{I}, \mathcal{BR} \rangle, \delta, \Omega, \Theta \rangle$ , where  ${}^{\circ}\text{desc}(\mathcal{BR}) = \mathcal{RD}$ ,  $\text{explain}(\delta, \mathcal{RD}) = \Omega$ , and  $\text{predict}(\delta, \mathcal{RD}) = \Theta$ , it holds that  $\forall \omega \in \Omega : [\exists \theta \in \Theta, \delta' \in \mathcal{L}_{\Delta} : [\theta = (\omega, \delta')]]$ .*

*Proof.* Def. 7 and Def. 12 show that  $\text{explain}$  and  $\text{predict}$  are both based on  ${}^{\circ}\chi$ . In case of  $\text{explain}(\delta, \mathcal{RD}) = \Omega$ , some  $\omega \in \Omega$  iff  $\exists (\omega, \delta'') \in {}^{\circ}\chi(\mathcal{RD}) : [\delta \preceq \delta'']$ . For  $\text{predict}(\delta, \mathcal{RD}) = \Theta$ , some  $(\omega, \delta') \in \Theta$  iff  $(\omega, \delta\delta') \in {}^{\circ}\chi(\mathcal{RD})$ . It follows from Def. 6 that if  $\delta \preceq \delta''$ , then  $\delta'' = \delta\delta'$  for some  $\delta'$  (possibly  $\delta' = \epsilon$ ), and therefore  $(\omega, \delta') \in \Theta$ . By definition of  $\forall$ , the proposition holds for  $\Omega = \emptyset$ .  $\square$

Note that Prop. 1 extends to  $\text{explain}_{trc}$ ,  $\text{explain}_{src}$  and  $\text{predict}_{no}$ , as these are directly based on  $\text{explain}$  and  $\text{predict}$ . This is a very welcome fact, because it ensures that for every explanation a corresponding prediction can be made, also in

the case of the context-dependent explanatory and predictive functions. Based on role-conformant explanation some explanation may come out as ‘top-ranked’. This can be considered the best explanation for the agent’s behavior, and corresponding predicted behavior be regarded as the most probable. Prediction of  $\epsilon$  — possibly indicating goal achievement — is outside of the scope of this approach, but very well worth further investigation.

It can occur that traces have overlapping segments. In such a case, the possibility exists that the Observer is able to explain a sequence of actions, but at a certain point observes an action that can neither be considered coherent with the currently presumed trace, nor as being the start of a new trace. This situation is visualized in Fig. 1, where  $\delta'$  is the suffix of some trace  $\delta\delta'$ , as well as the prefix of another trace  $\delta'\delta''$ . Let  $\alpha$  be the first action of  $\delta''$ , and  $\delta''' = \delta\delta'\alpha$ . If the Observer explains  $\delta\delta'$  as a coherent whole, then after perceiving  $\alpha$ , it may be that  $\text{explain}(\delta''', \mathcal{RD}) = \emptyset$  because  $\delta'''$  is not the prefix of any trace. It may, however, also be that  $\alpha$  itself is not the prefix of any trace, such that  $\text{explain}(\alpha, \mathcal{RD}) = \emptyset$ .



**Fig. 1.** Two traces,  $\delta\delta'$  and  $\delta'\delta''$ , with an overlapping part  $\delta'$ . The start of actual traces is denoted by  $\blacktriangleright$  and the end by  $\blacktriangleleft$ , explainable (segments of) potential traces end with  $\triangleleft$ , and  $\emptyset$  denotes a non-matching segment (ie. failure of explanation).

If the situation sketched in the previous paragraph occurs and explanation fails, then the Observer can backtrack along  $\delta'$ , starting at the end, until it finds a suffix  $\delta''' \succ \delta'$  that can be explained in coherence with the last observed action  $\alpha$ , such that  $\text{explain}(\delta'''\alpha, \mathcal{RD}) \neq \emptyset$ . Given the assumption that agents are assumed to be able to completely execute their plans, the maximum overlap of any two traces of an agent’s plans can be computed and used to give a measure of the maximum amount of backtracking the Observer has to perform. Let  $\text{len} : \mathcal{L}_\Delta \rightarrow \mathbb{N}$  be a function that maps a percept to its length, such that  $\text{len}(\alpha_1, \dots, \alpha_n) = n$ , and let  $\text{binds}$  be predicate denoting that some sequence  $\delta$  ‘binds’ the sequences  $\delta'$  and  $\delta''$  together with overlapping action sequences, defined as

$$\text{binds}(\delta, \delta', \delta'') \quad \text{iff} \quad [(\delta \succ \delta' \wedge \delta \preceq \delta'') \vee (\delta \succ \delta'' \wedge \delta \preceq \delta')]$$

Given the definitions of  $\text{len}$  and  $\text{binds}$ , let  $\text{overlap} : \mathcal{L}_\Delta \times \mathcal{L}_\Delta \rightarrow \mathbb{N}$  be a function that computes the overlap between two (distinct) action sequences, defined as

$$\text{overlap}(\delta', \delta'') = \begin{cases} \text{len}(\delta) & \text{if } \exists \delta \in \mathcal{L}_\Delta : \text{binds}(\delta, \delta', \delta''), \text{ and } \delta' \neq \delta'', \text{ and} \\ & \neg \exists \delta''' \in \mathcal{L}_\Delta : [\text{binds}(\delta''', \delta', \delta'') \wedge (\text{len}(\delta''') > \text{len}(\delta))] \\ 0 & \text{otherwise} \end{cases}$$

**Proposition 2 (backtrack with maximum trace overlap).** *Given an agent model  $\langle i, \langle r, \Gamma_{\leq}, \mathcal{I}, \mathcal{BR} \rangle, \delta, \Omega, \Theta \rangle$ , where  ${}^{\circ}\text{desc}(\mathcal{BR}) = \mathcal{RD}$  and  $\delta = \delta' \alpha_1 \cdots \alpha_n$ , for which it is the case that  $\text{explain}(\delta, \mathcal{RD}) = \emptyset$ ,  $\text{explain}(\delta' \alpha_1 \cdots \alpha_{n-1}, \mathcal{RD}) \neq \emptyset$ , and  $\delta'$  is the complete trace of an actual plan executed by the agent, there exists a non-empty suffix  $\delta'' \succ \delta$  such that  $\text{explain}(\delta'', \mathcal{RD}) \neq \emptyset$  and  $\text{len}(\delta'')$  is smaller than or equal to one, plus the maximum overlap of any two traces of any plan which is part of the rules in  $\mathcal{BR}$ .*

*Proof.* Let  $\Delta = \bigcup \{ \tau(\pi) \mid (\gamma \leftarrow \beta \uparrow \pi) \in \mathcal{BR} \}$  be the set of all (finite) observable traces of all plans part of the rules in the agent model. Because  $\delta'$  is a complete plan trace,  $\text{explain}(\delta', \mathcal{RD}) \neq \emptyset$  and  $\delta' \in \Delta$ . Let  $\delta'' \succ \delta$  such that  $\delta''$  is the prefix of a trace of the agent's latest plan and it holds that  $\delta'' \succ \alpha_1 \cdots \alpha_n$ . Then either  $\exists \delta''' \in \mathcal{L}_{\Delta} : [(\delta''' \preccurlyeq \alpha_1 \cdots \alpha_n) \wedge (\delta' \delta''' \in \Delta)]$  such that the actual 'old' trace  $\delta'$  is also the prefix of a misleading 'false' trace  $\delta' \delta'''$ , or not. If not, then  $\delta = \delta' \alpha$  such that  $\alpha \succ \delta$ ,  $\text{explain}(\alpha) \neq \emptyset$ , and  $\text{len}(\alpha) = 1$ .

If  $\exists \delta''' \in \mathcal{L}_{\Delta} : [(\delta''' \preccurlyeq \alpha_1 \cdots \alpha_n) \wedge (\delta' \delta''' \in \Delta)]$ , then  $\delta'''$  is the suffix of a 'false' trace  $\delta' \delta'''$  and the strict prefix of  $\alpha_1 \cdots \alpha_n$ , such that  $\text{len}(\delta''') < n$ . The 'new' trace, of which  $\delta''$  is the prefix, is started somewhere after the plan of which  $\delta'$  is a complete trace has finished, such that  $\delta'' \succ \delta' \alpha_1 \cdots \alpha_n$ . If the 'new' trace of which  $\delta''$  is the prefix is not started directly after  $\delta'$ , because inbetween a complete trace of yet another plan was executed which together with  $\delta'$  could be matched to a misleading trace, then  $\delta''$  is a strict suffix of  $\alpha_1 \cdots \alpha_n$ , such that  $\text{len}(\delta'') > n$ . If  $\delta''$  is started directly after  $\delta'$ , then the sequence  $\delta' \delta'''$  and the sequence  $\delta'' = \alpha_1 \cdots \alpha_n$  have an overlap of  $n - 1$ , which is the overlap of  $\delta' \delta'''$  and the trace of which  $\delta''$  is the prefix, such that  $\text{len}(\delta'') = n$ .

Let  $x$  be the maximum overlap of any two traces  $\delta_1, \delta_2 \in \Delta$ , such that  $(\text{overlap}(\delta_1, \delta_2) = x) \wedge (\neg \exists \delta_3, \delta_4 \in \Delta : [(\text{overlap}(\delta_3, \delta_4) = y) \wedge (y > x)])$ . Given that the trace  $\delta' \delta'''$  and the 'new' trace of which  $\delta''$  is the prefix are both in  $\Delta$  and have an overlap of  $n - 1$ , it holds that  $0 \leq n - 1 \leq x$ .  $\square$

Prop. 2 can be guaranteed if the Observer does 'forget' any percepts and agents complete their plans. Especially the latter condition is unmaintainable in certain environments. If agents can drop their plans, 'freak' scenarios can arise. Take, for example, the case where  $\alpha_1 \cdots \alpha_n$  is a plan trace, but the individual actions  $\alpha_1, \dots, \alpha_n$  are also the initial actions of individual plans. If an agent selects those plans in order, executing only the first action and then dropping the plan, the resulting sequence is indistinguishable from the trace.<sup>4</sup> However, because traces are finite and actions perfectly observable, Coroll. 1 still applies.

**Corollary 1.** *If agents drop their plans, then the Observer backtracks at most up to the length of the longest trace to find an explanation if  $\text{explain}(\delta, \mathcal{RD}) = \emptyset$ .*

*Proof.* Let  $\alpha_1 \cdots \alpha_n \in \Delta$  be the longest trace. As  $\text{explain}(\delta, \mathcal{RD}) = \emptyset$ , it must be that  $\exists \delta' \in \mathcal{L}_{\Delta} : [\delta' \succ \delta]$  and  $\delta'$  is the prefix of the agent's current plan. In worst case,  $\delta = \delta'' \alpha_1 \cdots \alpha_n$  for some  $\delta''$ , such that  $\text{explain}(\delta'' \alpha_1 \cdots \alpha_{n-1}) \neq \emptyset$ . After backtracking  $\text{len}(\alpha_1 \cdots \alpha_n) = n$  actions, Observer finds  $\text{explain}(\alpha_1 \cdots \alpha_n) \neq \emptyset$ .  $\square$

<sup>4</sup> One might ask whether explaining and predicting behavior has any benefit at all if such situations abound in some scenario, but that is not the point now.

## 5 Example

To illustrate the present approach, an example inspired by the popular role-playing game Oblivion [16] is introduced. Because of space limitations, some shorthand notation will be used. Lowercase predicate arguments represent ground atoms, and uppercase arguments represent variables. Our propositional language of course does not allow for variables, and therefore these are to be interpreted as a finite number of ground expressions, as should be clear from the context in which the notation is used. Spatial environments require moving around, and therefore `goto(Loc)` is defined, where the variable *Loc* stands for any valid location in the environment, and  $(\phi?; \pi) + (\neg\phi?; \pi')$  means **if  $\phi$  then  $\pi$  else  $\pi'$** .

$$\text{goto}(Loc) \equiv (\neg\text{nearby}(Loc)?; \text{walk\_towards}(Loc)) + (\text{nearby}(Loc)?)$$

The scene *S* takes place in a store and features ‘thief’ and ‘store owner’ roles. The norm in this scene forbids stealing any item, as expressed in shorthand notation, such that  $S = \langle s, \{R_t, R_{so}\}, \{F(\text{steal}(Item))\}, E \rangle$ . The ‘thief’ role prescribes the goal to have a particular item of interest ( $\gamma = \text{have}(item)$ ), and provides rules to achieve this goal. Also, a rule for exploring the store is provided.

$$\begin{aligned} R_t &= \langle \text{thief}, (\{\text{have}(item)\}, \{(\gamma, \gamma)\}), \mathcal{I}, \{\text{br}_{1_t}, \text{br}_{2_t}, \text{br}_{3_t}, \text{br}_{4_t}\} \rangle \\ \text{br}_{1_t} &= \text{have}(item) \leftarrow \text{distracted}(owner) \wedge \text{in}(Cabinet, item) \uparrow \\ &\quad \text{goto}(Cabinet); \text{open}(Cabinet); \text{steal}(item); \text{close}(Cabinet) \\ \text{br}_{2_t} &= \text{have}(item) \leftarrow \neg\text{distracted}(owner) \uparrow \text{goto}(owner); \text{distract}(owner) \\ \text{br}_{3_t} &= \text{explored}(store) \leftarrow \neg\text{explored}(cabinet_1) \wedge \dots \wedge \neg\text{explored}(cabinet_n) \uparrow \\ &\quad \text{goto}(cabinet_1); \text{inspect}(cabinet_1); \dots; \text{goto}(cabinet_n); \text{inspect}(cabinet_n) \\ \text{br}_{4_t} &= \text{ensured}(safety) \leftarrow \neg\text{nearby}(Person) \uparrow \text{double} - \text{check\_if\_nearby}(Person) \end{aligned}$$

The ‘store owner’ role  $R_{so}$  is left unspecified, except that it is stated she wants to protect her merchandise. In this paper procedural rules have not been discussed, but they may be allowed if only used for goal generation on grounds of events. If a customer breaks an object in the store, the perception of this fact prompts the store owner to adopt the goal to demand money from the culprit. This high-level approach remedies shortcomings in scripted character behavior in a natural way; in the game Oblivion it is possible, for example, to jump on the store counter or to smash objects without repercussions, because the store owner is only scripted to react to theft, and apparently not to vandalism.

The scene *S* transitions to some new scene *S'* upon entry of a city guard, as mentioned in Sect. 3.1, such that  $S' = \langle s', \{R'_t, R'_{so}, R_{cg}\}, \{F(\text{steal}(Item))\}, E' \rangle$ . The ‘city guard’ role  $R_{cg}$  is left unspecified, but it should suffice to say that the guard has merely come into the store to buy some item or chat with the store owner. If he becomes aware that someone is breaking the law (possibly the thief stealing the item), he may come into action and arrest the perpetrator. In the new scene the thief shows more cautious behavior because of the presence of the guard. This is illustrated by a change in the thief’s role specification, such

that  $R'_t = \langle \text{thief}', (\{\gamma, \gamma'\}, \{(\gamma, \gamma), (\gamma', \gamma'), (\gamma, \gamma')\}), \mathcal{I}, \{\text{br}_{1_t}, \text{br}_{2_t}, \text{br}_{3_t}, \text{br}_{4_t}\} \rangle$ . In this slightly changed role specification, it still is the case that  $\gamma = \text{have}(\text{item})$ , but there is another role-prescribed goal  $\gamma' = \text{have}(\text{item}) \wedge \text{ensured}(\text{safety})$  for which it is the case that  $\gamma' >_I \gamma$ .

Various possibilities exist for improvement, but lack of space forces us to skim over subtleties. More interesting is it to see how the Observer comes into play. Let  $\mathcal{G}$  be the scenes of the game, and  $A_{gent} = \langle \text{gent}, R_t, \text{walk\_towards}(\text{cabinet}_1) \rangle$  the model of some agent called *gent*, such that the Observer has observed *gent*, in its role of *thief*, to perform the action of walking towards a certain cabinet. Let the Observer configuration for the first scene be  $\langle \mathcal{G}, s, \{A_{gent}\} \rangle$ . Given the rules  $\mathcal{BR}_t$  for the thief role,  $\text{explain}(\text{walk\_towards}(\text{cabinet}_1), \mathcal{BR}_t)$  maps to a set of explanations  $\Omega = \{\omega_1, \omega_2\}$ , such that  $g(\omega_1) = \{\text{goal}(\text{have}(\text{item}))\}$  and  $g(\omega_2) = \{\text{goal}(\text{explored}(\text{store}))\}$ . Based on role-conformant explanation, either loose or strict,  $\omega_1 > \omega_2$  because having the item is a role-derived goal.

Prop. 1 states that every explanation is matched by a prediction. For  $\omega_1$ , the tuples  $(\omega_1, [\text{open}(\text{Cabinet}), \text{steal}(\text{item}), \text{close}(\text{Cabinet})])$  (with the percept in Prolog-style list notation) are in the set of predictions. Given the small scenario and limited set of rules, this is the only possible prediction for the goal of having the item. If the thief is assumed not to be norm-obedient, the prediction that he will open the cabinet and steal the item comes out, which is plausible in this context. Assuming the thief actually is norm-obedient (which would be plausible in the scene  $S'$  where the guard is also present) gives a different picture. In that case another explanation for walking towards the cabinet can be considered best, if the corresponding predicted action sequence doesn't contain any forbidden actions. In this example only the goal to explore the store qualifies, but in a more extensive case this could include the goal to choose and purchase some item located in the cabinet.

The example in this section is inspired by an actual commercial role-playing game. It serves mainly to illustrate some of the focal points of the approach presented in this paper, and is necessarily limited in its scope and detail. Nevertheless, it should be sufficiently rich to convince the reader of the fact that the high-level concepts of organizational modeling and agent programming apply transparently to the complex world of role-playing games. Moreover, the use of high-level social/intentional concepts has the additional benefit that these concepts can be reused for modeling, programming, and inter-character explanation and prediction of behavior.

## 6 Conclusion and Future Work

In this paper mental state abduction in the context of an agent-based game was described. A declarative game specification based on organizational principles such as roles, norms, and scenes, was introduced, and it was mentioned how it can be employed to have a system of autonomous agents behave in accordance with an intended storyline. An abstract Observer was said to observe the behavior of agents and provide explanations that take into account role-conformant



behavior, making the abduction process more efficient because it is based on a subset of rules, and ensuring that explanations are relevant to context. The Observer can also predict agents' future actions based on previously observed behavior, taking norm-obedience into account if the situation warrants this assumption. Role-conformant explanation and norm-obedient prediction have been shown to be complementary.

Future research should focus on explicitly taking models of the environment and agents' presumed mental states into account in the abductive process. Depending on the nature of the environment, it could be possible for an observer, be it an abstract entity or situated agent, to actively check whether specific (predicted) actions are possible, or whether an agent has achieved its goal or has some particular belief. Also, making use of norms in the explanatory process as well as the predictive process is to be considered. Finally, the path of formally investigating the multi-agent game as an operational system seems promising.

## References

1. Loyall, A.B.: Believable Agents. PhD thesis, Carnegie Mellon University (1997)
2. Laird, J.E.: It knows what you're going to do: Adding anticipation to a Quakebot. In: AGENTS. (2001)
3. Albrecht, D.W., Zukerman, I., Nicholson, A.E.: Bayesian models for keyhole plan recognition in an adventure game. *User Modeling & User-Adapted Interaction* **8**(1-2) (1998) 5–47
4. Goultiaeva, A., Lespérance, Y.: Incremental plan recognition in an agent programming framework. *Proceedings of PAIR* (2007)
5. Sindlar, M.P., Dastani, M.M., Dignum, F., Meyer, J.-J.Ch.: Mental state abduction of BDI-based agents. In: *Proceedings of DALI*. (2008) 110–125
6. Scott, B.: Architecting a Game AI. In: *AI Game Programming Wisdom*. Charles River Media (2002) 285–289
7. Dignum, V.: A Model for Organizational Interaction. PhD thesis, SIKS Dissertation Series (2004)
8. Coutinho, L.R., Sichman, J.S., Boissier, O.: Modeling organization in MAS. In: SEAS. (2005)
9. Dastani, M.: 2APL: A practical agent programming language. *Autonomous Agents and Multi-Agent Systems* **16** (2008) 214–248
10. Pokahr, A., Braubach, L., Lamersdorf, W.: Jadex: A BDI reasoning engine. In Dastani, M., Dix, J., El Fallah Seghrouchni, A., eds.: *Multi-Agent Programming*. Springer (2005) 149–174
11. Tozour, P.: The Perils of AI Scripting. In: *AI Game Programming Wisdom*. Charles River Media (2002) 541–547
12. Dastani, M., Dignum, V., Dignum, F.: Role-assignment in open agent societies. In: *Proceedings of AAMAS*. (2003)
13. Dastani, M., Riemdsijk, M.B.V., Hulstijn, J., Meyer, J.-J.Ch.: Enacting and detecting roles in agent programming. In: *Proceedings of AOSE*. (2004) 189–204
14. Dignum, F.: Autonomous agents with norms. *Artificial Intelligence and Law* **7**(1) (1999) 69–79
15. Castelfranchi, C., Dignum, F., Jonker, C.M., Treur, J.: Deliberative normative agents: Principles and architecture. In: *ATAL*. (1999) 364–378
16. Bethesda Game Studios: *The Elder Scrolls IV: Oblivion* (2006)

# Correctness Properties for Multiagent Systems

Munindar P. Singh<sup>1</sup> and Amit K. Chopra<sup>2</sup>

<sup>1</sup> North Carolina State University, Raleigh, USA  
singh@ncsu.edu

<sup>2</sup> Università degli Studi di Trento, Trento, Italy  
akchopra.mail@gmail.com

**Abstract.** What distinguishes multiagent systems from other software systems is their emphasis on the interactions among autonomous, heterogeneous agents. This paper motivates and characterizes correctness properties for multiagent systems. These properties are centered on commitments, and capture correctness at a high level. In contrast to existing approaches, commitments underlie key correctness primitives formalized in terms of meaning; for example, commitment *alignment* maps to *interoperability*; commitment *discharge* maps to *compliance*. This paper gives illustrative examples and characterizations of these and other properties. The properties cover the specification of protocols, roles, and agents, the principal artifacts of an interaction-based approach to designing multiagent systems, and thus provide the formal underpinnings of any such approach.

## 1 Introduction

Interaction is the key distinguishing feature of multiagent systems. We investigate the science of interaction as it underlies the engineering of multiagent systems whose constituent agents are *heterogeneous* (independently designed) and *autonomous* (independently motivated). In such systems, the internal workings of the agents take backstage to the interactions among them.

We begin from a simple yet profound question: *How may we treat interactions as first-class citizens in modeling and analyzing multiagent systems?* The usual objectives of engineering—modularly specifying, developing, composing, verifying, and validating parts—apply for interactions just as for traditional software approaches. However, existing solutions, which are designed for components such as objects, do not readily lift to interactions: an interaction somehow must simultaneously accommodate more than one perspective. Thus, importantly, the novelty of the interactive setting yields fresh and crucial technical challenges, which offer a great opportunity for multiagent systems research.

Of the many applications of multiagent systems, those in *cross-organizational* business processes provide the happy mix of practical value, theoretical subtlety, and opportunity (in the form of interest in industry) that our research community needs to sustain this research effort. Cross-organizational processes fundamentally differ from conventional software in that they are naturally modeled via interactions among *heterogeneous* and *autonomous* agents [1]. The interactions of interest are of an arms-length nature, and thus naturally understood as *communications*. In our study, we assume the existence

of suitable approaches for the transmittal of information and therefore concentrate on communication understood at the level of meaning.

To engineer a multiagent system based on interactive principles presupposes a notion of the correctness of interactions among agents—in particular, here, of communications. Given such a notion, we ask if an agent is *compliant* with its expected behavior. Further, we can ask if the given agents are *interoperable* meaning that they are able to work together as expected. We can ask the above questions from the perspective of the system as a whole or of any of the participants. To formalize interoperability and compliance in interactive terms requires that we develop a theory of types using which we might modularize communications into *protocols*. We might then create repositories of protocols; determine if one protocol refines another or aggregates two or more protocols; modularly validate the protocols; modularly verify agents with respect to each relevant protocol; and so on. Notice that interfaces in object-oriented computing correspond to protocols and support constructs such as refinement and aggregation as well as the usual forms of type inference.

### 1.1 Approach

The meaning of an interaction thus lies at the crux of the question of its correctness. When we think at levels above the transmission of information, the meaning of communication is grounded in the relationships among the parties involved. Communication then is based on conventions by which such relationships are created, progressed (or otherwise altered), and ended. We concentrate on the contractual relationships expressed through the notion of commitments. A *commitment* involves a debtor, a creditor, an antecedent, and a consequent; it is represented as  $C(\text{debtor}, \text{creditor}, \text{antecedent}, \text{consequent})$ . Roughly, the debtor stakes a claim or makes a promise to the creditor about the specified consequent provided that the antecedent holds (Section 2 explains this further). Commitments naturally express the *what* of business relationships, and minimally constrain the *how*. For example, a commitment to pay for goods received may be discharged by paying directly, or delegated to someone who would discharge or delegate it, and so on (for any finite sequence of delegations).

In our approach, a *protocol* specifies business interactions primarily by stating how messages affect the participants' commitments. For example, returning purchased goods unopened may release the buyer from a commitment to pay. Thus many possible enactments may result from the same protocol. This is how commitments yield both rigor and flexibility. Because of its naturalness, the commitment-based approach has attracted the attention of finance and health care industry groups [2].

*Protocols* are interfaces: they constrain how agents interact, not how they are implemented. Protocols are *doubly* modular: in terms both of functionality and autonomy. For example, for functionality, an ORDER protocol between a customer and a merchant would specify only interactions dealing with order placement, leaving other functionalities to separate protocols, e.g., one for INVENTORY FULFILLMENT. Our approach enables composing protocols to yield more complex protocols, of enhanced functionality. Further, for autonomy, ORDER would specify the interactions, leaving to each the autonomous decision making of whether and how to interact, which could depend on

its goals [3]. We define a *process* as the aggregation of the behaviors of the parties involved, including both their interactions and their local reasoning.

To model a process, we identify the protocols using which the different participants interact [1]. For example, a merchant and a customer may interact with each other using a NEGOTIATION protocol; the merchant, customer, and payment agency may interact via an ESCROW protocol; and, the merchant, customer, and shipper may interact through some specialized LOGISTICS protocol. When each participant acts according to its local reasoning but respecting the stated protocols, they jointly enact a multiparty business process. The contractually significant parts of the process would have been encoded in the commitments specified in the protocols; the other parts may feature only in the local policies of the participants and need not be visible externally. An agent's policies could be geared to optimize its outcomes. For example, policies would help decide what item to order, what price to quote, and so on.

Notice that the above approach obviates defining a monolithic global flow that specifies the actions of each party. Each protocol could be refined to capture additional requirements, e.g., adding receipts or guarantees to SHIPPING or PAYMENT to produce new refined protocols. Notice that protocols can involve more than two parties; in typical usage, one partner would play multiple roles in multiple protocols [4]. For example, a purchase process may be defined as a composition of ORDER, SHIPPING, and PAYMENT protocols where the buyer in ORDER is the receiver in SHIPPING and the payer in PAYMENT.

*The potential benefits* of our protocol-based approach over traditional approaches (reviewed in Section 4) include the following. One, for process *design*, protocols are naturally reusable whereas complete processes are not. More importantly, protocols lend themselves to modeling abstractions such as refinement and aggregation. Two, for process *implementation*, implementations of agents playing multiple roles can be more readily assembled from specifications of the roles. Three, for process *enactment*, flexible protocols enable each agent to exercise discretion via its policies or preferences even as it follows a protocol. For example, a merchant may accept only cash for discounted goods and a customer may prefer to pay for goods early or late depending upon private considerations such as of fiscal year accounting. This flexibility also enables us to capture and handle business exceptions and opportunities in a natural manner at the level of protocols. Four, for process *monitoring*, protocols provide a clean basis for determining that the interacting agents are complying with their roles in the given protocols.

## 1.2 Contributions

We motivate and characterize the key properties that would enable engineering multiagent systems with a special emphasis on applications such as cross-organizational processes. Compared to traditional formal approaches, the emphases on communications and commitments give us a novel start. By assigning meaning to communications in terms of commitments, we accomplish the following. First, we reconstruct the correctness of behaviors by characterizing *compliance* as the eventual discharge of commitments. Second, we characterize the *interoperability* of agents as the alignment of their commitments, meaning that a creditor's expectations about a commitment are met

by the debtor. Third, we expand the treatment of design artifacts such as protocols by viewing them as communication types and showing how to refine and aggregate them. Using the above, we characterize the *conformance* of an agent with a role in a protocol. Further, we characterize important properties of a protocol such as its *transparency* in terms of the ability of the parties involved to verify each other's compliance. By contrast, traditional approaches (formal or otherwise) are largely confined to details such as message ordering and occurrence, and thus miss the forest for the trees. Throughout the paper, we discuss the properties with appropriate rigor, but without introducing formal notation.

Importantly, unlike most other multiagent systems work, our approach is undergirded by key ideas of distributed computing, especially dealing with the fact that key information is not immediately shared by all parties (even if they wish to share it). In fact, this is why protocols are important beyond plain commitments. This paper characterizes the above concepts under realistic assumptions, including multiparty settings with asynchronous communication (which aren't accommodated even in fairly recent interoperability research, e.g., [5–7]). Hence, this paper reflects crucial basic research not being addressed elsewhere. Its relevance to declarative agent languages and techniques arises from the fact that declarative representations for interaction are central to engineering robust, flexible multiagent systems, and this paper introduces and illustrates correctness criteria based on such declarative representations.

The rest of this paper is organized as follows. Section 2 introduces commitments and protocols in greater detail. Section 3 characterizes the correctness properties for interactions. Section 4 describes our contributions in relation to the most relevant literature.

## 2 Background on Protocols and Commitments

In classical software engineering methodologies, information modeling involves the application of key abstractions such as classification, aggregation, and association among components. It would be valuable to develop similar abstractions for interactions. Notice that traditional flow-based process models don't readily support such abstractions. One, existing work imposes severely limiting assumptions to support such abstractions—refinement is specified for Petri nets restricted to one input and one output place [8], which are not as expressive as general Petri nets needed to express real processes. Two, absent a business-level semantics, the models are rigid and any deviation would be potentially erroneous, thus making it difficult to refine or generalize processes.

By contrast, protocols focus on interactions, not on implementations. Our commitment-based semantics of protocols enables us to determine if a protocol refines another protocol, and how protocols may be aggregated into other protocols. Further, we specify a protocol primarily in terms of the vocabulary for communication that it defines and only secondarily in terms of (generally, ad hoc) constraints on the ordering and occurrence of messages. By basing correctness on the discharge of commitments, we enable agents to behave flexibly. For example, a merchant may ship before receiving payment if it wishes; a customer may pay directly or via a third party; and so on. On occasion, an application may impose an otherwise ad hoc constraint. For example, in a

(sit-down) restaurant, the protocol is to pay after food has been received and consumed; in a drive-through, payment precedes delivery. Such constraints often are merely guidelines for the participants and have no bearing on correctness unless they are enshrined in commitments. For example, a restaurant patron may pay early; a drive-through clerk may hand over the food before taking payment from the customer.

Flexible enactment and modeling in terms of refinement and aggregation are possible only because our formal semantics establishes the correctness criteria by which legitimate enactments, refinements, and aggregations can be identified [4]. Commitments express how contractual relationships form and progress during the agents' interactions. The commitment-based semantics is readily grounded via operational or messaging-level constraints [9].

*Commitments.* Contracts are key to flexible interoperation. Hohfeld [10] clarified a legal notion of contracts. Commitments, as introduced on Page 2, cover the relevant aspects of Hohfeld's notions [11], and thus naturally represent the contractual relationships of interest.

Two main forms of commitments arise [12]: *practical* commitments are about bringing about a future condition (i.e., oriented toward tasks), whereas *dialectical* commitments [13] are about staking a claim (as in argumentation) about the past, present, or future (i.e., oriented toward assertions). The distinction between them is significant even when directed to the future. For example, I might commit dialectically that the postman will ring twice, without committing practically to ensure that the postman rings twice. This paper deals with practical commitments except where specified otherwise. For example, the customer's agreement to pay the price for the book after it is delivered is a practical commitment that the customer (as debtor) has towards the bookstore (as creditor) to ensure the price is paid.

Using commitments enables us to model interactions *computation independently* (using this term as in Model-Driven Architecture (MDA) [14]). On the one hand, commitments describe the evolving state of the ongoing business interaction and how it evolves due to the participants' communications. On the other hand, commitments help express the expectations that participants have of one another: this is the fundamental purpose of a protocol. Jointly, these enable us to readily detect and accommodate business exceptions and opportunities. Consequently, commitments lend coherence to interactions [15].

Commitments can be manipulated through a small set of operations, including create, discharge, cancel, release, assign, and delegate [11], which we lack the space to discuss here. With additional assumptions, commitments can be enforced—by penalizing agents who do not comply with their commitments. Commitments have a formal semantics [12], and reasoning about commitment operations in distributed settings is addressed in [9].

*Protocols and commitments.* An advantage of incorporating commitments in our models is that they directly represent contractual relationships, are flexible, and lend coherence to the interactions of the participants in a process. The formalization of the specialization and generalization hierarchy of protocols is made the more interesting and

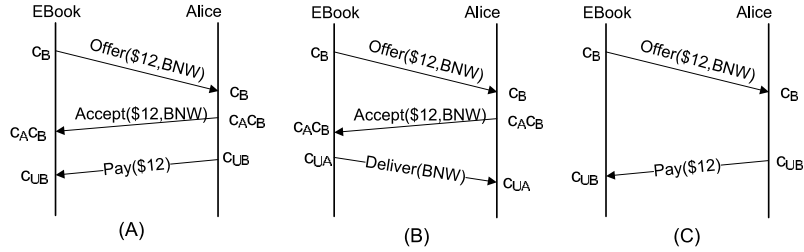
useful because of the presence of commitments and roles in our model. Instead of considering uninterpreted runs (of actions and states), we consider how the commitments of the various roles evolve over different runs. The use of commitments enables more sophisticated reasoning about meaning than in traditional approaches. In particular, it enables us to characterize the similarity of states and refinement protocols in potentially subtle ways. An example is when a participant from its local perspective considers two states as interchangeable simply because it features as the creditor and debtor in the same commitments regardless of the other parties. For instance, in some settings, Alice may care only of her total accounts receivable, and not care if it is Bob or Charlie who is committed to paying her the money. In other words, instead of merely considering raw computations, it makes sense to “normalize” them in terms of commitments so as to make more precise judgments about how protocols relate to one another.

**Table 1.** A purchase protocol (customer is  $c$  and merchant is  $m$ )

---

$Offer(m, c, payment, book)$	means	$Create(m, c, payment, book)$
$Accept(c, m, payment, book)$	means	$Create(c, m, book, payment)$
$Reject(c, m, payment, book)$	means	$Release(m, c, payment, book)$
$Deliver(m, c, book)$	means	$Inform(m, c, book)$
$Pay(c, m, payment)$	means	$Inform(c, m, payment)$

---



**Fig. 1.** Three possible enactments of protocol in Table 1

Table 1 shows the key messages in a purchase protocol and their meanings. Figure 1 shows some possible enactments of the protocol between a customer Alice and a merchant EBook (in the figure,  $c_A$  is  $C(Alice, EBook, BNW, 12)$ ;  $c_{UA}$  is  $C(Alice, EBook, \top, 12)$ ;  $c_B$  is  $C(EBook, Alice, 12, BNW)$ ;  $c_{UB}$  is  $C(EBook, Alice, \top, BNW)$ ). The meanings of the messages are crucial, because they help characterize the protocol declaratively. The meanings are systematically formalized in a declarative action language. Our language and technique are introduced in [16–18].

Traditional approaches force a tradeoff: checking compliance is simple with rigid automaton-based representations and difficult with flexible unconstrained reasoning agents. Commitments help us find the happy middle: protocols maximize flexibility by constraining the participants’ interactions at the business level, yet provide a crisp notion of compliance: a party complies if its commitments are discharged, no matter if delegated or otherwise manipulated.

*Protocols and computations.* In essence, each protocol allows a set of computations or runs, each run being an alternative that requires a specific sequence of executions upon the participants. Two basic intuitions about protocol refinement are that (1) a more general protocol includes additional runs (more ways to satisfy) beyond those in a less general protocol; and (2) a more general protocol includes shorter runs (fewer steps to satisfy) than a less general protocol.

Our commitment-based semantics yields a rigorous basis for protocol *refinement* and *aggregation* [19]. In principle, these properties enable reusing protocols from a repository. For example, PAYMENT BY CHECK refines PAYMENT. Further, ORDER, PAYMENT, and SHIPPING can be combined into a new protocol for PURCHASE. This composed protocol would capture the reusable interactions and service agreements that underlie a business process. For example, PURCHASE would specify how orders may be placed, payments made, and shipping arranged. When protocols are composed, so are the roles; e.g., the payer in PAYMENT may be composed with the receiver in SHIPPING. Multiple copies of the same protocol may be composed: in an ARBITRAGE protocol, the arbitrageur role would compose the seller role in one copy of PAYMENT with the buyer role in the second copy.

As in other formal semantics, the runs are merely abstract entities used to establish logical properties. We would never explicitly enumerate the potentially infinite number of possible runs, but we can use the abstract definition to show important algebraic relationships. Mallya & Singh [19] show important progress, but their approach is far from complete. Specifically, it deals with sets of runs, but does not apply directly on protocol specifications as one would find in a repository.

### 3 Correctness Properties

We begin by motivating some key definitions. Notice that although the above discussion uses protocols as design artifacts, compliance and interoperability apply without regard to any protocol. Although our main definitions and methods are oriented toward commitments, they are undergirded by considerations of distributed computing, especially of asynchrony in messaging.

#### 3.1 Interoperability

The *interoperability* of a set of roles or agents, regardless of protocol, means that they jointly meet the expectations they jointly place on each other. Some aspects of interoperability depend on meanings; others on the messaging system that underlies communications.

We assume that messaging is asynchronous, reliable, and pairwise (for each sender and receiver) order-preserving: this matches what emerging middleware standards [20] offer. Thus in two-party cases, each party would eventually learn of the relevant moves and expectations of the other: the only kind of pathology possible is that the parties may view some pairs of messages in opposite orders. In multiparty cases, the messaging conditions can become more subtle: e.g., a party would lack direct information about



messages exchanged among other parties. Mostly, this is a good thing because the parties can proceed with minimal mutual dependencies. However, when such information materially affects a desired property, we would need to change either the requirements (so information about remote events becomes irrelevant) or the specification (so that the necessary information flows to the appropriate parties).

Interoperation classically is treated as a conjunction of liveness and safety. To these we add alignment.

**Liveness** means that progress will take place: desirable states will be visited infinitely often. Liveness can fail if a receiver blocks (awaiting a message that is never sent). For example, let Buyer-A demand delivery before payment and Seller-A demand payment before delivery. Now, Buyer-A and Seller-A would deadlock, each awaiting the other's message.

**Safety** means that the system doesn't enter an undesirable state: agents must be ready to receive the messages being sent to them. Safety is best understood in a multiparty setting. If a buyer expects to receive a confirmation before a shipment but receives them in the opposite order, its resultant state is not defined. We should ensure the messages occur in only those orders that the buyer accepts.

We apply causality [21] to model the above concepts. The sending of a message is causally prior to its receipt; for any two locally ordered events (sends or receives), the first is (potentially) causally prior to the second: "potential" because from external observations we cannot infer if the two events are truly related. We can infer true causality from the agents' specifications, in settings where the specifications are available. We can characterize liveness and safety in terms of the compatibility among causal orders involving receives and sends. We conjecture that the above will yield superior solutions to those in the recent distributed computing literature, e.g., [5–7]. The literature considers two-party cases or violates substitutability: that substituting an agent with a conforming agent must preserve interoperability.

**Alignment** is interoperability with respect to expectations at the level of meaning: do the participants agree about the states of their commitments to each other? A set of agents or roles is *aligned* provided throughout any enactment, whenever one concludes it is the creditor of a commitment, the corresponding debtor  $x$  concludes that  $x$  is the debtor of the commitment [22]. In other words, the debtor recognizes a commitment that the creditor expects of it. How commitments are created, discharged, and manipulated depends on the messages sent and received.

From the point of view of interoperability, interesting agent specifications are of two kinds: constitutive and regulative [22]. An agent's constitutive specification deals only with the meaning of messages. In other words, it specifies what messages *count as* for the agent. An agent's regulative specification, in contrast, describes agent behavior; i.e., it describes the conditions under which the agent sends and receives particular messages. Regulative specifications are thus closer to implementations.

Judging the constitutive alignment of a set of agents by statically analyzing their specifications is nontrivial because message meanings are conditional, and thus potentially affected by how other messages change the relevant conditions. For example, if one message constitutes an authorization and the meaning of a second message relies

upon that authorization, the commitments resulting from the second message would depend upon whether the first message precedes it. For a two-party example, if a buyer and seller interpret the above quote message as different commitments, they would be misaligned [22] even though they satisfy safety.

In multiparty situations, a debtor and creditor's conclusions about a commitment may conflict because they see different messages occurrences or orders. Delegations and assignments of commitments inherently involve three parties: thus these cases are crucial. We propose to extract causality information from the specifications to use as explained above.

A specification may fail safety or liveness without failing alignment. For example, let protocol FLEXIBLE PURCHASE allow a payment to occur before or after the delivery of goods. It is easy to see that Buyer-A and Seller-A (introduced above), respectively, conform to the customer and merchant roles in FLEXIBLE PURCHASE. We saw above that Buyer-A and Seller-A fail liveness. However, they never disagree about their commitments and hence would satisfy alignment.

### 3.2 Conformance and Operability

Conformance and operability apply to each interoperability property: liveness, safety, and alignment. A role *conforms* to, i.e., is a subtype of, another role provided the first role meets all expectations placed on the second and holds no expectations of others beyond what the second does. Similarly, an agent conforms to, i.e., instantiates, a role. Conformance is important because it helps us build a library of roles without which engineering would lapse into one-off solutions. To handle conformance properly would require considering the semantics of protocols not in terms of simple runs, but in terms of the choices they afford each role. Echoing the intuition of alternating refinement [23], expectations placed on a role correspond to “external” choices; expectations held by a role correspond to “internal” choices.

A protocol is *operable*, i.e., potentially enactable, if the roles it specifies are interoperable. A protocol may fail to be operable when it requires a role to act based on events that the role cannot observe. Operability is an important quality criterion for protocols: ideally, the protocols in a library should be operable, so developers may implement selected roles conformantly, and be assured of interoperation.

Recall that Buyer-A and Seller-A together fail liveness even though they conform to roles in FLEXIBLE PURCHASE. Hence FLEXIBLE PURCHASE is not operable for liveness. Conversely, let PREPAID PURCHASE require payment to occur before delivery. Then, any pair of conforming customer and merchant would be live and safe. Hence, PREPAID PURCHASE is operable. Buyer-A is nonconformant with the customer role, whereas Seller-A is conformant with the merchant role of PREPAID PURCHASE. Seller-A and Buyer-A failing liveness doesn't mean PREPAID PURCHASE is inoperable: it is Buyer-A that is messed up.

### 3.3 Compliance and Transparency

*Compliance* means that each agent performs as expected by others, by discharging its commitments. We can prove compliance only when we know each agent's specifica-

tion and relevant assumptions about the environment hold. That is, compliance can be verified for specific runs but not proved in general for open systems [24]. Notice that alignment and compliance are independent of each other: e.g., an interoperable buyer may be committed to pay, but may refuse to do so. An agent may *verify* a debtor’s compliance based on its observations in a specific enactment. Assuming that the discharge of a commitment is observable (e.g., occurs via a message), verifying compliance is simple in two-party cases. If a debtor complies, the creditor would eventually know. If a debtor does not comply, then the creditor would eventually know—provided the commitment includes a deadline. In multiparty cases, a creditor may lack some important observations, and hence special techniques would be required to verify alignment.

A protocol is *transparent* if each role in it can verify the compliance of its debtors. However, not all protocols enable each role to verify compliance at runtime: a protocol may be such that “news” relevant to a commitment might not be propagated to the creditor. Transparency is an important quality criterion for protocols: it ensures that participants can verify if others are not complying.

### 3.4 Refinement and Compatibility

The *refinement* of a protocol by another protocol means that the second protocol generates only computations that are allowed by the first. Modeling via commitments enables us to finesse the intuitions about protocol refinement. For example, a simple PAYMENT protocol might require that the payer transfer funds to the payee. A particular refinement of this might be PAYMENT WITH A CHECK. To pay with a check, the payer would send a check to the payee, who would deposit the check at his bank, which would present it to the payer’s bank, which would send the funds to the payee’s bank, which would make those funds available to the payee. Thus PAYMENT BY CHECK is a specialization of PAYMENT, but it involves additional roles and steps, and skips some of the steps of PAYMENT, e.g., direct transfer. With a commitment-based definition, we can formally establish that PAYMENT BY CHECK refines PAYMENT—something that would not be possible with traditional approaches because of the above differences between the two protocols. The key intuition is that the commitments at critical states line up correctly. This is a significant departure from traditional notions of refinement which, because they lack commitments, insist upon the computations to match in their detailed steps.

Notice that an agent designed to play a role in a refined protocol may not comply with any role in the original protocol. This is because the agent may not interpret messages in a way compatible with the original protocol. For example, in PAYMENT BY CHECK, a merchant may interpret a check as being adequate as a receipt (once it is cleared and returned to the customer by the customer’s bank), but the customer may not interpret it like that and may continue to expect a separate receipt as in PAYMENT. Further, the agent may fail to interoperate with roles defined in the original protocol. This is because it may send messages that are not defined in the original protocol. In general we would not be able to substitute a role from a refined protocol for a role in the original protocol. The foregoing is motivation for the property of *compatibility*, which determines if roles in one protocol conform to roles in another protocol.

Table 2 summarizes the above properties. With the exception of compliance, these properties can be verified by a static analysis of the appropriate declarative specifica-

**Table 2.** The properties summarized

Property	Of What?
Refinement, compatibility, operability, transparency	Protocols
Interoperability (safety, liveness, or alignment)	Agents and roles
Conformance	Roles
Compliance	Agents

tions. Compliance can realistically be determined only at runtime: we expect that for even the simplest practically useful agents, the complexity of checking agent specifications for compliance would make it infeasible.

#### 4 Discussion: Relevant Literature

Our main contribution in this paper is in characterizing the key correctness properties that would support an interaction-oriented approach to building software systems, particularly cross-organizational business processes. In particular, the correctness properties reflect high-level requirements of such systems.

Interestingly, Parnas [25] proposed early in the study of software architectures that connectors be treated not as control or data flow constructs but as *assumptions* made by each component about the others. Arguably, much of the subsequent work on software architecture regressed from Parnas' insight: it has primarily considered connectors at the level of flow, e.g., dealing exclusively with message order and occurrence [26]. In formulating the assumptions at a high level, we see a great opportunity for multiagent systems research to address some of the long-standing challenges in software.

*Conventional formal methods.* Current modeling formalisms, such as finite state machines and Petri Nets, originated in distributed computing and apply at lower levels of abstraction than needed for flexible business interactions [27, 8]. When applied to business protocols, these formalisms result in specifications that are over-constrained to the level of specific sequences of actions. Recent approaches have sought to express scheduling requirements declaratively, via temporal logic [28–30]. Although they are more flexible and modular than operational representations, these approaches do not express business semantics.

FIPA, the Foundation for Intelligent and Physical Agents (now part of IEEE) recognized the importance of reusable interaction protocols in the late 1990s [31]. Odell *et al.* [32] give one of the earliest uses of UML for protocols. They show how various UML diagrams can be applied for modeling agent interactions. This work shows about how far you can go in a conventional software framework, and has inspired our work. The present paper is about fundamental enhancements to conventional models to capture protocols and their commitment-based semantics.

Leading approaches model conversations via finite-state machines and establish properties such as how roles may realize a protocol or a protocol subsumes another [33, 34]. Dastani *et al.* [35] show how to model a rich family of coordination connectors for

multiagent systems. Honda *et al.* [36] develop a type theory that would support multi-party sessions: in essence this would help robustly generate roles. The above works are valuable but lack the business-level semantics that distinguishes our work.

The term “protocol” sometimes refers to the constraints on one agent’s behavior [5]—this is a *role* in our terminology. Because existing approaches do not characterize the business semantics as such, they do not address the challenges described above. However, their treatment of messages and computations at a low level is useful, and complementary to our work.

Whereas deontic logic only deals with what is obligatory or permissible and thus disregards an agent’s obligations *to* another agent, commitments are directed and context sensitive. Commitments include support for a variety of operations [11, 37]. Foster *et al.* [38] seek to capture the semantics of process interactions via the notion of obligation policies. They represent choreographies via message sequence charts and state machines. Obligations are rather weak in their formulation, however. Specifically, obligations are not reified, and cannot be manipulated to capture flexible interactions among independent parties. However, Foster *et al.*’s work illustrates an encouraging trend in software engineering, namely, the expression of richer semantics in understanding interactions among autonomous components. The present paper contributes to this trend. Lomuscio *et al.* [39] formalize correctness properties in a temporal logic and show how to verify them. They consider obligations but do not consider commitments as here. Lomuscio *et al.* also do not consider the rich set of properties defined above, concentrating on only one, which is somewhat like compliance.

*Business processes.* The MIT Process Handbook (MITPH) [40] is of great relevance intellectually. MITPH includes an extensive classification and systematic organization of business processes based on two dimensions of process hierarchies, one that composes the *uses* of a process out of its constituent *parts*, and another that subclasses *generalizations* of a process into *specializations*. Our work can provide the rigorous underpinnings for work such as the MITPH. Grosz and Poon [41] develop a system to represent and execute business rules from MITPH. Wyner and Lee [42] study specialization for data flow diagrams. Their approach can form the basis of the processes identified in MITPH. These concepts turn out to be complex and not readily applied to entire business processes. Further, since Wyner and Lee do not capture the content through a high-level representation such as commitments, the results are not intuitive.

Our approach agrees with the newer declarative forms of artifacts-based process modeling [43] in terms of deemphasizing low-level operational details in favor of business semantics. However, these approaches do not have a central organizing principle on par with commitments, and thus do not offer a generic and flexible basis for determining the properties we introduced above.

*Agent communications.* Fornara and Colombetti [44] describe how commitments relate to FIPA messages, demonstrating this with an example. Rovatsos [45] proposes a formal operational semantics for communications based on commitments under synchronous messaging. His approach violates autonomy by legislating agent behaviors from within the language specification: this level of prescription is ill-suited to most multiagent applications.

Yolum and Singh [46] [47] offer one of the first accounts of the use of commitments in modeling protocols to improve flexibility for participating agents, which was enhanced by Winikoff *et al.* [48]. Johnson *et al.* [49] develop a scheme for identifying when two commitment-based protocols are equivalent. Their scheme, however, is simplistic, classifying protocols based solely on their syntactic structure. Our work provides stronger results from an application point of view and relates better to Web services.

Commitments have found application in formalizing argumentation, e.g., [50, 51]. Usually, though, this work makes simplifying assumptions such as (1) maintaining a unique commitment store; (2) informally specifying the meanings of communicative acts as effects on the store; (3) assuming synchronous two-party communications.

*Agent-oriented software engineering (AOSE)*. A number of useful software methodologies for building multiagent systems for IT applications have emerged that incorporate rich metamodels and describe how to build a series of software artifacts [52, 53, 3]. Garcia-Ojeda *et al.* [54] synthesize existing metamodels into a comprehensive metamodel of organizations geared toward process modeling. We recently developed Amoeba, a protocol-based methodology compatible with the ideas of this paper [1].

The above methodologies address the challenges of autonomy and heterogeneity by giving prominence to communication. Such works are clearly valuable and worthwhile. However, current approaches do not consider the full subtleties both of meaning and of distribution. By contrast, this paper addresses the foundations for business interactions understood in terms of commitments. The proposed definitions will offer a foundations for building a new family of tools that, in principle, could be used within any of the above methodologies, because they all support aspects of interaction and of agents playing roles in interactions.

## 5 Conclusions and Directions

This paper presents a key step in our program of research to develop underpinnings of multiagent systems—and indeed, of all software—on interactive grounds with an emphasis on declarative formulations. The main point to take away is the richness of the correctness properties. These properties echo well-known conventional properties but their formulation in a declarative, interactive setting adds a lot of subtlety that traditional approaches cannot express. The foregoing leads to two broad questions.

- *Theory*. What are practical decision algorithms for these properties? How can we specify agents who may play specified roles (while applying their local policies)? How can we determine that agents (supposedly) enacting a protocol are complying with the protocol? What are practical algorithms for judging the varieties of interoperability, conformance, operability, compliance, and transparency?
- *Suitability and applicability*. Does representing meaning via commitments provide a sufficiently natural basis for business interoperation? How readily can meaning be associated with tools to engineer and use protocols? Can we specify commitments sufficiently precisely in real-life business settings? How can we use the above properties and algorithms to enable protocol design and agent implementation?

Our future work lies in answering the above questions.

## References

1. Desai, N., Chopra, A.K., Singh, M.P.: Amoeba: A methodology for modeling and evolution of cross-organizational business processes. *ACM Transactions on Software Engineering and Methodology (TOSEM)* (2009) To appear.
2. Desai, N., Chopra, A.K., Arrott, M., Specht, B., Singh, M.P.: Engineering foreign exchange processes via commitment protocols. In: *Proceedings of the 4th IEEE International Conference on Services Computing (SCC)*, IEEE Computer Society Press (2007) 514–521
3. Bresciani, P., Perini, A., Giorgini, P., Giunchiglia, F., Mylopoulos, J.: Tropos: An agent-oriented software development methodology. *Journal of Autonomous Agents and Multi-Agent Systems* **8** (2004) 203–236
4. Desai, N., Mallya, A.U., Chopra, A.K., Singh, M.P.: Interaction protocols as design abstractions for business processes. *IEEE Transactions on Software Engineering* **31** (2005) 1015–1027
5. Baldoni, M., Baroglio, C., Martelli, A., Patti, V.: A priori conformance verification for guaranteeing interoperability in open environments. In: *Proceedings of the 4th International Conference on Service-Oriented Computing (ICSOC)*. Volume 4294 of LNCS., Springer (2006) 339–351
6. Fournet, C., Hoare, C.A.R., Rajamani, S.K., Rehof, J.: Stuck-free conformance. In: *Proceedings of the 16th International Conference on Computer Aided Verification (CAV)*. Volume 3114 of LNCS., Springer (2004) 242–254
7. Giordano, L., Martelli, A.: Verifying agent conformance with protocols specified in a temporal action logic. In: *Proceedings of the 10th Congress of the Italian Association for Artificial Intelligence (AI\*IA)*. Volume 4733 of LNCS., Springer (2007) 145–156
8. van der Aalst, W., van Hee, K.: *Workflow Management Models, Methods, and Systems*. MIT Press, Cambridge, MA (2002)
9. Chopra, A.K., Singh, M.P.: Multiagent Commitment Alignment. *Proceedings of the 8th International Joint Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*, (2009)
10. Hohfeld, W.N.: *Fundamental Legal Conceptions as Applied in Judicial Reasoning and other Legal Essays*. Yale University Press, New Haven, CT (1919).
11. Singh, M.P.: An ontology for commitments in multiagent systems: Toward a unification of normative concepts. *Artificial Intelligence and Law* **7** (1999) 97–113
12. Singh, M.P.: Semantical considerations on dialectical and practical commitments. In: *Proceedings of the 23rd Conference on Artificial Intelligence (AAAI)*, (2008) 176–181
13. Walton, D.N., Krabbe, E.C.W.: *Commitment in Dialogue: Basic Concepts of Interpersonal Reasoning*. State University of New York Press, Albany (1995)
14. OMG: The Object Management Group’s Model Driven Architecture (MDA) (2006) <http://www.omg.org/mda/>.
15. Jain, A.K., Aparicio IV, M., Singh, M.P.: Agents for process coherence in virtual enterprises. *Communications of the ACM* **42** (1999) 62–69
16. Chopra, A.K., Singh, M.P.: Contextualizing commitment protocols. In: *Proceedings of the 5th International Joint Conference on Autonomous Agents and Multiagent Systems*. (2006) 1345–1352
17. Desai, N., Singh, M.P.: A modular action description language for protocol composition. In: *Proceedings of the 22nd Conference on Artificial Intelligence (AAAI)*, (2007) 962–967
18. Desai, N., Singh, M.P.: On the enactability of business protocols. In: *Proceedings of the 23rd Conference on Artificial Intelligence (AAAI)*, (2008) 1126–1131
19. Mallya, A.U., Singh, M.P.: An algebra for commitment protocols. *Journal of Autonomous Agents and Multi-Agent Systems* **14** (2007) 143–163

20. AMQP: Advanced message queuing protocol (2007) <http://www.amqp.org>.
21. Lamport, L.: Time, clocks, and the ordering of events in a distributed system. *Communications of the ACM* **21** (1978) 558–565
22. Chopra, A.K., Singh, M.P.: Constitutive interoperability. In: *Proceedings of the 7th International Joint Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*. (2008) 797–804
23. Alur, R., Henzinger, T.A., Kupferman, O., Vardi, M.Y.: Alternating refinement relations. In: *Proceedings of the 9th International Conference on Concurrency Theory (CONCUR)*. Volume 1466 of LNCS., Springer (1998) 163–178
24. Venkatraman, M., Singh, M.P.: Verifying compliance with commitment protocols: Enabling open Web-based multiagent systems. *Autonomous Agents and Multi-Agent Systems* **2** (1999) 217–236
25. Parnas, D.L.: Information distribution aspects of design methodology. In: *Proceedings of the International Federation for Information Processing Congress*. Volume TA-3., Amsterdam, North Holland (1971) 26–30
26. Hohpe, G., Woolf, B.: *Enterprise Integration Patterns: Designing, Building, and Deploying Messaging Solutions*. Signature Series. Addison-Wesley, Boston (2004)
27. Harel, D., Gery, E.: Executable object modeling with statecharts. *IEEE Computer* **30** (1997) 31–42
28. Singh, M.P.: Distributed enactment of multiagent workflows: Temporal logic for service composition. In: *Proceedings of the 2nd International Joint Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*, (2003) 907–914
29. Pesic, M., Schonenberg, M.H., Sidorova, N., van der Aalst, W.M.P.: Constraint-based workflow models: Change made easy. In: *Proceedings of the On the Move to Meaningful Internet Systems (Confederated International Conferences CoopIS, DOA, ODBASE, GADA, and IS)*, Part I. Volume 4803 of LNCS., Springer (2007) 77–94
30. Wu, Q., Pu, C., Sahai, A., Barga, R.S.: Categorization and optimization of synchronization dependencies in business processes. In: *Proceedings of the 23rd International Conference on Data Engineering (ICDE)*, IEEE (2007) 306–315
31. FIPA: FIPA interaction protocol specifications (2003) FIPA: The Foundation for Intelligent Physical Agents, <http://www.fipa.org/repository/ips.html>.
32. Odell, J., Parunak, H.V.D., Bauer, B.: Representing agent interaction protocols in UML. In: *Proceedings of the 1st International Workshop on Agent-Oriented Software Engineering (AOSE)*. (2001)
33. Benatallah, B., Casati, F., Toumani, F.: Representing, analysing and managing web service protocols. *Data & Knowledge Engineering* **58** (2006) 327–357
34. Bultan, T., Fu, X., Hull, R., Su, J.: Conversation specification: A new approach to design and analysis of e-service composition. In: *Proceedings of the Twelfth International World Wide Web Conference (WWW)*. (2003) 403–410
35. Dastani, M., Arbab, F., de Boer, F.S.: Coordination and composition in multi-agent systems. In: *Proceedings of the 4rd International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, (2005) 439–446
36. Honda, K., Yoshida, N., Carbone, M.: Multiparty asynchronous session types. In: *Proceedings of the 35th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, ACM (2008)
37. Singh, M.P., Chopra, A.K., Desai, N.: Commitment-based SOA. *IEEE Computer* **42** (2009) Accepted; available from <http://www.csc.ncsu.edu/faculty/mpsingh/papers/>.
38. Foster, H., Uchitel, S., Magee, J., Kramer, J.: Model-based analysis of obligations in web service choreography. In: *Proceedings of the Advanced International Conference on Telecommunications and International Conference on Internet and Web Applications and Services (AICT-ICIW)*. (2006) 149–156



39. Lomuscio, A., Qu, H., Solanki, M.: Towards verifying compliance in agent-based web service compositions. In: Proceedings of the 7th International Joint Conference on Autonomous Agents and MultiAgent Systems (AAMAS) (2008) 265–272
40. Malone, T.W., Crowston, K., Herman, G.A., eds.: Organizing Business Knowledge: The MIT Process Handbook. MIT Press, Cambridge, MA (2003)
41. Grosz, B.N., Poon, T.C.: SweetDeal: Representing agent contracts with exceptions using XML rules, ontologies, and process descriptions. In: Proceedings of the 12th International Conference on the World Wide Web. (2003) 340–349
42. Wyner, G.M., Lee, J.: Defining specialization for process models. In: [40]. MIT Press (2003) 131–174
43. Hull, R.: Artifact-centric business process models: Brief survey of research results and challenges. In: On the Move to Meaningful Internet Systems. Volume 5332 of LNCS., Springer (2008) 1152–1163
44. Fornara, N., Colombetti, M.: Defining interaction protocols using a commitment-based agent communication language. In: Proceedings of the 2nd International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS), (2003) 520–527
45. Rovatsos, M.: Dynamic semantics for agent communication languages. In: Proceedings of the 6th international Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS). (2007) 1–8
46. Yolum, P., Singh, M.P.: Commitment machines. In: Proceedings of the 8th International Workshop on Agent Theories, Architectures, and Languages (ATAL-01). Volume 2333 of LNAI., Springer-Verlag (2002) 235–247
47. Yolum, P., Singh, M.P.: Flexible protocol specification and execution: Applying event calculus planning using commitments. In: Proceedings of the 1st International Joint Conference on Autonomous Agents and MultiAgent Systems (AAMAS), (2002) 527–534
48. Winikoff, M., Liu, W., Harland, J.: Enhancing commitment machines. In: Proceedings of the 2nd International Workshop on Declarative Agent Languages and Technologies (DALT). Volume 3476 of LNAI., Berlin, Springer-Verlag (2005) 198–220
49. Johnson, M.W., McBurney, P., Parsons, S.: When are two protocols the same? In Huget, M.P., ed.: Communication in Multiagent Systems: Agent Communication Languages and Conversation Policies. Volume 2650 of LNAI. Springer-Verlag, Berlin (2003) 253–268
50. Amgoud, L., Maudet, N., Parsons, S.: An argumentation-based semantics for agent communication languages. In: Proceedings of the 15th European Conference on Artificial Intelligence (ECAI), IOS Press (2002) 38–42
51. Norman, T.J., Carbogim, D.V., Krabbe, E.C.W., Walton, D.N.: Argument and multi-agent systems. In Reed, C., Norman, T.J., eds.: Argumentation Machines. Kluwer (2004)
52. Bergenti, F., Gleizes, M.P., Zambonelli, F., eds.: Methodologies and Software Engineering for Agent Systems. Kluwer, Boston (2004)
53. Henderson-Sellers, B., Giorgini, P., eds.: Agent-Oriented Methodologies. Idea Group, Hershey, PA (2005)
54. Garcia-Ojeda, J.C., DeLoach, S.A., Robby, Oyenon, W.H., Valenzuela, J.: O-MaSE: A customizable approach to developing multiagent processes. In: Proceedings of the 8th International Workshop on Agent Oriented Software Engineering (AOSE 2007). (2007)

# Reasoning and Planning with Cooperative Actions for Multiagents Using Answer Set Programming

Tran Cao Son<sup>1</sup> and Chiaki Sakama<sup>2</sup>

<sup>1</sup> Dept. of Computer Science, New Mexico State University, Las Cruces, NM 88003, USA  
t.son@cs.nmsu.edu

<sup>2</sup> Computer and Communication Sciences, Wakayama University, Wakayama 640-8510, Japan  
sakama@sys.wakayama-u.ac.jp

**Abstract.** In this paper, we investigate the multiagent planning problem in the presence of cooperative actions and agents, which have their own goals and are willing to cooperate. To this end, we extend the action language  $\mathcal{A}$  in [12] to represent and reason about plans with cooperative actions of an individual agent operating in a multiagent environment. We then use the proposed language to formalize the multiagent planning problem and the notion of a joint plan for multiagents in this setting. We discuss a method for computing joint plans using answer set programming and provide arguments for the soundness and completeness of the implementation.

## 1 Introduction

*Cooperative actions* are actions of an agent which can be executed only if the agent is operating in a multiagent environment. They can be actions for soliciting something from other agents or actions for setting up some conditions for other agents. They differ from individual actions in that they might affect other agents. Cooperative actions are important not only in situations where multiple agents have to work together to accomplish a common goal but also in situations where each agent has its own goal. This can be seen in the following story, a modified version of the story in [21]:

*Example 1.* Three new students  $A$ ,  $B$ , and  $C$  are moving in a shared apartment and planning to decorate their rooms. Each would like to hang one of their objects on the wall, e.g.,  $A$  would like to hang a mirror,  $B$  a diploma, and  $C$  a painting.  $A$  and  $B$  know how to use either a nail or a screw to complete their job but  $C$  knows to use the screw only.  $A$  has neither a nail or a screw.  $B$  has both.  $C$  has only a nail. To use a nail, one will need a hammer. Among three, only  $B$  has a hammer.

Do the students have a joint-plan that allows each of them to achieve his/her goal?

Intuitively, we can see that only  $B$  can accomplish her job independent of  $A$  and  $C$ . The three can achieve their goals if  $B$  uses the hammer and the nail to hang her diploma then gives  $A$  the hammer and  $C$  the screw, respectively.  $C$ , on the other hand, gives  $A$  the nail and uses the screw to hang her painting.  $A$  uses the nail (from  $C$ ) and the hammer (from  $B$ ) to hang her mirror. Of course, to avoid unpleasant moments,  $A$  should ask for the nail (from  $C$ ) and the hammer (from  $B$ ) and  $C$  should ask for the screw (from  $B$ ).

However, it is easy to see that if either  $B$  or  $C$  does not want to give out anything, then only  $B$  can achieve her goal. Furthermore, if  $B$  decides to use the screw instead of using the nail in hanging her diploma, then  $C$  has no way of achieving her goal.  $\square$

In the above example, the action of giving a nail, a hammer, or a screw between the students can be considered as cooperative actions. The action of requesting something from others can also be considered as cooperative actions. It is obvious that without some cooperative actions, not all students can achieve their own goals. Even with the cooperative actions at their disposal, the students might still need to coordinate in creating their corresponding plans.

In Example 1, agents (the students) maintain their own local worlds and their actions do generally not affect others' worlds. It should be emphasized that the fact that agents have their own world representation does not exclude the situations in which the worlds of different agents overlap and the execution of one agent's individual actions might affect others as well or the execution of their joint-action.

*Example 2.* Let us consider  $A$  and  $B$  who are in one room and studying at their tables. Each of them sits next to a switch which can control the lamp in the room. Flipping either switch will change the status of the light.

Assume that  $A$  and  $B$  maintain their world representation separately. (They might use the same theory for this purpose but we will not impose this.) Obviously, if  $A$  flips the switch next to her, the world in which  $B$  is in will also change.

Similarly, if  $A$  and  $B$  lift a table and place it at different location, their joint-action change the world of both as well.  $\square$

In this paper, we will consider multiagent planning problems in which each agent maintains its own representation about the world and its capabilities, which includes individual actions and cooperative actions; and has its own goal. We are mainly interested in the process of creating a joint plan prior to its execution. We will begin by extending the language  $\mathcal{A}$  in [12] to allow cooperative actions for a single agent. The semantics of the new language is defined by a transition function which maps pairs of actions and states to states. We then define the multiagent planning problems and the notion of a joint plan for multiagents in presence of cooperative actions. Finally, we discuss a method for computing joint plans using answer set programming [18, 19].

## 2 An action language with cooperative actions

In this section, we present a language for representing and reasoning about plans for an agent in the multiagent environment with cooperative actions. To this end, we extend the language  $\mathcal{A}$  in [12] to allow cooperative actions<sup>1</sup>. In this paper, we consider cooperative actions as actions that an agent would not have if she were in a single agent environment. Specifically, we consider two types of cooperative actions, one that requests the establishment of a condition in an agent's world and another establishes some conditions in the world of another agent. We will assume an arbitrary but fixed set of agent identifiers  $\mathcal{AG}$ . A planning problem of an agent in  $\mathcal{AG}$  is defined over a set of fluents (or state variables)  $F$ , a set of *individual actions*  $A$ , and a set of *cooperative*

---

<sup>1</sup> The choice of  $\mathcal{A}$  will be discussed in Section 5.

actions C. We will assume that A always contains a special action `wait` which does not have any effect on the agent's world<sup>2</sup>. Furthermore, we will require that actions in C do not appear in A. This highlights the fact that the cooperative actions are presented due to the presence of other agents.

A *fluent literal* (or *literal*) is either a fluent or its negation. Fluent formulas are propositional formulas constructed from literals and propositional connectives.

## 2.1 Specifying Individual Actions

A *domain specification*  $DI$  over  $F$  and  $A$  describes the individual actions of an agent and consists of laws of the following form:

$$a \text{ causes } l \text{ if } \phi \quad (1)$$

$$a \text{ executable } \phi \quad (2)$$

where  $a$  is an individual action (in  $A$ ),  $l$  is a fluent literal and  $\phi$  is a set of fluent literals.

A law of the form (1), called a *dynamic law*, states that if  $a$  is executed when  $\phi$  is true then  $l$  becomes true. (2) is an *executability condition* and says that  $a$  can be executed only if  $\phi$  is true. The semantics of a domain specification is defined by the notion of *state* and by a *transition function*  $\Phi$ , that specifies the result of the execution of an action  $a$  in a state  $s$ .

A set of literals  $S$  satisfies a literal  $l$  ( $l$  holds/is true in  $S$ ), denoted by  $S \models l$ , if  $l \in S$ . For a set of literals  $\phi$ ,  $S \models \phi$  if  $S \models l$  for every  $l \in \phi$ . A *state*  $s$  is a set of fluent literals that is *consistent* and *complete*, i.e., for every  $f \in F$ , either  $f \in s$  or  $\neg f \in s$  but  $\{f, \neg f\} \not\subseteq s$ . In the following,  $\bar{l}$  denotes the negation of  $l$ , i.e., if  $l = f$  and  $f \in F$ , then  $\bar{l} = \neg f$ ; if  $l = \neg f$  for some  $f \in F$ , then  $\bar{l} = f$ . For a set of literals  $S$ ,  $\bar{S} = \{\bar{l} \mid l \in S\}$ .

An action  $a$  is *executable* in a state  $s$  if there exists an executability condition ( $a$  **executable**  $\phi$ ) in  $DI$  such that  $s \models \phi$ .

Let  $e_a(s) = \{l \mid \exists (a \text{ causes } l \text{ if } \phi) \in DI. [s \models \phi]\}$ . The result of the execution of  $a$  in  $s$  is defined by

- $\Phi(a, s) = \text{fails}$  if  $a$  is not executable in  $s$ ; and
- $\Phi(a, s) = (s \setminus \bar{e}_a(s)) \cup e_a(s)$  if  $a$  is executable in  $s$ .

A domain specification  $DI$  is *consistent* if  $\Phi(a, s) \neq \text{fails}$  holds for every pair of action  $a$  and state  $s$  such that  $a$  is executable in  $s$ .

$\Phi$  is extended to reason about effect of a sequence of actions as follows.

**Definition 1 (Transition function).** Let  $DI$  be a domain specification,  $s$  be a state, and  $\alpha = [a_1; \dots; a_n]$  be a sequence of actions.

- $\hat{\Phi}(\alpha, s) = s$  if  $n = 0$ ;
- $\hat{\Phi}(\alpha, s) = \Phi(a_n, \hat{\Phi}([a_1; \dots; a_{n-1}], s))$ , otherwise

where  $\Phi(a, \text{fails}) = \text{fails}$ .

An agent can use the transition function to reason about effects of its actions and to planning. An action sequence  $\alpha$  is a *plan* achieving a set of literals  $O$  from a state  $I$  iff  $O$  is true in  $\hat{\Phi}(\alpha, I)$ .

<sup>2</sup> We envision a multiagent environment where agents may have to wait for other agents to finish some actions before they can go on with their course of actions.

*Example 3.* The domain specification  $DI_A$  for  $A$  in Example 1 is defined over  $F_A = \{h\_nail, h\_screw, mirror\_on, h\_ham\}$  and  $A_A = \{hw\_nail, hw\_screw\}$  with the set of laws<sup>3</sup>:

$hw\_nail$ <b>causes</b> $mirror\_on$	$hw\_screw$ <b>causes</b> $mirror\_on$
$hw\_nail$ <b>causes</b> $\neg h\_nail$	$hw\_screw$ <b>causes</b> $\neg h\_screw$
$hw\_nail$ <b>executable</b> $h\_nail, h\_ham$	$hw\_screw$ <b>executable</b> $h\_screw$

In all of the above, the prefix “ $hw$ ” stands for “hang with” and “ $h$ ” stands for “has.”  $\square$

## 2.2 Specifying Cooperative Actions

The specification of the set of cooperative actions of an agent, denoted by  $DC$ , is defined over  $\mathcal{C}$  and  $\mathcal{F}$  and consists of laws of the following form:

$$r \text{ requests } \gamma \text{ from } \mathcal{A}_i \text{ may\_cause } \phi \text{ if } \psi \text{ and} \quad (3)$$

$$p \text{ provides } \gamma \text{ for } \mathcal{A}_i \text{ causes } \phi \text{ if } \psi \quad (4)$$

$r$  and  $p$  are action names in  $\mathcal{C}$ ,  $\gamma$ ,  $\phi$ , and  $\psi$  are sets of literals and  $\gamma \subseteq \phi$ , and  $\mathcal{A}_i$  is a set of agent identifiers in  $\mathcal{AG}$ .  $r$  is called a *request* for  $\gamma$  and  $p$  an *offer* for  $\gamma$ . Since these actions are intended to address other agents, we require that the identifier of the agent having  $r$  and/or  $p$  does not belong to  $\mathcal{A}_i$ . Furthermore, for a request-action, we require that  $\phi \cap \psi \neq \emptyset$  which indicates that an agent will only request for something that he/she does not have.

Intuitively, (3) represents a set of requests that can be made by the agent; if the agent makes the request for  $\gamma$  (which is the action  $r$ ) directed to an agent in  $\mathcal{A}_i$  then  $\phi$  might become true. The condition  $\gamma \subseteq \phi$  guarantees that requested literals ( $\gamma$ ) are true if the request is satisfied ( $\phi$ ). Furthermore, the action can only be executed if  $\psi$  is true. For this reason, we call  $r(\gamma, i)$ ,  $i \in \mathcal{A}_i$ , an instance of a request (3). Similarly, (4) represents the set of offers  $p(\gamma, i)$ ,  $i \in \mathcal{A}_i$ . This offer addresses a request made to the agent by establishing  $\gamma$  (for the requestor). This action is similar to the individual actions in  $A$  of an agent. The main difference is that they also change the worlds of other agents. It can only be executed if  $\psi$  is true and its effects is  $\phi$ .

For simplicity of the presentation, we will assume that each action in  $\mathcal{C}$  occurs in at most one law of the form (3) or (4). We use *cooperative action* to refer to either a request- or an offer-action. When  $\mathcal{A}_i$  is the set of all other agents, we often omit the part ‘**from**  $\mathcal{A}_i$ ’ from (3) and ‘**for**  $\mathcal{A}_i$ ’ from (4).

*Example 4.* In Example 1, it is reasonable for  $A$  to request and/or offer other agents on the literal  $h\_nail$ . An action for requesting for (offering of)  $h\_nail$  for  $A$  can be specified by

$give\_me\_nail$ <b>requests</b> $h\_nail$ <b>from</b> $\{B, C\}$ <b>may\_cause</b> $h\_nail$ <b>if</b> $\neg h\_nail$
$get\_this\_nail$ <b>provides</b> $h\_nail$ <b>for</b> $\{B, C\}$ <b>causes</b> $\neg h\_nail$ <b>if</b> $h\_nail$

where  $give\_me\_nail$  is a request-action and  $get\_this\_nail$  is an offer-action. If the agent  $A$  wants to ask for help, then her set of cooperative actions needs to include the

<sup>3</sup> To simplify the representation, we often write  $l_1, \dots, l_n$  instead of  $\{l_1, \dots, l_n\}$  in describing the domain.

action *give\_me\_nail*. On the other hand, if she wants to help others, then it should include the action *get\_this\_nail*.  $\square$

**Definition 2 (Planning problem with cooperative actions).** A planning problem with cooperative actions<sup>4</sup>  $\mathcal{P}$  is a tuple  $\langle DI, I, O, DC \rangle$  where  $DI$  is a domain specification,  $I$  is a state representing the initial state,  $O$  is a set of literals representing the goal, and  $DC$  is a set of laws of the form (3) and (4).

Given a planning problem  $\mathcal{P} = \langle DI, I, O, DC \rangle$ , we need to specify what is a “plan” achieving  $O$  in the presence of the cooperative actions. Intuitively, we could consider these actions as the actions of the agent and use the notion of a plan mentioned in the previous subsection. This is, however, not enough since an agent, when executes a request, might or might not receive an offer satisfying his/her request. For example, a request for a nail from  $A$  to  $C$  might not result in  $A$  having the nail because  $C$  has already given the nail to  $B$ .

We will therefore extend the transition function  $\Phi$  of the domain specification  $DI$  to consider cooperative actions. We will use  $\Phi_D$  to denote the transition function of  $DI \cup DC$ . By assuming that cooperative actions are different from the individual actions (i.e.,  $A \cap C = \emptyset$ ), it suffices to specify what is the result of the execution of a request/offer-action in a given state.

For simplicity of the presentation, we assume that each individual agent executes only one action at a time. The method presents in this paper can be easily extended to the case where individual agents can execute parallel actions.

Let  $s$  be a state. We say that an instance  $r(\gamma, i)$  of a request-action specified by the law

$$r \text{ requests } \gamma \text{ from } A_i \text{ may\_cause } \phi \text{ if } \psi$$

in  $DC$  is *executable* in  $s$  if  $\psi$  is true in  $s$ . Executing the action  $r(\gamma, i)$  in  $s$  does not guarantee that the agent will obtain  $\phi$  in the resulting state. This is because the agent, whom the request was made to, might not have the capability to establish  $\phi$  for the requestor. We say that the execution of  $r(\gamma, i)$  in  $s$  might or might not succeed. As such, the result of executing  $r(\gamma, i)$  in  $s$  is either  $s$ , representing the case when the request is not satisfied (by the agent whom the request was made to); or  $(s \setminus \bar{\phi}) \cup \phi$ , representing the case when the request is satisfied.

*Remark 1.* Observe that under the assumption that an agent will execute a request-action only when it is necessary (i.e.,  $\bar{\phi} \cap \psi \neq \emptyset$ ), we have that  $s \neq (s \setminus \bar{\phi}) \cup \phi$  for every instance  $r(\gamma, i)$ . This allows us to recognize when a request is satisfied.

An instance  $p(\gamma, i)$  of an offer-action specified by the law

$$p \text{ provides } \gamma \text{ for } A_i \text{ causes } \phi \text{ if } \psi$$

in  $DC$  is *executable* in  $s$  if  $\psi$  is true in  $s$ . The state resulting from executing  $p(\gamma, i)$  in  $s$  is given by  $(s \setminus \bar{\phi}) \cup \phi$ .

<sup>4</sup> For simplicity of presentation, we will use planning problem instead of planning problem with cooperative actions whenever no confusion is possible.

**Definition 3 (Transition function).** The transition function  $\Phi_D$  over  $DI \cup DC$ , a mapping from pairs of actions and states to sets of states, is defined as follows. Let  $s$  be a state.

- For  $a \in \mathbf{A}$ ,  $\Phi_D(a, s) = \{\Phi(a, s)\}$  if  $\Phi(a, s) \neq \text{fails}$ ; otherwise,  $\Phi_D(a, s) = \emptyset$ .
- For an instance of an offer-action  $p(\gamma, i)$ ,  $\Phi_D(p(\gamma, i), s) = \{(s \setminus \bar{\phi}) \cup \phi\}$  if  $p$  is executable in  $s$ ; otherwise,  $\Phi_D(p, s) = \emptyset$ .
- For an instance of a request-action  $r(\gamma, i)$ ,  $\Phi_D(r(\gamma, i), s) = \{s, (s \setminus \bar{\phi}) \cup \phi\}$  if  $r(\gamma, i)$  is executable in  $s$ ; otherwise,  $\Phi_D(r(\gamma, i), s) = \emptyset$ .

*Remark 2.* The definition of  $\Phi_D$  assumes that each cooperative action occurs in only one law of the form (3) or (4). The definition can be extended to remove this restriction by (i) defining a set  $ec_{r(\gamma, i)}(s)$  (resp.  $ec_{p(\gamma, i)}(s)$ ), similar to the definition of the set of effects of an action  $e_a(s)$  and (ii) changing the definition accordingly.

The transition function is extended to reason about plans as follows.

**Definition 4 (Plan with cooperative actions).** Let  $\mathcal{P}$  be a planning problem  $\langle DI, I, O, DC \rangle$ . We define

- A sequence  $s_0, a_0, s_1, \dots, a_{n-1}, s_n$ , where  $s_i$ 's are states and  $a_j$ 's are actions, is a trajectory if  $s_{i+1} \in \Phi_D(a_i, s_i)$  for  $0 \leq i < n$ .
- A trajectory  $s_0, a_0, s_1, \dots, a_{n-1}, s_n$  is a possible plan achieving  $O$  (or a solution of  $\mathcal{P}$ ) if  $s_0 = I$  and  $s_n \models O$ .
- An occurrence of a request  $r(\gamma, i) = a_j$  in a trajectory  $s_0, a_0, s_1, \dots, a_{n-1}, s_n$  is satisfied if  $s_{j+1} \neq s_j$ ; otherwise, the request is said to be unsatisfied.

Notice that the third item in the above definition is sensible due to Remark 1. A trajectory satisfying the goal  $O$  of the planning problem is a solution of  $\mathcal{P}$  if all satisfied requests assumed in the trajectory indeed materialized, i.e., for each satisfied  $r(\gamma, i)$  in the trajectory, the agent  $i$  executes the action  $p(\gamma, j)$  ( $j$  is the identifier of the agent issuing the request). The topic of coordination between agents will be discussed in the next section.

*Example 5.* Let  $\mathcal{P}_A = \langle DI_A, I_A, O_A, DC_A \rangle$  be the planning problem for  $A$  with  $DI_A$  (Example 3),  $I_A = \{\neg h\_nail, \neg h\_screw, \neg h\_ham, \neg mirror\_on\}$  and  $O_A = \{mirror\_on\}$ , and  $DC_A$  is the set of actions *give\_me\_nail* and *get\_this\_nail* whose specifications are given (Example 4) and the two actions

*give\_me\_ham* **requests**  $h\_ham$  **from**  $\{B, C\}$  **may cause**  $h\_ham$  **if**  $\neg h\_ham$ ,  
*get\_this\_ham* **provides**  $h\_ham$  **for**  $\{B, C\}$  **causes**  $\neg h\_ham$  **if**  $h\_ham$ .

We can easily check the following:

- for  $n \leq 2$ , the problem has no possible plan.
- for  $n = 3$ ,  $\mathcal{P}_A$  has a possible plan which is the following trajectory:  
 $s_0^A, give\_me\_nail(h\_nail, C), s_1^A, give\_me\_ham(h\_ham, B), s_2^A, hw\_nail, s_3^A$   
where  $s_0^A = \{\neg h\_nail, \neg h\_ham, \neg h\_screw, \neg mirror\_on\}$ ,  
 $s_1^A = \{h\_nail, \neg h\_ham, \neg h\_screw, \neg mirror\_on\}$ ,  
 $s_2^A = \{h\_nail, h\_ham, \neg h\_screw, \neg mirror\_on\}$ ,  
 $s_3^A = \{\neg h\_nail, h\_ham, \neg h\_screw, mirror\_on\}$ . □

### 3 Planning for Multiagents

In a multiagent environment, each agent needs to know her capabilities. She also needs to know from whom she can ask for some favors or to whom she could offer helps. Furthermore, it is also common that groups of agents need to know about their joint capabilities. It is also possible that agents might talk the same language. This can be summarized as follows.

- Each agent has its own planning problem, which is described in the previous section.
- The agent might or might not share the same world representation. By default, the world representation of the agent is local. For example, the three agents in Example 1 can use the same set of fluents and actions; and  $A$  has  $\neg h\_nail$  in her initial state whereas  $B$  has  $h\_nail$  in hers, yet this is not a contradictory statement about the world since the fluents are local. On the other hand, the two agents in Example 2 share certain features (e.g. the light) and therefore the fluents encoding these features should have the same value in their representations.
- An agent might request another agent to establish certain conditions in her own world. For example,  $A$  might request  $B$  to establish  $h\_nail$  to be true for her.
- An agent might execute some actions that change the local world of another agent. For example,  $B$  can give  $A$  the nail, thus establishing  $h\_nail$  in the world of  $A$ .
- There might be actions that a set of agents should not execute in parallel. For example, two cars— one goes north-south and another east-west— cannot cross an intersection at the same time.
- There might be actions that a set of agents need to execute in parallel. For example, the action of lifting a table by two agents need to be done in parallel.

It turns out that the language developed in the previous section can be extended to represent and reason about plans/actions of agents in a multiagent environment. With the help of the notion of a planning problem with cooperative actions, a multiagent planning problem can be defined as follows.

**Definition 5 (Multiagent planning problem).** A multiagent planning problem  $\mathcal{M}$  is a tuple  $\langle \mathcal{AG}, \{\mathcal{P}_i\}_{i \in \mathcal{AG}}, \mathcal{F}, \mathcal{IC}, \mathcal{C} \rangle$  where

- $\mathcal{AG}$  is a set of agents,
- $\mathcal{P}_i$  is a planning problem with cooperative actions for each agent  $i \in \mathcal{AG}$ ,
- $\mathcal{F}$  is the set of tuples of the form  $(i, j, f_i, f_j)$  where  $i, j \in \mathcal{AG}$  and  $f_i \in \mathbf{F}_i$  and  $f_j \in \mathbf{F}_j$ , and
- $\mathcal{IC}$  and  $\mathcal{C}$  are sets of sets of agent action pairs of the form  $(i, a_i)$  where  $i$  is an agent and  $a_i$  is an action in  $\mathbf{A}_i$ .

Intuitively, each tuple  $(i, j, f_i, f_j)$  indicates that  $f_i$  and  $f_j$  represent the same state variable in the worlds of two agents  $i$  and  $j$  and can be changed by either  $i$  or  $j$ . This means that they should have the same value in every state of  $i$  and  $j$ . A set of agent-action pairs  $\{(i_1, a_{i_1}), \dots, (i_k, a_{i_k})\} \in \mathcal{IC}$  indicates that the agents  $i_1, \dots, i_k$  cannot execute the actions  $a_{i_1}, \dots, a_{i_k}$  at the same time. On the other hand, a set of agent-action pairs  $\{(i_1, a_{i_1}), \dots, (i_k, a_{i_k})\} \in \mathcal{C}$  indicates that the agents  $i_1, \dots, i_k$  must execute the actions  $a_{i_1}, \dots, a_{i_k}$  concurrently for their effects to be materialized. The sets  $\mathcal{F}, \mathcal{IC}$ , and  $\mathcal{C}$  are called constraints of  $\mathcal{M}$ .



*Example 6.* The planning problem in Example 1 can be represented by

$\mathcal{M}_1 = \langle \{A, B, C\}, \{\mathcal{P}_A, \mathcal{P}_B, \mathcal{P}_C\}, \emptyset, \emptyset, \emptyset \rangle$  where

- $A, B,$  and  $C$  are the students from Example 1;
- $\mathcal{P}_A$  is defined as in Example 5;
- $\mathcal{P}_B = \langle DI_B, I_B, O_B, DC_B \rangle$  where  $DI_B$  is defined over  
 $F_B = \{h\_nail, h\_screw, diploma\_on, h\_ham\}$  and  $A_B = \{hw\_nail, hw\_screw\}$   
with the set of laws:  
 $hw\_nail$  **causes**  $diploma\_on$                        $hw\_nail$  **causes**  $\neg h\_nail$   
 $hw\_nail$  **executable**  $h\_ham, h\_nail$                $hw\_screw$  **causes**  $diploma\_on$   
 $hw\_screw$  **causes**  $\neg h\_screw$                        $hw\_screw$  **executable**  $h\_screw$   
 $I_B = \{h\_nail, h\_screw, h\_ham, \neg diploma\_on\}$  and  $O_B = \{diploma\_on\}$ , and  
 $DC_B$  contains cooperative actions similar to that in  $DC_A$  and  $DC_C$  (below).
- $\mathcal{P}_C = \langle DI_C, I_C, O_C, DC_C \rangle$  where  $DI_C$  is defined over

$$F_C = \{h\_nail, h\_screw, painting\_on\}$$

$$A_C = \{hw\_screw\}$$

with the set of laws:  $hw\_screw$  **causes**  $painting\_on$   
 $hw\_screw$  **causes**  $\neg h\_screw$                        $hw\_screw$  **executable**  $h\_screw$   
 $I_C = \{h\_nail, \neg h\_screw, \neg painting\_on\}$ ,  $O_C = \{painting\_on\}$ , and  $DC_C$  contains the following laws:

$give\_me\_screw$  **requests**  $h\_screw$  **from**  $\{A, B\}$  **may\\_cause**  $h\_screw$  **if**  $\neg h\_screw$   
 $get\_this\_screw$  **provides**  $h\_screw$  **for**  $\{A, B\}$  **causes**  $\neg h\_screw$  **if**  $h\_screw$      $\square$

We now define the notion of a solution for a planning problem.

**Definition 6 (Joint plan for multiagents).** Let  $\mathcal{M} = \langle \mathcal{AG}, \{\mathcal{P}_i\}_{i \in \mathcal{AG}}, \mathcal{F}, \mathcal{IC}, \mathcal{C} \rangle$  be a multiagent planning problem. For each  $i \in \mathcal{AG}$ , let  $S_i = [s_0^i a_0^i, \dots, a_{n-1}^i s_n^i]$  be a possible plan of  $\mathcal{P}_i$ . We say that  $\{S_i\}_{i \in \mathcal{AG}}$  is a joint plan (or solution) of length  $n$  for  $\mathcal{M}$  if for every  $0 \leq k \leq n$ :

- for each instance of a request  $a_k^i = r(\gamma, j)$  that is satisfied in  $S_i$ , we have that  $a_k^j = p(\gamma, i)$ ;
- for each  $(i, j, f_i, f_j) \in \mathcal{F}$ ,  $f_i \in s_k^i$  iff  $f_j \in s_k^j$ ;
- for each  $S \in \mathcal{IC}$ , there exists some  $(i, a) \in S$  such that  $a_k^i \neq a$ ; and
- for each  $S \in \mathcal{C}$ , either  $\{a \mid (i, a) \in S \text{ and } a = a_k^i\} = \{a \mid (i, a) \in S\}$  or  $\{a \mid (i, a) \in S \text{ and } a = a_k^i\} = \emptyset$ .

Intuitively, a joint plan is composed of individual plans which allow the agents to achieve their own goals and satisfy the various constraints of the problem. In the process, agents can help each other in establishing certain conditions. However, if a request of an agent is assumed (by the requestor) to be satisfied within a joint plan then the joint plan must also contain an agent who actually executes an offer action satisfying the request (first item). The second item states that the individual plans must agree with each other on their effects of shared fluents, i.e., it enforces the constraints in  $\mathcal{F}$ . The third and fourth items make sure that non-parallel and parallel constraints in  $\mathcal{IC}$  and  $\mathcal{C}$  are maintained by the joint plan.

*Example 7.* For the multiagent planning problem  $\mathcal{M}_1$  from Example 6, We can easily check the following:

- for  $n \leq 2$ ,  $\mathcal{M}_1$  has no solution.
- for  $n = 3$ , it has a solution consisting of the following plans
  - $S_A = [s_0^A, \text{give\_me\_nail}(h\_nail, C), s_1^A, \text{give\_me\_ham}(h\_ham, B), s_2^A, \text{hw\_nail}, s_3^A, \text{wait}, s_4^A]$
  - $S_B = [s_0^B, \text{hw\_nail}, s_1^B, \text{get\_this\_ham}(h\_ham, A), s_2^B, \text{get\_this\_screw}(h\_screw, C), s_3^B, \text{wait}, s_4^B, ]$
  - $S_C = [s_0^C, \text{get\_this\_nail}(h\_nail, A), s_1^C, \text{wait}, s_2^C, \text{give\_me\_screw}(h\_screw, B), s_3^C, \text{hw\_screw}, s_4^C]$

where all requests are satisfied and the states are uniquely determined by the initial states and the executed actions.  $\square$

The joint plan for the agents in Example 7 requires that each agent executes some cooperative actions. It is easy to see that any joint plan for the two agents in the problem  $\mathcal{M}_2$  requires that only one agent to flip the switch next to her and other agent to wait.

## 4 Computing Joint Plans

In this section, we will present different approaches to computing joint plans. Our approaches utilize answer set programming [18, 19], a declarative programming paradigm that has recently emerged from the study of logic programming under the answer set semantics [11].

### 4.1 Answer Set Semantics of Logic Programs

A logic program  $\Pi$  is a set of rules of the form

$$a_0 \leftarrow a_1, \dots, a_m, \text{not } a_{m+1}, \dots, \text{not } a_n \quad (5)$$

where  $0 \leq m \leq n$ , each  $a_i$  is an atom of a propositional language<sup>5</sup> and *not* represents *negation-as-failure*. A negation as failure literal (or naf-literal) is of the form *not a* where  $a$  is an atom. For a rule of the form (5), the left (right) hand sides of the rule are called the *head* (*body*), respectively. The head and the body can be empty (but not at the same time). A rule is a *constraint* if its head is empty; it is a *fact* if its body is empty.

Consider a set of ground atoms  $X$ . The body of a rule of the form (5) is *satisfied* by  $X$  if  $\{a_{m+1}, \dots, a_n\} \cap X = \emptyset$  and  $\{a_1, \dots, a_m\} \subseteq X$ . A rule of the form (5) with nonempty head is satisfied by  $X$  if either its body is not satisfied by  $X$  or  $a_0 \in X$ . In other words,  $X$  satisfies a rule of the form (5) if its head belongs to  $X$  whenever  $X$  satisfies its body. A constraint is *satisfied* by  $X$  if its body is not satisfied by  $X$ .

For a set of ground atoms  $S$  and a program  $\Pi$ , the *reduct* of  $\Pi$  w.r.t.  $S$ , denoted by  $\Pi^S$ , is the program obtained from the set of all ground instances of  $\Pi$  by deleting

1. each rule that has a naf-literal *not a* in its body with  $a \in S$ , and
2. all naf-literals in the bodies of the remaining rules.

$S$  is an *answer set* of  $\Pi$  if it satisfies the following conditions.

<sup>5</sup> Rules with variables are viewed as a shorthand for the set of its ground instances.

1. If  $\Pi$  does not contain any naf-literal (i.e.  $m = n$  in every rule of  $\Pi$ ) then  $S$  is the smallest set of atoms that satisfies all the rules in  $\Pi$ .
2. If the program  $\Pi$  does contain some naf-literal ( $m < n$  in some rule of  $\Pi$ ), then  $S$  is an answer set of  $\Pi$  if  $S$  is the answer set of  $\Pi^S$ . (Note that  $\Pi^S$  does not contain naf-literals, its answer set is defined in the first item.)

A program  $\Pi$  is said to be *consistent* if it has an answer set. Otherwise, it is inconsistent. To make answer set style programming easier, Niemelä et al. [20] introduce a new type of rules, called *cardinality constraint rule* (a special form of the *weight constraint rule*) of the following form:

$$A_0 \leftarrow A_1, \dots, A_m, \text{not } A_{m+1}, \dots, \text{not } A_n$$

where each  $A_i$  is a *choice atom* of the form  $l\{b_1, \dots, b_k\}u$  with  $b_j$  are atoms and  $l$  and  $u$  are two integers,  $l \leq u$ ; and  $A_0$  can be empty. An atom  $l\{b_1, \dots, b_k\}u$  is said to be true wrt. a set of literals  $S$  iff  $l \leq |S \cap \{b_1, \dots, b_k\}| \leq u$ . The satisfaction of a rule wrt. a set of atoms is extended in the usual way. Using rules of this type, one can greatly reduce the number of rules of programs in answer set programming. The semantics of logic programs with such rules is given in [20].

#### 4.2 Finding a Possible Plan for One Agent

We will represent each individual problem of each agent  $\mathcal{P}_i$  by a logic program. The program will consist of rules describing the effects of actions, the initial knowledge of the agent, and the goal of the agent. Answer set planning [16] refers to the use of answer set programming in planning. This method has been applied to a variety of problems [10, 25]. Let  $\mathcal{P} = \langle DI, I, O, DC \rangle$  be a planning problem. We will now describe the program  $\Pi(\mathcal{P})$  that encodes  $\mathcal{P}$ . We adapt the conventional style in logic programming: terms starting with lower-case letter are constant and others are variables. It also has a parameter denoting the maximal length of the plan that the agent considers permissible. The key predicates of  $\Pi(\mathcal{P})$  are:

- $h(l, t)$  – fluent literal  $l$  holds at the time step  $t$ ; and
- $o(a, t)$  – action  $a$  is executed (by the agent) at the time step  $t$ ;
- $poss(a, t)$  – action  $a$  can be executed at the time step  $t$ .

$h(l, t)$  can be extended to define  $h(\phi, t)$  for an arbitrary fluent formula  $\phi$ , which states that  $\phi$  holds at the time moment  $t$ . In writing the program, we use  $h(\{l_1, \dots, l_k\}, T)$  as a shorthand for  $h(l_1, T), \dots, h(l_k, T)$ . In addition, we write  $h(ok(r(\gamma, i)), t)$  to denote that the request-action  $r(\gamma, i)$  is satisfied at the time step  $t$ . The rules of the program is divided into groups:

- *Group 1:* The program contains the following facts:

$$\begin{aligned} & \{fluent(f) \mid f \in \mathbf{F}\} \cup \{action(a) \mid a \in \mathbf{A}\} \cup \\ & \{action(r(\gamma), i) \mid r \text{ occurring in a law of form (3), } i \in \mathcal{A}_i\} \cup \\ & \{action(p(\gamma), i) \mid p \text{ occurring in a law of form (4), } i \in \mathcal{A}_i\} \end{aligned}$$

These facts declare the fluents and the actions of the problem.

- *Group 2:* rules for reasoning about effects of actions. For each action  $a \in \mathbf{A}$ ,
  - if  $DI$  contains the law ( $a$  **executable**  $\phi$ ) then  $\Pi(\mathcal{P})$  contains the rules

$$poss(a, T) \leftarrow h(\phi, T) \tag{6}$$

$$\leftarrow o(a, T), \text{not } poss(a, T) \tag{7}$$

- if  $DI$  contains the law ( $a$  **causes**  $l$  **if**  $\phi$ ) then  $\Pi(\mathcal{P})$  contains the rule

$$h(l, T + 1) \leftarrow o(a, T), h(\phi, T) \quad (8)$$

- *Group 3*: rules for reasoning about request-actions. For each statement of the form

$r$  **requests**  $\gamma$  **from**  $\mathcal{A}_i$  **may\_cause**  $\phi$  **if**  $\psi$

and each  $i \in \mathcal{A}_i$ ,  $\Pi(\mathcal{P})$  contains the rules

$$poss(r(\gamma, i), T) \leftarrow h(\psi, T) \quad (9)$$

$$\leftarrow o(r(\gamma, i), T), not\ poss(r(\gamma, i), T) \quad (10)$$

$$0 \{h(ok(r(\gamma, i)), T + 1)\} 1 \leftarrow o(r(\gamma, i), T). \quad (11)$$

$$h(\phi, T) \leftarrow h(ok(r(\gamma, i)), T) \quad (12)$$

where (12) is a shorthand for the collection of rules  $\{h(l, T) \leftarrow h(ok(r(\gamma, i)), T) \mid l \in \phi\}$ . Observe that atoms of the form  $h(ok(\gamma, i), T)$  are used to record the satisfaction of the request  $r(\gamma, i)$  and there might be different ways for a condition  $\gamma$  to be satisfied. Hence, (11) and (12) need to be separated even though it looks like they could have been merged into one.

- *Group 4*: rules for reasoning about offer-actions. For each statement of the form

$p$  **provides**  $\gamma$  **for**  $\mathcal{A}_i$  **causes**  $\phi$  **if**  $\psi$

and  $i \in \mathcal{A}_i$ ,  $\Pi(\mathcal{P})$  contains the rules

$$poss(p(\gamma, i), T) \leftarrow h(\psi, T) \quad (13)$$

$$\leftarrow o(p(\gamma, i), T), not\ poss(p(\gamma, i), T) \quad (14)$$

$$h(\phi, T + 1) \leftarrow o(p(\gamma, i), T). \quad (15)$$

These rules are similar to the rules encoding the effect of individual actions of the agent. The difference between the encoding of a request-action and the encoding of an offer-action lies in that we do not need to introduce an atom of the form  $h(ok(p(\gamma, i)), T)$  to record the execution of  $p(\gamma, i)$ , i.e., effects of offer-actions are deterministic.

- *Group 5*: rules describing the initial state. For each literal  $l \in I$ ,  $\Pi(\mathcal{P})$  contains the fact  $h(l, 0)$ .
- *Group 6*: rules encoding the goal state. For each literal  $l \in O$ ,  $\Pi(\mathcal{P})$  contains

$$\leftarrow not\ h(l, n). \quad (16)$$

where  $n$  is the desired length of the plan.

- *Group 7*: rules for reasoning by inertial. For each fluent  $F \in \mathbf{F}$ ,  $\Pi(\mathcal{P})$  contains

$$h(F, T + 1) \leftarrow h(F, T), not\ h(\neg F, T + 1). \quad (17)$$

$$h(\neg F, T + 1) \leftarrow h(\neg F, T), not\ h(F, T + 1). \quad (18)$$

$$\leftarrow h(F, T), h(\neg F, T). \quad (19)$$

- *Group 8*: rules for generating action occurrences.  $\Pi(\mathcal{P})$  contains the rule

$$1 \{o(A, T) : action(A)\} 1 \leftarrow T < n. \quad (20)$$

which states that at any time step, the action must execute one of its actions<sup>6</sup>.

Let  $\mathcal{P} = \langle DI, I, O, DC \rangle$  be a planning problem and  $\Pi(\mathcal{P}, n)$  denote the set of ground rules of  $\Pi(\mathcal{P})$  in which the variable  $T$  is instantiated with integers between 0 to  $n$ . Let  $M$  be an answer set of  $\Pi(\mathcal{P}, n)$ . Let  $s_t[M] = \{l \mid l \text{ is a fluent literal and } h(l, t) \in M\}$ . By  $\alpha[M]$  we denote the sequence  $s_0[M], a_0, s_1[M], \dots, a_{n-1}, s_n[M]$  where  $o(a_i, i) \in M$ . We can show the following:

**Theorem 1.** *Let  $\mathcal{P}$  be a planning problem. Then,*

- *for each possible plan  $\alpha$  of  $\mathcal{P}$  there exists an  $n$  and an answer set  $M$  of  $\Pi(\mathcal{P}, n)$  such that  $\alpha = \alpha[M]$ ;*
- *for each  $n$ , if  $\Pi(\mathcal{P}, n)$  is inconsistent then  $\mathcal{P}$  does not have a solution of length less than or equal to  $n$ ; and*
- *for each  $n$ , if  $\Pi(\mathcal{P}, n)$  has an answer set  $M$  then  $\alpha[M]$  is a solution of  $\mathcal{P}$ .*

### 4.3 Compatible Answer Sets and Joint Plan

Individual possible plans can be computed using the program  $\Pi(\mathcal{P}_i)$ . We will now discuss an approach for combining them to create a plan for all the agents. Intuitively, we need to make sure that if a request is assumed to be satisfied by an agent then there exists an instance of an offer-action matching this request. This can be easily characterized by the notion of a compatible answer sets.

**Definition 7 (Compatible answer sets).** *Let  $\mathcal{M} = \langle \mathcal{AG}, \{\mathcal{P}_i\}_{i \in \mathcal{AG}}, \mathcal{F}, \mathcal{IC}, \mathcal{C} \rangle$  be a multiagent planning problem and  $M = \langle M_i \rangle_{i \in \mathcal{AG}}$  be a sequence of answer sets of  $\langle \Pi(\mathcal{P}_i, n) \rangle_{i \in \mathcal{AG}}$  where the constant  $n$  is fixed.  $M$  is a set of compatible answer sets if for each  $k \leq n$ ,*

- *for each  $i \in \mathcal{AG}$ , if  $h(ok(r(\gamma, j)), k + 1) \in M_i$  then  $o(p(\gamma, i), k) \in M_j$ ;*
- *for each  $i \in \mathcal{AG}$ , if  $o(p(\gamma, j), k) \in M_i$  then  $h(ok(r(\gamma, i)), k + 1) \in M_j$ ;*
- *for each  $(i, j, f_i, f_j)$  in  $\mathcal{F}$ ,  $h(f_i, k) \in M_i$  iff  $h(f_j, k) \in M_j$ ;*
- *for each  $S \in \mathcal{IC}$  there exists some  $(i, a_i) \in S$  such that  $o(a_i, k) \notin M_i$ ; and*
- *for each  $S \in \mathcal{C}$ , either  $\{a_i \mid (i, a_i) \in S \text{ and } o(a_i, k) \in M_i\} = \{a \mid (i, a) \in S\}$  or  $\{a_i \mid (i, a_i) \in S \text{ and } o(a_i, k) \in M_i\} = \emptyset$ .*

Intuitively, a set of compatible answer sets corresponds to a joint plan (as we will prove in the next theorem) similar to the correspondence between answer sets and plans in the case of a single agent. Observe also that  $h(ok(\cdot), T)$  is present only due to the successfulness of a request-action, not an offer-action. The conditions imposed on a set of compatible answer sets make sure that the collection of individual plans extracting from them satisfies the constraints of the planning problem and the requirement that satisfying requests must be matched with offers.

<sup>6</sup> Since we assume that `wait` always belongs to the set of actions of an agent, this is not a strict requirement as it might sound.

**Theorem 2.** Let  $\mathcal{M} = \langle \mathcal{AG}, \{\mathcal{P}_i\}_{i \in \mathcal{AG}}, \mathcal{F}, \mathcal{IC} \rangle$  be a multiagent planning problem and  $n$  be an integer:

- if  $\langle \Pi(\mathcal{P}_i, n) \rangle_{i \in \mathcal{AG}}$  does not have a set of compatible answer sets then  $\mathcal{M}$  does not have a solution with length  $n$ .
- a sequence of answer sets  $M = \langle M_i \rangle_{i \in \mathcal{AG}}$  is compatible iff there exists a solution  $S = \langle \alpha_i \rangle_{i \in \mathcal{AG}}$  such that  $\alpha[M_i] = \alpha_i$  for every  $i \in \mathcal{AG}$ .

*Example 8.* Let  $\mathcal{M}_1$  be the multiagent planning problem from Example 6. We can easily check the following:

- $\{\Pi(\mathcal{P}_i, n)\}_{i \in \{A, B, C\}}$  for  $n \leq 2$  does not have compatible answer sets,
- for  $n = 3$ , the three answer sets  $M_A, M_B$ , and  $M_C$  of  $\Pi(\mathcal{P}_A, 3), \Pi(\mathcal{P}_B, 3)$ , and  $\Pi(\mathcal{P}_C, 3)$ , where
  - $M_A$  contains  $o(\text{give\_me\_nail}(h\_nail, c), 0), h(\text{ok}(\text{give\_me\_nail}(h\_nail, c)), 1), o(\text{give\_me\_ham}(h\_ham, b), 1), h(\text{ok}(\text{give\_me\_ham}(h\_ham, b))), 2), o(\text{hw\_nail}, 2)$ , and  $o(\text{wait}, 3)$ .
  - $M_B$  contains  $o(\text{hw\_nail}, 0), o(\text{get\_this\_ham}(h\_ham, a), 1), o(\text{get\_this\_screw}(h\_screw, c), 2), o(\text{wait}, 3)$ ; and
  - $M_C$  contains  $o(\text{get\_this\_nail}(h\_nail, a), 0), o(\text{wait}, 1), o(\text{give\_me\_screw}(h\_screw, b), 2), h(\text{ok}(\text{give\_me\_screw}(h\_screw, b))), 2)$ , and  $o(\text{hw\_screw}, 3)$ .

These answer sets are compatible and correspond to the solution in Example 5.  $\square$

The notion of joint plan can be specialized as follows.

**Definition 8 (Optimal Joint Plan).** Let  $\mathcal{M} = \langle \mathcal{AG}, \{\mathcal{P}_i\}_{i \in \mathcal{AG}}, \mathcal{F}, \mathcal{IC}, \mathcal{C} \rangle$  be a multiagent planning problem and  $\{S_i\}_{i \in \mathcal{AG}}$  be a plan for  $\mathcal{M}$ . We say that  $\{S_i\}_{i \in \mathcal{AG}}$  is optimal if there exists no unsatisfied request actions in  $\{S_i\}_{i \in \mathcal{AG}}$ .

*Remark 3.* The program  $\Pi(\mathcal{P}_i)$  can be easily adapted to generate only optimal plans. Indeed, the only modification that needs to be done is to replace the rule (11) with

$$h(\text{ok}(r(\gamma, i)), T + 1) \leftarrow o(r(\gamma, i), T).$$

Intuitively, this rule states that the request  $r(\gamma, i)$  is satisfied. Thus, if a joint plan is found it will not contain any unsatisfied requests, i.e., it must be optimal.

Definitions 6 and 7 provide us with a way for computing joint plans of length  $n$  for a planning problem  $\mathcal{M}$ . The process involves (i) computing a set  $\{M_i\}_{i \in \mathcal{AG}}$  of answer sets, where  $M_i$  is an answer set of  $\Pi(\mathcal{P}_i, n)$ ; and (ii) checking the compatibility of  $\{M_i\}_{i \in \mathcal{AG}}$ . In what follows, we discuss a method for doing it. This method computes a joint plan by (a) forming a program representing  $\mathcal{M}$  from the programs representing the individual plans and the set of constraints in  $\mathcal{M}$ ; and (b) extracting joint plan from answer sets of the new program. This method is useful if the planning problem  $\mathcal{M}$  is known to an agent or a manager.

#### 4.4 Computing Joint Plans by Answer Set Programming

Let  $\mathcal{M} = \langle \mathcal{AG}, \{\mathcal{P}_i\}_{i \in \mathcal{AG}}, \mathcal{F}, \mathcal{IC}, \mathcal{C} \rangle$  be a planning problem. We will define a program  $\Pi(\mathcal{M})$  whose answer sets represent the solutions of  $\mathcal{M}$ .  $\mathcal{M}$  is constructed from the

programs  $\Pi(\mathcal{P}_i)$  for  $i \in \mathcal{AG}$  as follows. For each  $i \in \mathcal{AG}$ , let  $\Pi^i(\mathcal{P}_i)$ , referred as the *tagged version* of  $\Pi(\mathcal{P}_i)$ , be the program obtained from  $\Pi(\mathcal{P}_i)$  by replacing every literal  $x$  in  $\Pi(\mathcal{P}_i)$  with the atom  $x^i$  (e.g.,  $action(a)^i$  for  $action(a)$ ,  $h(f, t)^i$  for  $h(f, t)$ , etc.). The program  $\Pi(\mathcal{M})$  consists of

- for each  $i \in \mathcal{AG}$ , the tagged version  $\Pi^i(\mathcal{P}_i)$  of  $\Pi(\mathcal{P}_i)$ ;
- for each tuple  $(i, j, f^i, f^j)$  in  $\mathcal{F}$ , the constraints

$$\leftarrow h^i(f^i, T), h^j(\neg f^j, T) \quad (21)$$

$$\leftarrow h^i(\neg f^i, T), h^j(f^j, T) \quad (22)$$

ensure that shared variables maintain their consistency.

- for each set  $S = \{(i_1, a_1), \dots, (i_k, a_k)\}$  in  $\mathcal{C}$ , the constraint

$$\leftarrow 0 \{o^{i_1}(a_1, T), \dots, o^{i_k}(a_k, T)\} k - 1 \quad (23)$$

which makes sure that if a part of  $S$  is executed, i.e.,  $o(i_j, a_j)$  belongs to an answer set, then the whole set  $S$  is executed.

- for each set  $\{(i_1, a_1), \dots, (i_k, a_k)\}$  in  $\mathcal{IC}$ , the constraints

$$\leftarrow o^{i_1}(a_1, T), \dots, o^{i_k}(a_k, T) \quad (24)$$

This constraint guarantees that not all the actions  $a_1, \dots, a_k$  are executed at the same time.

- for every pair of instance  $r(\gamma, j)$  and  $p(\gamma, i)$  of a request-action  $r$  (for  $\gamma$ ) of an agent  $i$  and an offer-action  $p$  (for  $\gamma$ ) of an agent  $j$ , the following constraints

$$\leftarrow o^i(r(\gamma, j), T), h^i(ok(r(\gamma, j)), T + 1), not\ o^j(p(\gamma, i), T) \quad (25)$$

$$\leftarrow o^j(p(\gamma, i), T), not\ o^i(r(\gamma, j), T) \quad (26)$$

$$\leftarrow o^j(p(\gamma, i), T), not\ h^i(ok(r(\gamma, j)), T + 1) \quad (27)$$

The first constraint makes sure that if  $i$  requests for  $\gamma$  from  $j$  and it is satisfied then  $j$  does indeed offer the service. The last two rules guarantee the converse.

For a set  $X$  of literals in the language of  $\Pi(\mathcal{M})$ , let  $X|_i = \{a \mid a \text{ is a literal in the language of } \Pi(\mathcal{P}_i) \text{ and } a^i \in X\}$ . We have:

**Theorem 3.** *Let  $\mathcal{M} = \langle \mathcal{AG}, \{\mathcal{P}_i\}_{i \in \mathcal{AG}}, \mathcal{F}, \mathcal{IC}, \mathcal{C} \rangle$  be a multiagent planning problem.  $M$  is an answer set of  $\Pi(\mathcal{M}, n)$  iff there exists a set of compatible answer sets  $\{M_i\}_{i \in \mathcal{AG}}$  such that  $M|_i = M_i$ .*

The proof of Theorem 3 relies on the Splitting Theorem for logic programs [17]. It is divided into two steps. First, it is proved for program without the constraints (21)-(27). The significance of this proposition is that it allows us to compute the solution of a multiagent planning problem by computing a single answer set of  $\mathcal{P}(\mathcal{M})$ . Since the problem of determining whether a propositional program has an answer set or not is NP-complete, the following holds.

**Corollary 1.** *Determining whether a solution of polynomial bounded length of a multiagent planning problem  $\mathcal{M}$  exists or not is NP-complete.*

## 5 Related Works

Multiagent planning could be viewed as a special case of distributed problem solving [9]. In this respect, our work could be viewed as one in the Centralized Planning for Distributed Plans group according to the classification in [9]. This is achieved by the program  $\Pi(\mathcal{M})$ . Alternatively, the individual plans can also be computed distributedly and coordinated using the program consisting of the constraints (21)-(27) and the tagged versions of the individual answer sets.

Our main goal is to generate a joint plan for the agents before its execution. In this regards, our work differs from many distributed continual planning systems that were discussed in the survey [7] and many papers presented in the recent AAMAS conferences which concentrate on planning and replanning or dealing with unexpected events during the plan execution.

Our approach to generating a joint plan in this paper blends the two components “planning” and “coordination” in the equation

$$\boxed{\text{Multiagent planning} = \text{Planning} + \text{Coordination}}$$

presented in [6] into a single step. Furthermore, we employ a plan representation that allows for the coordination to be done by using time-steps presented in individual plans. This is different from several other systems in which partial order plans are used for plan representation and refinement planning is used for coordination (e.g., [4, 3] or earlier works such as the Patial Global Planning framework).

We use answer set programming [16], a method that has been used for single agent planning [10, 25], in computing the joint plan. The declarativeness and modularity of answer set programming make the process of computing the joint plan fairly simple and simplify the coordination of the plans<sup>7</sup>. Our work is similar to the spirit of that in [8] where an attempt is made to construct joint plan using SAT-based single agent planner. Nevertheless, our use of answer set programming does not require the development of additional algorithms to assemble the final joint plan.

In [2], a language based on PDDL for modeling multiagent planning problems has been proposed that allows for the specification of and reasoning about several features important for multiagent planning and execution such as concurrency, individual and mutual beliefs of agents, planning and acting with incomplete information, communication, continuous events, etc. A special agent, called `env`, is present in all problems for modeling the environment which may act “unpredictably”. Our language is less expressive than the above mentioned language as our focus is solely on the generation of a joint plan prior to its execution. On the other hand, the semantics provided in this paper can be used to prove formal properties of plans as well as the correctness of the logic program encoding of multiagent planning problem.

We note that collaborative actions presented in this paper is also suitable for the modeling of multiagent planning with resources. Requesting for a resource and offering a resource can be modeled in a similar fashion to that of asking for and offering of a nail (Example 4). Since our focus is the generation of joint plan before execution, the proposed language is different from the resource logic introduced in [5], whose focus

---

<sup>7</sup> Recall that this is achieved by simply adding the rules (21)-(27).



was on the plan merging phase. The requesting/offering actions can be seen as special case of *negotiation actions* discussed in [26].

We would like to point out that we use  $\mathcal{A}$  because of its simple semantics and its close relationship to PDDL, the language developed for describing planning problems [14]. This means that other extensions or variations of  $\mathcal{A}$  (e.g.,  $\mathcal{B}$ ,  $\mathcal{C}$  [13],  $\mathcal{E}$  [15]) could also be extended to formalize cooperative actions. Observe that there are subtle differences between request actions and non-deterministic actions. First, a cooperative action changes the world of other agents while a non-deterministic action does not. Second, a cooperative action does not change the world of the agent executing this action, while a non-deterministic action does. In this sense, a cooperative action of an agent is like an exogenous action for other agents. Thus, modelling cooperative actions using non-deterministic actions might not be the most natural way.

Finally, we would like to note that an extension of the STRIPS language has been considered for multiagent planning in [1]. In this framework, a multiagent planning problem is formulated as a *single problem* and agent identifiers are attached to the actions, which is different from what we proposed here. As such, the framework in [1] is only appropriate for domains where no privacy among agents is required. This is not an issue in our formulation.

## 6 Conclusions and Future Works

We extend the action language  $\mathcal{A}$  to define a language for representing and reasoning about actions and their effects in presence of cooperative actions between agents. We define the notion of a plan with cooperative actions and use it in formalizing the notion of a joint plan. We use answer set programming to generate joint plans. We introduce the notion of a set of compatible answer sets and provide a translation of a multiagent planning problem to a logic program whose answer sets represent joint plans.

The work so far has focused on the development of a theoretical framework for generating joint plans using answer set programming. The encoding of the examples are available in the technical report version of this paper [24]. Our immediate goal for the future is to investigate the scalability and efficiency of the proposed method. The use of answer set programming allows us to easily incorporate preferences or domain knowledge in the generation of the joint plans [22, 23]. Additionally, we would like to explore the use of more expressive languages (e.g., action languages with constraints and sensing actions) in representing and reasoning about joint-plans of multiagents by addressing various questions mentioned in [2]. This is because the method provides in Section 2 has proved to be very effective in the single agent case (e.g. [25]).

## References

1. C. Boutilier and R.I. Brafman. Partial-order planning with concurrent interacting actions. *JAIR*, 14:105–136, 2001.
2. M. Brenner. Planning for Multiagent Environments: From Individual Perceptions to Coordinated Execution. In *Work. on Multiagent Planning & Scheduling, ICAPS*, 80–88. 2005.
3. J. S. Cox and E. H. Durfee. An efficient algorithm for multiagent plan coordination. *AAMAS 2005*, 828–835.
4. J. S. Cox and E. H. Durfee and T. Bartold. A Distributed Framework for Solving the Multiagent Plan Coordination Problem. In *AAMAS*, pages 821–827. ACM Press, 2005.

5. M. de Weerd, A. Bos, H. Tonino, and C. Witteveen. A resource logic for multi-agent plan merging. *Ann. Math. Artif. Intell.*, 37(1-2):93–130, 2003.
6. M. de Weerd, A. ter Mors, and C. Witteveen. Multi-agent planning: An introduction to planning and coordination. In *Handouts of the Euro. Agent Summer School*, 1–32, 2005.
7. M. desJardins, E. H. Durfee, C. L. Ortiz, and M. Wolverton. A survey of research in distributed, continual planning. *AI Magazine*, 20(4):13–22, 1999.
8. Y. Dimopoulos and P. Moraitis. Multi-agent coordination and cooperation through classical planning. In *IEEE/WIC/ACM/IAT*, 398–402. IEEE Comp. Society, 2006.
9. E. Durfee. Distributed Problem Solving and Planning. In *Muliagent Systems (A Modern Approach to Distributed Artificial Intelligence)*, pages 121–164. MIT Press, 1999.
10. T. Eiter, W. Faber, N. Leone, G. Pfeifer, and A. Polleres. Answer Set Planning under Action Costs. *Journal of Artificial Intelligence Research*, 19:25–71, 2003.
11. M. Gelfond and V. Lifschitz. The stable model semantics for logic programming. In *Int. Conf. on Logic Programming*, 1070–1080, 1988.
12. M. Gelfond and V. Lifschitz. Representing actions and change by logic programs. *Journal of Logic Programming*, 17(2,3,4):301–323, 1993.
13. M. Gelfond and V. Lifschitz. Action languages. *ETAI*, 3(6), 1998.
14. M. Ghallab, A. Howe, C. Knoblock, D. McDermott, A. Ram, M. Veloso, D. Weld, and D. Wilkins. PDDL — the Planning Domain Definition Language. Ver. 1.2. TR1165, Yale, 1998.
15. A. C. Kakas, R. Miller, F. Toni: E-RES. Reasoning about Actions, Events and Observations. In *LPNMR*, 254–266, 2001.
16. V. Lifschitz. Action languages, answer sets and planning. In *The Logic Programming Paradigm: a 25-Year Perspective*, pages 357–373. Springer Verlag, 1999.
17. V. Lifschitz and H. Turner. Splitting a logic program. In *ICLP*, 23–38, 1994.
18. V. Marek and M. Truszczyński. Stable models and an alternative logic programming paradigm. In *The Log. Prog. Paradigm: a 25-year Perspective*, 375–398, 1999.
19. I. Niemelä. Logic programming with stable model semantics as a constraint programming paradigm. *AMAI*, 25(3,4):241–273, 1999.
20. I. Niemelä, P. Simons, and T. Soinen. Stable model semantics for weight constraint rules. In *Proc. Logic Programming and NonMonotonic Reasoning*, 315–332, 1999.
21. S. Parsons, C. Sierra, and N. R. Jennings. Agents that reason and negotiate by arguing. *J. of Log. and Comp.*, 8(3):261–292, 1998.
22. T. C. Son, C. Baral, N. Tran, and S. McIlraith. Domain-Dependent Knowledge in Answer Set Planning. *ACM Transactions on Computational Logic*, 7(4), 2006.
23. T. C. Son and E. Pontelli. Planning with Preferences using Logic Programming. *Journal of Theory and Practice of Logic Programming (TPLP)*, 6:559–607, 2006.
24. T. C. Son and C. Sakama. Reasoning and Planning with Cooperative Actions for Multiagents Using Answer Set Programming. Technical Report, 2008.
25. P. H. Tu, T. C. Son, C. Baral. Reasoning and Planning with Sensing Actions, Incomplete Information, and Static Causal Laws using Logic Programming. *TPLP*, 7:1–74, 2006.
26. M. Wooldridge and S. Parsons. Languages for negotiation. In *Proceedings of ECAI*, 2000.

# Social commitments in time: satisfied or compensated (preliminary report)

Paolo Torroni, Federico Chesani, Paola Mello, and Marco Montali

DEIS, University of Bologna. V.le Risorgimento 2, 40136 Bologna, Italy

**Abstract.** We define a framework based on computational logic technology and on a reactive axiomatization of the Event Calculus to formalize the evolution of commitments in time. We propose a new characterization of commitments with time that enables a rich modeling of the domain, various forms of reasoning, and run-time and static verification.

## 1 Introduction

Social commitments are commitments made from an agent to another agent to bring about a certain property. In broad terms, a social commitment represents the commitment that an agent, called *debtor*, has towards another agent, called *creditor*, to bring about some property or state of affairs, which is the *subject* of the commitment. In some instantiations of this idea, such as [6, 14], the subject of a commitment is a temporal logic formula.

Commitments are a well-known concept in Multi-Agent Systems (MAS) research [2, 12]. Representing the commitments that the agents have to one another and specifying constraints on their interactions in terms of commitments provides a principled basis for agent interactions [13]. From a MAS modelling perspective, a role can be modelled by a set of commitments. For example, a seller in an online market may be understood as committing to its price quotes and a buyer may be understood as committing to paying for goods received. Commitments also serve as a natural tool to resolve design ambiguities. The formal semantics enables verification of conformance and reasoning about the MAS specifications [5] to define core interaction patterns and build on them by reuse, refinement, and composition.

Central to the whole approach is the idea of manipulation of commitments: their creation, discharge, delegation, assignment, cancellation, and release, since commitments are stateful objects that change in time as events occur. Time and events are, therefore, essential elements. Some authors distinguish between *base-level* commitments, written  $C(x, y, p)$ , and *conditional* commitments, written  $CC(x, y, p, q)$  ( $x$  is the debtor,  $y$  is the creditor, and  $p/q$  are properties).  $CC(x, y, p, q)$  signifies that if  $p$  is brought out,  $x$  will be committed towards  $y$  to bring about  $q$ .

In this work we give emphasis to temporal aspects of commitments. We build from previous research by Mallya et al. [10, 9]. In our opinion, they represent the

best articulated research on time-enhanced commitments to date. The main idea in these articles is to extend a commitment framework with a way to describe time points and intervals, and alternative outcomes due to commitments extending into the uncertain future. The perspective on commitment-related temporal issues proposed by [10] mainly aims to capture the idea of validity of a commitment. Thus the previous notation  $C(x, p, y)$  is extended with existential and universal *temporal quantifiers*, which become prefixes of  $p$ . There are two types of quantification. By an existential quantification,  $[t_1, t_2]p$ , we mean that  $p$  is true at one or more moments in the interval beginning at  $t_1$  and ending at  $t_2$ . By a universal quantification,  $\overline{[t_1, t_2]}p$ , we indicate instead that  $p$  is true at every moment in the interval beginning at  $t_1$  and ending at  $t_2$ . Nestings are possible. For example, a commitment from  $x$  to  $y$  that  $q$  is going to hold at every moment in a week beginning on some day between the 1st and 20th of February could be written as follows:  $C(x, y, [01.02.2009, 20.02.2009]([t_{start}, t_{start} + 7days]q))$ .

This is an elegant approach which decouples the temporal quantification from the proposition, enabling reasoning about the temporal aspect without regard to the propositions' meaning. However, there are still some cases in which such a characterization is difficult to use in practical applications. The main problems are due to the lack of variables in temporal logic expressions, and from the separation between such expressions and the other parts of the represented knowledge. Among the aims of this work there is our intention to identify such cases and discuss them.

Along with a notation to express commitments, we need a language to express operations on commitments. For example, Yolum and Singh propose a notation based on the Event Calculus temporal representation language to describe commitment manipulation inside an operational framework [14]. Moreover, from a design perspective, we need an architecture in which a commitment notation, a temporal representation language and a specification and verification framework are given a specific role.

In this paper, we discuss our ongoing research about commitment frameworks. We start by introducing some issues regarding social commitment modeling, and define a number of desiderata for social commitment frameworks. We then propose a new notation for commitments and commitment specification programs: the Commitment Modeling Language ( $\mathcal{CML}$ ). Finally, we outline an abstract commitment framework architecture and a concrete instance of it that supports  $\mathcal{CML}$ . In such an instance, temporal reasoning with commitments is operationalized using a reactive implementation of the Event Calculus and various verification tasks can be accomplished thanks to an underlying declarative, computational logic-based framework.

## 2 Some issues regarding modeling

The following informal discussion is example-driven. Examples are mainly taken from the literature. We start by observing that in some cases Mallya et al.'s notation can be simplified, considering that to represent an existentially quantified

time interval it is sufficient to represent a time point using a variable with a domain. We then sketch a new possible notation that accommodates variables with domains and temporal modeling in several dimensions. Again, based on literature examples, we demonstrate the possible usage of rules to model conditional commitments. Finally we discuss time associated with commitment state changes and the issue of compensation.

## 2.1 Time variables, rules and constraints

Let us analyze the scenario proposed by Mallya et al. in [10].

*Example 1.* A travel agent wishes to book an airline ticket to a certain destination, a rental car to use while there, and a hotel room at which to stay. Consider four situations:

- *Situation 1.1.* The travel agent wants the passenger to fly on a particular day while still reserving the right to choose any flight on that day. If the airline offers such a deal, it becomes committed to maintaining a condition—a booked ticket—over an extended time period.
- *Situation 1.2.* The car rental company offers a one-week free rental in January.
- *Situation 1.3.* A hotel offers an electronic discount coupon that expires today, but text on the coupon states that it can only be used during a future spring break. Note that in this case the commitment violates a constraint about time. In fact, the coupon expires before it can be used.
- *Situation 1.4.* The car rental company guarantees that its cars will not break down for at least two days, promising an immediate replacement if one does. However, if the company is closed on weekends, then a customer who rents a car on a Friday would not benefit from the warranty if the car broke down on Saturday. Thus in this case the car rental company offers a warranty that cannot be used during the period in which the warranty is valid.  $\square$

Following [10], we use the symbols  $h$  for hotel,  $g$  for guest,  $r$  for rental company,  $c$  for customer,  $a$  for airline and  $p$  for the proposition, subject of the commitment. How do we model the situations above using commitments?

*Situation 1.1.* Let  $p$  represent that a ticket booking is guaranteed. Thus, using an existential temporal quantifier,  $[t_1, t_1 + 24hrs]p$ , we express that a ticket booking is guaranteed for 1 day, as of  $t_1$  [10].

However, in practical applications, it may be interesting to explicitly model the time at which the commitment is satisfied (e.g., when the ticket is issued). To this end, we could use an alternative notation, which associates  $p$  with a variable, and binds such a variable to a domain interval:  $[T]p, T \in [t_1, t_1 + 24hrs]$ . We write the commitment as follows:

$$C(a, g, [T]p), t_1 \leq T, T \leq t_1 + 24hrs. \quad (1)$$

In this way, the commitment is satisfied if there is a possible value of  $T$  which falls in the range  $[t_1, t_1 + 24hrs]$ , and such a value can be used for further inferences.

*Situation 1.2.* Let  $p$  denote free rental, and  $t_1$  January 1st. Thus, using a universal temporal quantifier, guaranteed free rental for 7 days as of time  $t_3$  is denoted by  $\overline{[t_3, t_3 + 7days]p}$ . Then to express that such an interval  $\overline{[t_3, t_3 + 7days]}$  is inside January, Mallya et al. [10] use a condition on  $t_3$ , namely  $t_1 \leq t_3 \leq t_1 + 24days$ , and they attach an existential temporal quantifier outside of the quantifier above:  $[t_1, t_1 + 31days](\overline{[t_3, t_3 + 7days]p}, t_1 \leq t_3 \leq t_1 + 24days)$ .

Let us now use the notation introduced above, instead of existential temporal quantification. We obtain  $[T, T + 7days]p, t_1 \leq T, T \leq t_1 + 24days$ . Note that we simplified the notation. In particular, we do not need to distinguish any more between existentially/universally quantified time intervals, because all intervals are universally quantified, and we can drop the over-line notation. The resulting commitment is:

$$C(r, c, [T, T + 7days]p), t_1 \leq T, T \leq t_1 + 24days. \quad (2)$$

*Situation 1.3.* Mallya et al. propose the following solution:

$$C(h, c, [t_1, t_1 + 24hrs](\overline{[t_3, t_3 + 7days]p}), t_1 + 24hrs < t_3,$$

where  $t_1, t_1 + 24hrs$  is “today” (before spring break) and spring break starts on  $t_3$  and lasts one week. In this way, we obtain two disjoint intervals. The commitment should be resolved before the end of the first interval in order not to be breached, however it can only be resolved during the second interval, which implies that it will be necessarily breached. An alternative notation for the same commitment is the following:

$$C(h, t, [T_s, T_e]p), t_1 \leq T_s, T_e \leq t_1 + 7days, t_3 \leq T_s, T_e \leq t_3 + 24hrs. \quad (3)$$

In this way, we eliminate the need for nested intervals, and unresolvability can automatically be discovered by basic CLP inference [7].

Situation 1.3 shows that in a specific case, we can do away with nesting. In general, all existential temporal quantifiers can be mapped onto CLP domain restrictions, so the need for nesting intervals is only due to nested universal temporal quantifiers. An example of such a situation is the following:

*Example 2.* The car rental company offers a one-week free rental every month, for the whole 2009.  $\square$

In this case, we cannot do away with nested intervals. It is possible to extend Mallya et al.’s Solution 2 and write

$$\overline{[t_1, t_1 + 12months](\overline{[t_3, t_3 + 7days]p}), t_1 \leq t_3 \leq t_1 + 24days,$$

however that does not capture the “every month” concept, due to lack of domain variables. A different notation is needed. For example, we may use nested

commitments, instead of nested intervals. Alternatively, if the “every month” concept is often used in the domain, we could define a new (non-binary) constraint and a dedicated propagation algorithm which ensures a very compact notation and an efficient inference process. Non-binary (global) constraints are one of the prominent features of constraint programming frameworks. It may well be that a global constraints that we can use is already available off-the-shelf.<sup>1</sup>

*Situation 1.4.* Mallya et al. propose the following solution:

$$C(r, c, (\overline{[t_1, t_1 + 2days]}great\_car \vee [t_1, t_2]replace\_car)), t_2 < t_1 + 2days,$$

where *great\_car* means that the car has not broken down, and *replace\_car* represents the warranty that the rental company gives on the quality of the car,  $t_1$  denotes the instant at which the car is rented on Friday and  $t_2$  denotes the closing of the rental company on Friday. Using the framework presented in [10] it is possible to reason on this “warranty paradox” using CTL and realize that the warranty cannot be enjoyed if the car is rented on a Friday and it breaks down on Saturday.

Note that this modeling, however intuitive, may give rise to some counter-intuitive results. For example,  $c$  may decide to satisfy the commitment by replacing a perfectly functioning car with a broken car.

If we wish to follow Situation 1.4’s description more literally, we should opt for a different formalization. For example, the commitment about the replacement car should only be a consequence of the car breaking down:

$$C(r, c, [T]replace\_car), t_1 \leq T, T \leq t_2 \leftarrow \mathbf{H}(break\_down, T), T \leq t_1 + 2days \quad (4)$$

where by  $\mathbf{H}(break\_down(T))$  we denote an event occurred (“Happened”) at time  $T$ . Again, it is possible to reason on the “warranty paradox” using basic CLP inference. The result of such a reasoning would be a “real” validity interval for the warranty, which excludes Saturday.

Thus using a rule-based notation we can express many situations in a faithful way. In particular, it would be possible to express conditional commitments. However, there is another possible solution, based on the concept of *compensation*. A compensation is an action to be taken to recover from a situation of violation. To this end, we need to explicitly denote the state of commitments. For instance, we can write  $viol(C(x, y, p), t)$  to indicate that a commitment has been violated at time  $t$  (due to an event occurring at time  $t$  which falsifies  $p$ , or due to the elapsing at time  $t$  of a time interval in which  $p$  was supposed to be verified). We obtain:

$$C(r, c, [t_1, t_2]great\_car). \quad (5)$$

$$C(r, c, [T_r]replace\_car) \leftarrow viol(C(r, c, [T_s, T_e]great\_car), T_b), \quad (6)$$

$$T_s \leq T_b, T_b \leq T_s + 2days, T_b \leq T_e, T_b \leq T_r.$$

<sup>1</sup> See Beldiceanu and Carlsson’s global constraints catalog, <http://www.emn.fr/x-info/sdemasse/gccat/>

( $T_b$  is the time of break down,  $T_r$  is the time of replacement). Alternatively, using an extended notation, we could write:

$$\begin{aligned} & \text{compensate}(r, c, [T_r]\text{replace\_car}, C(r, c, [T_s, T_e]\text{great\_car}), T_b) \leftarrow \\ & T_s \leq T_b, T_b \leq T_s + 2\text{days}, T_b \leq T_e, T_b \leq T_r. \end{aligned} \quad (7)$$

More on compensations below.

Note that we can easily refine the rules above to specify what “immediate replacement” means (1 day? 3 hours?), by posing another constraint between  $T_b$  and  $T_r$ , other than  $T_b \leq T_r$ .

## 2.2 Time of commitments

Another temporal dimension could be interesting in many applications. It is the time of the commitment itself. To the best of our knowledge, this dimension has not been considered by previous research.

*Example 3.* Consider a university-like agent organization, in which agent  $x$  and faculty  $f$  belong to. There are roles with social responsibilities, which we can express by way of commitments. One such role is that of director of studies (*dos*).  $x$  has been appointed *dos* at Faculty  $f$  on October 29, 2008.

We can express that  $x$  has been appointed *dos* for 2009 at Faculty  $f$ , using a notation like:

$$C(x, f, [01.01.2009, 31.12.2009]\text{dos}). \quad (8)$$

But how can we express that  $x$  has been appointed *dos* on October 29th 2008? This could be an important element of the domain. Consider a regulation that says that *a Faculty member that has been appointed director of studies cannot take more commitments in the Faculty*. The notation above does not permit to reason at this level. The closest approximation is probably: *a Faculty member cannot take more commitments in the Faculty while he is director of studies*. Or, we could resort to an additional commitment to express the appointment, beside the *dos* commitment. But this would complicate the model by increasing the number of commitments. A simple solution is to attach the duration of the commitment to the commitment itself:

$$[29.10.2008, T_{end}]\text{active}(C(x, f, [01.01.2009, 31.12.2009]\text{dos})). \quad (9)$$

## 2.3 Compensations

Contracts often involve deadlines and compensations. Usually, compensation actions are possibilities given to recover from a situation of violation. In a typical setting, a commitment not satisfied in time will not become satisfied by actions taken after the deadline, but it will instead incur in a further commitment from the debtor’s side to take a compensation action. The extent of the compensation required may be subject to context-dependent conditions and be directly related to the time spent after the deadline before the compensation action is taken.



*Example 4.* According to the Italian civil code, the owners of real estate must pay taxes to the municipality (I.C.I.) between June 1 and June 16 of every tax year, for the property owned during the previous solar year. The Law allows the debtor who did not pay by the deadline to recover, by their own initiative, from such a violation by a procedure called *spontaneous revision*. The spontaneous revision procedure is permitted only if the debtor has not yet been officially notified about ongoing investigations related to such a violation. A spontaneous revision’s compensation of a previous violation amounts to 3% of the amount not paid, which counts as a sanction, plus the legal interests on the amount not paid, which depend on the time elapsed between the I.C.I. payment deadline and the payment by spontaneous revision.  $\square$

To model this example, we can resort to a *compensate* notation like we did above. Let  $t_1$  be June 1st,  $t_2$  be June 16th,  $c$  a citizen,  $m$  a municipality, and let domain-dependent knowledge such as the interest rate  $IR$  and the amount of taxes to be paid by a citizen be defined by rules or facts such as  $interest\_rate(0.025)$  and  $ici(c, 100euro)$ . A possible solution of Example 4 is the following:

$$C(c, m, [T]pay\_ICI(Amt)), t_1 \leq T, T \leq t_2 \leftarrow ici(c, Amt). \quad (10)$$

$$compensates(c, m, [T_p]pay\_ICI(Amt), C(c, m, [T_r]s\_rev(Amt_{new}))) \leftarrow \quad (11)$$

$$interest\_rate(IRate), Amt_{new} = Amt \times (1.03 + IRate \times (T_r - T_p)).$$

( $s\_rev$  stands for spontaneous revision,  $Amt$  for amount, and “=” is a CLP equality constraint).

Such examples are ubiquitous in legislation bodies, and in many domains in which contracts are used to establish rights and obligations of interacting parties. To be able to model such situations and reason about them, a notation should accommodate variables inside commitments and allow us to relate such variables with domains and expressions containing other variables.

Note that in this case *compensates* is syntactic sugar for an alternative and equally expressive notation. One could resort for instance to conditional commitments, and include original content (*pay\_ICI*) and compensating content (*s\_revision*) in one single CC-like fact. However, compensations could be defined in different ways depending on the domain. For example, various degrees of violation can be defined, such as mild/serious violation if before/after official notification. A commitment modeling framework should be flexible enough to accommodate all these needs. This notation helps to abstract away from the specific compensation semantics.

### 3 Desiderata for a commitment modeling framework

The considerations made above suggest a number of desiderata for a commitment modeling framework that enables reasoning with time. The first two desiderata are taken from [10].

- Time intervals** Contracts often involve time bounds. It should be possible to express such time bounds, in order to enable reasoning about satisfaction or breach of commitments in general.
- Achievement and maintenance** Two kinds of commitment conditions are possible: achievement conditions (a ticket will be issued by the end of the day), and maintenance conditions (the car will work without breaking down for 2 days). They should both be accommodated.
- Degrees of violation** It should be possible to reason about the extent of breach of a commitment, to capture ideas such as partial or mild violation of a contract.
- Compensation** The language should enable associating commitments with compensation actions.
- Time of commitment state changes** It should be possible to reason about the time a commitment assumes a particular state, e.g., the time a commitment is created, violated or discharged. The framework should enable reasoning about the state of commitments along time.
- Meta-level reasoning** There could be commitments about commitments (and further nesting). The notation should accommodate contracts in which a commitment is about another commitment that will be created at some later point, or about some existing commitment.
- Simplicity** The notation should be easy and at the same time rigorous. It should be possible to run automated reasoning tasks on commitment-based contract specifications. Some interesting reasoning tasks are: contract analysis, commitment tracking, compliance verification.
- Modularity** It should be possible to extend the commitment notation or modify the underlying theories and reasoning procedures in a modular way. Moreover, it should be possible to integrate a commitment framework with other domain knowledge, so as to enable reasoners and agents using commitments to reason using all available knowledge, possibly including ontological knowledge. Such an integration should preserve the modularity principle.

#### 4 A new notation for social commitments: $\mathcal{CML}$

We propose a new notation for social commitments. We call it  $\mathcal{CML}$  (Commitment Modeling Language). To enable reasoning, we consider commitments as a part of a knowledge base. In particular, social commitments are specified inside  $\mathcal{CML}$  programs ( $\mathcal{C}Programs$ ), which could describe for example a contract.

A  $\mathcal{CML}$  program is made of rules. A rule in the form

$$CRuleHead \leftarrow CRuleBody \tag{12}$$

is used to define effects of events on the state of commitments. More specifically, the user can use such rules to define for instance which events create, discharge, or break which commitments, in the style of [14]. The body defines the context, i.e., the conditions that must hold in order for an event to have an effect on the state of a commitment. If no context is defined, the rule is said to be a fact.

Atoms in the body of rules may be external predicates, not defined in the commitment specification program (this should be allowed because of the modularity principle), or they can be fluents modeling the state of commitments. Such a state can be associated with existentially quantified temporal variables (*holds\_e* notation) or with universally quantified intervals (*holds\_u* notation).

The *CMC* syntax, shown in Figure 1, is easily extensible to accommodate

```

CMLProgram ::= CRules
CRules ::= CRule[CRules]
CRule ::= CRuleHead.|CRuleHead ← CRuleBody.
CRuleHead ::= OPC(Terms,Commitment)
OPC ::= = create|discharge|cancel|release|assign|delegate|...
CRuleBody ::= CRuleBodyElem[|,CRuleBody]
CRuleBodyElem ::= holds_e(State(Commitment,Time),Time)|
holds_u(Interval State(Commitment,Time))|
Atom|Constraint
Commitment ::= C(Agent,Agent,[Interval]CAtom)
Interval ::= [TExpr[,TExpr]]
TExpr ::= Time OPT Time|Time OPT Duration
OPT ::= +|-
Time ::= Date|Numeral|Variable|TExpression
Duration ::= Numeral Granularity
Granularity ::= hrs|days|weeks|months|...
Agent ::= Term
Atom ::= Ident|Ident(Terms)
Term ::= Atom|Numeral|Variable
Terms ::= Term[,Term]
CAtom ::= Atom|Commitment
Constraint ::= Variable ∈ Domain|Variable OPCLP TExpr
State ::= viol|p_viol|active|...
OPCLP ::= = | ≠ | ≤ | ≥ | < | >
Domain ::= Interval|Set

```

**Fig. 1.** Syntax of *CMC* programs

various models of commitment evolution. For example, *p\_viol* has been introduced alongside *viol* and *active* as a possible new commitment state, and a new

type of commitment operation, e.g., *compensate*, could be added to the language to provide the user with such a high-level abstraction.

A sample *CProgram*, modeling Situation 1.4, is the following:

$$\text{create}(\text{rent\_a\_car}(T_c, T_e), \mathbf{C}(r, c, [T_c, T_c + 2\text{days}]\text{great\_car})). \quad (13)$$

$$\begin{aligned} \text{create}(\text{car\_broken}(T_b), \mathbf{C}(r, c, [T_r]\text{replace\_car})), T_r \leq T_b + 24\text{hours} \leftarrow \\ \text{holds\_e}(\text{viol}(\mathbf{C}(r, c, [T_s, T_e]\text{great\_car}), T_b), T_b). \end{aligned} \quad (14)$$

Renting a car at time  $T_c$  until  $T_e$  creates a commitment that for 2 days as of  $T_c$  the car does not break down. The car breaking down at a time  $T_b$  creates a commitment that the car must be replaced within 24 hours of the incident, if the breakdown has caused a breach of commitment.

While  $\mathcal{CML}$  provides very expressive constructs such as variables with domains, constraints and rules, on the other hand it does not explicitly accommodate temporal logic expressions, such as  $p\mathcal{U}q$  or  $\bigcirc p$ . We are currently investigating whether and how temporal logic formulae can be mapped onto  $\mathcal{CML}$  expressions.

## 5 Commitment manipulation and reasoning

Two fundamental aspects of commitment frameworks are manipulation and reasoning [13]. Manipulation operates on the state of commitments.

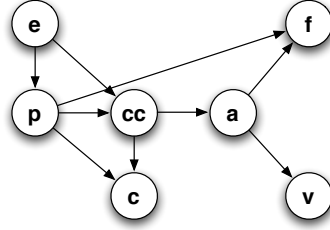
### 5.1 States of commitments

Recently, many authors proposed different possible evolutions of the commitment state, from an initial state after creation, down to its satisfaction through discharge, delegation or cancellation operations, or else to its violation due to the occurrence of events that contradict the subject of the agreement represented by the commitment itself. For instance, in [6], Fornara and Colombetti propose the following set of states for a commitment: empty (e), cancelled (c), precommitment (p), conditional (cc), active (a), fulfilled (f), and violated (v).

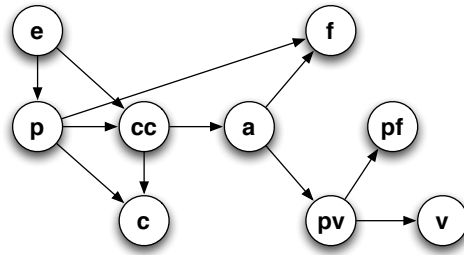
The states and their transitions are depicted in Figure 2.

Usually, once the conditions specified in the commitment are either satisfied (for an achievement commitment) or violated (for a maintenance commitment), the commitment assumes a final state. However, as discussed above, if we consider also relevant temporal aspects, such as deadlines, we could define a finer-grained characterization of the state of a commitment. For example, after a deadline has passed, the debtor may still opt for a belated commitment satisfaction. It may be interesting to distinguish among (1) commitment satisfied in time, (2) commitment “violated” before the deadline but “satisfied” after the deadline (partial violation/partial satisfaction), and (3) violated commitment.

Such a distinction is depicted in Figure 3 as an extension to Fornara and Colombetti’s transitions.



**Fig. 2.** Fornara & Colombetti's commitment state transitions [6]



**Fig. 3.** A possible extension to a commitment state transitions accounting for partial violation (pv) and partial fulfillment (pf) of commitments.

## 5.2 Reasoning about Commitments

Yolum and Singh [14] propose to reason about commitments using the Event Calculus ( $\mathcal{EC}$ ) [8]. The  $\mathcal{EC}$  was introduced by Kowalski and Sergot as a logic programming framework for representing and reasoning about events and their effects. Basic concepts are that of *event*, happening at a point in time, and *property* (or *fluent*), holding during time intervals. Fluents are initiated/terminated by occurring events. There are many different formulations of the  $\mathcal{EC}$  axioms. A simple one, taken from [4], is the one below ( $F$  stands for *Fluent*,  $Ev$  for *Event*).

$$\begin{aligned} \text{holds\_at}(F, T) \leftarrow & \text{initiates}(Ev, F, T_{\text{start}}) \\ & \wedge T_{\text{start}} < T \wedge \neg \text{clipped}(T_{\text{start}}, F, T). \end{aligned} \quad (ec_1)$$

$$\begin{aligned} \text{clipped}(T_1, F, T_3) \leftarrow & \text{terminates}(Ev, F, T_2) \\ & \wedge T_1 < T_2 \wedge T_2 < T_3. \end{aligned} \quad (ec_2)$$

$$\begin{aligned} \text{initiates}(Ev, F, T) \leftarrow & \text{happens}(Ev, T) \wedge \text{holds}(F_1, T) \\ & \wedge \dots \wedge \text{holds}(F_N, T). \end{aligned} \quad (ec_3)$$

$$\begin{aligned} \text{terminates}(Ev, F, T) \leftarrow & \text{happens}(Ev, T) \wedge \text{holds}(F_1, T) \\ & \wedge \dots \wedge \text{holds}(F_N, T). \end{aligned} \quad (ec_4)$$

Axioms  $ec_1$  and  $ec_2$  are the general ones of  $\mathcal{EC}$ , whereas  $ec_3$  and  $ec_4$  are user-defined, domain-specific axioms. The  $\mathcal{EC}$  is a suitable formalism to specify the effects of commitment manipulation, and reason from such operations. As a sample fragment of Yolum and Singh’s formalization, consider a *create* operation, whose purpose is to establish a commitment, and can only be performed by the debtor. To express that an event  $e(x)$  carried out by  $x$  at time  $t$  creates a commitment  $C(x, y, p)$ , Yolum and Singh define the operation  $create(e(x), C(x, y, p))$  in terms of  $happens(e(x), t) \wedge initiates(e(x), C(x, y, p), t)$ .

In the same way, the semantics of  $\mathcal{CML}$  can be given in terms of  $\mathcal{EC}$  programs. This helps *simplicity*, because the language of  $\mathcal{EC}$  is very simple, and *modularity*, because for different domains we can define different theories of commitments.

The  $\mathcal{EC}$  is an effective framework for temporal reasoning. It has been extensively used in the literature to carry out two main reasoning tasks: deductive *narrative verification*, to check whether a certain fluent holds given a narrative (set of events), and abductive *planning*, to simulate a possible narrative which satisfies some requirements [11]. Chittaro and Montanari [3] proposed a way to use the  $\mathcal{EC}$  for run-time monitoring and verification. It is based on a mechanism to cache the outcome of the inference process every time the knowledge base is updated by a new event. In a nutshell, the *Cached Event Calculus (CEC)* computes and stores fluents’ *maximum validity intervals (MVIs)*, which are the maximum time intervals in which fluents hold, according to the known events. The set of cached validity intervals is then extended/revised as new events occur or get to be known. Therefore, the  $\mathcal{EC}$  can be used as a basis for reasoning on commitments in many ways, including not only planning and static verification, but also tracking, depending on the  $\mathcal{EC}$  implementation used (abductive, deductive, reactive).

## 6 Social commitment framework architecture

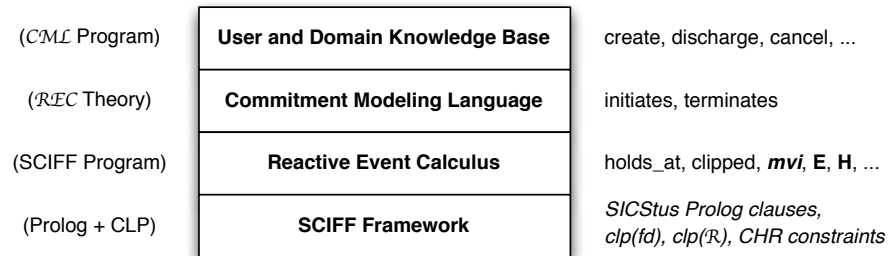
We propose an abstract, layered architecture that enables modeling and reasoning with social commitments. It consists of:

- a user application layer;
- a commitment modeling layer;
- a temporal representation and reasoning layer;
- a reasoning and verification layer.

On the top layer, the user can define contracts or agent social interaction rules using commitments. Such definitions are based on a language provided by the layer below. The commitment modeling language is implemented using a temporal representation and reasoning framework, which is in turn built on top of a more general reasoning and verification framework, which lies at the bottom layer. It is important to rely on a formal framework that accommodates various forms of verification, because in this way commitments can be operationalized and the user can formally analyze commitment-based contracts, reason on the

state of commitments, plan for actions needed to reach states of fulfillment, and track the evolution of commitments at run-time.

Indeed, the underlying reasoning and verification layer must be powerful enough to implement a temporal representation and reasoning layer. We propose a concrete instance of such an architecture, represented in Figure 4.



**Fig. 4.** Social commitment framework architecture

At the bottom layer, we find a number of Prolog+CLP modules which implement the SCIFF family of proof-procedures and provide the SCIFF language to the layer above [1]. The SCIFF framework is based on abductive logic programming and it consists of a declarative specification language and a family of proof-procedures for reasoning from SCIFF specifications. Some kinds of reasoning are: deduction, hypothetical reasoning, static verification of properties, compliance checking and run-time monitoring. In general, SCIFF comes in hand for a number of useful tasks in the context of agent interaction. A high-level description of SCIFF and of its usage is given in [13], also in relation with commitments. The CLP solvers integrated in SCIFF can work with discrete and dense domains, depending on the application needs, and they are particularly useful for reasoning along the temporal dimension.

On top of the SCIFF layer there is the SCIFF implementation of the  $\mathcal{EC}$ ,  $\mathcal{REC}$ . There are several implementations for  $\mathcal{EC}$ . One of them which uses ideas taken from  $\mathcal{CEC}$  and thus enables runtime verification, is called the Reactive Event Calculus ( $\mathcal{REC}$ ) and it is implemented in SCIFF. Thus  $\mathcal{REC}$  is implemented as a SCIFF program. This layer provides to the layer above the  $\mathcal{REC}$  language, which consists of the domain-dependent axioms  $ec_3$  and  $ec_4$ .

In the third layer, the constructs that define the Commitment Modeling Language (*CML*), i.e., the notation proposed above, are written by way of  $\mathcal{REC}$  theories. Thus this layer will provide the language to write a *CProgram* to the top layer. The top layer consists of user and domain-dependent knowledge encoded into a *CProgram*. An example of a program for the top layer was given in Section 4.

We believe that such an architecture, and its instantiation based on SCIFF,  $\mathcal{REC}$ , and the  $\mathcal{CML}$ , can successfully address the desiderata identified above. Modularity is achieved in two directions: in the vertical direction, by making  $\mathcal{CML}$  programs,  $\mathcal{EC}$  theory, and commitment theories independent of each other, and in the horizontal direction, by letting the user refer to external inputs by way of simple atoms. Atoms can be mapped into function calls via suitable interfaces such as those available in most Prolog engines.  $\mathcal{CML}$  is a simple and easily extensible language, which consists of primitives such as *create*, *discharge*, etc., in the style of [14]. The language is expressive enough to express time intervals, achievement and maintenance conditions, and time of commitment state change. Thanks to the expressivity of the language and to the modularity of the architecture, it is possible to extend the framework to model different kinds of violation and powerful new constructs such as compensation. In fact, the states of commitments and manipulation operations are not hard-wired in the architecture, but they can be (re)defined by the user. Finally,  $\mathcal{CML}$  accommodates meta-level reasoning on commitments, and the underlying  $\mathcal{REC}$  engine can reason about commitments at all levels by treating a commitment state as a fluent which holds for a period of time.

## 7 Conclusion

We identified some issues regarding the representation of commitments which are still open, and we formulated a number of desiderata for a commitment modeling framework. To the best of our knowledge, in the state of the art there is no framework that satisfies all the desiderata. We believe that a possible answer could come from a commitment framework organized into four layers.

On top of the stack, at the user level, contracts can be specified by way of commitment programs. We identified in SCIFF a potential candidate for the bottom layer and we defined a notation for top-level programs. A working implementation exists of the two bottom layers.<sup>2</sup> At the current stage, the notation is still a preliminary proposal.

Our discussion was informal and example-driven. We gave emphasis to temporal aspects of commitments in relation with deadlines. However, deadlines are only a special case of temporal constraints and CLP constraints in general. Surely there are many other types of constraint that could be very useful for modeling the domain correctly and compactly. In particular, global constraints capture common patterns and help specify complex and recurring constraints in a simple way. Each global constraint comes with an effective and efficient propagation algorithm capable of powerful inference. A useful activity could be to isolate a subset of CLP constraints of interest for commitment-related domains. A concrete implementation of the commitment modeling framework should include a library of such CLP constraints. To the best of our knowledge, the inference potential of such a technology, unlocked by the architecture we propose, is unprecedented in the domain of commitments.

<sup>2</sup> <http://lia.deis.unibo.it/research/sciff/>



Other future work will aim to give a concrete instantiation and implementation of the top two layers, and on studying the formal properties of the overall framework and the expressivity of the top-level commitment programming language. We also intend to investigate how  $\mathcal{CML}$  and the abstract architecture fit into concrete application domains. Among them, the domain of institutions. Commitments or obligations are also included in their modeling, and often they work with deadlines of events taking place instead of time.

**Acknowledgments.** We thank the anonymous reviewers for their useful comments. This work has been partially supported by the FIRB project *TOCAI.IT*.

## References

1. M. Alberti, F. Chesani, M. Gavanelli, E. Lamma, P. Mello, and P. Torroni. Verifiable agent interaction in abductive logic programming: the SCIFF framework. *ACM Transactions on Computational Logic*, 9(4):1–43, 2008.
2. C. Castelfranchi. Commitments: From individual intentions to groups and organizations. In V. R. Lesser and L. Gasser, editors, *Proceedings of the First International Conference on Multiagent Systems*, pages 41–48. The MIT Press, 1995.
3. L. Chittaro and A. Montanari. Efficient temporal reasoning in the cached event calculus. *Computational Intelligence*, 12(2):359–382, Aug. 1996.
4. L. Chittaro and A. Montanari. Temporal representation and reasoning in artificial intelligence: Issues and approaches. *Annals of Mathematics and Artificial Intelligence*, 28(1-4):47–106, 2000.
5. M. Fisher, R. H. Bordini, B. Hirsch, and P. Torroni. Computational logics and agents: A road map of current technologies and future trends. *Computational Intelligence*, 23(1):61–91, Feb. 2007.
6. N. Fornara and M. Colombetti. Operational specification of a commitment-based agent communication language. In *Proc. 1st AAMAS*, pages 536–542. ACM, 2002.
7. J. Jaffar and M. Maher. Constraint logic programming: a survey. *Journal of Logic Programming*, 19-20:503–582, 1994.
8. R. A. Kowalski and M. Sergot. A logic-based calculus of events. *New Generation Computing*, 4(1):67–95, 1986.
9. A. U. Mallya and M. N. Huhns. Commitments among agents. *IEEE Internet Computing*, 7(4):90–93, 2003.
10. A. U. Mallya, P. Yolum, and M. P. Singh. Resolving commitments among autonomous agents. In F. Dignum, editor, *Advances in Agent Communication*, volume *Lecture Notes in Computer Science 2922*, pages 166–182. Springer, 2004.
11. M. Shanahan. An abductive event calculus planner. *Journal of Logic Programming*, 44(1-3):207–240, 2000.
12. M. P. Singh. An ontology for commitments in multiagent systems: Toward a unification of normative concepts. *Artificial Intelligence and Law*, 7:97–113, 1999.
13. P. Torroni, P. Yolum, M. P. Singh, M. Alberti, F. Chesani, M. Gavanelli, E. Lamma, and P. Mello. Modelling interactions via commitments and expectations. In V. Dignum, editor, *Handbook of Research on MAS: Semantics and Dynamics of Organizational Models*, pages 263–284. Hershey, Pennsylvania, Mar. 2009. IGI Global.
14. P. Yolum and M. Singh. Flexible protocol specification and execution: applying event calculus planning using commitments. In *Proc. 1st AAMAS*, pages 527–534. ACM Press, 2002.

# Verifying Dribble Agents

Doan Thu Trang, Brian Logan, and Natasha Alechina

The University of Nottingham  
School of Computer Science

**Abstract.** We describe a model-checking based approach to verification of programs written in the agent programming language Dribble. We define a logic (an extension of the branching time temporal logic CTL) which describes transition systems corresponding to a Dribble program, and show how to express properties of the agent program in the logic and how to encode transition systems as an input to a model-checker. We prove soundness and completeness of the logic and a correspondence between the operational semantics of Dribble and the models of the logic.

## 1 Introduction

BDI-based agent-oriented programming languages [5] facilitate the implementation of cognitive agents by providing programming constructs to implement concepts such as beliefs, goals, and (pre-defined) plans. In such languages, an agent selects a plan to achieve one or more goals based on its beliefs about the environment. However in anything other than toy environments, selecting an appropriate plan does not guarantee that it can be successfully executed. The beliefs used to select a particular plan for a given goal is only a heuristic, and cannot capture the preconditions of all the actions in the plan (some of which may be false when the plan is selected and will only be made true by actions in the plan). Moreover, in dynamic environments, an agent's beliefs (and hence the best way of achieving a goal) may change in unanticipated ways during plan execution, and a rational agent must be prepared to revise its plans at run time to take advantage of 'fortuitous' changes in the environment (e.g., which allow some steps in the plan to be skipped) or to recover from 'adverse' changes in the environment (e.g., when a precondition of an action is discovered not to hold).

Many BDI-based agent programming languages provide facilities to drop plans if the corresponding goal is 'unexpectedly' achieved or when execution of the plan fails [6, 16, 17]. More advanced languages, e.g., [18, 9] provide support for arbitrary modification of plans during their execution. However, while such meta-level capabilities simplify the development of rational agents, they make it more difficult to reason about the execution of agent programs, e.g., to verify their correctness. In addition to reasoning about the agent's beliefs, goals and plans, we need to model the current state of plan execution, and the evolution of the agent's program at run time in response to interactions between the effects the agent's actions in its environment and its plan revision capabilities.

In this paper we present an approach to verifying agent programs which admit arbitrary revisions at run time. We focus on the BDI agent programming language Dribble

introduced in [18]. Dribble allows the implementation of agents with beliefs, (declarative) goals, actions, abstract actions (procedural goals), plans, and rules for selecting and revising plans. Although relatively simple and abstract, it is representative of a wider class of BDI agent programming languages which support plan revision, and presents significant challenges for verification. Our approach is based on model-checking. We define a logic (an extension of the branching time temporal logic CTL) which describes transition systems corresponding to a Dribble program, and show how to express properties of the program in the logic and how to encode transition systems as an input to a model-checker. We prove soundness and completeness of the logic and a correspondence between the operational semantics of Dribble and the models of the logic.

The rest of the paper is organised as follows. In the next section, we describe the syntax and operational semantics of Dribble. In section 3 we introduce a logic for expressing properties of Dribble programs, and give a complete axiomatisation of the set of models corresponding to the operational semantics of a Dribble agent program. We discuss the use of the logic for verification in section 4, where we describe model-checking of Dribble programs and give a simple example of a program and a property which can be model-checked. We give a brief survey of related work in section 5.

## 2 Dribble

In this section, we briefly review the syntax and operational semantics of Dribble.

### 2.1 Beliefs and goals

Let  $Prop$  be a finite set of propositional variables and  $\mathcal{L}$  the set of propositional formulas. In order to make  $\mathcal{L}$  finite, we allow  $\mathcal{L}$  to contain only formulas in Disjunctive Normal Form (DNF). A formula is said to be in DNF iff it is a disjunction of conjunctive clauses in which a conjunctive clause is a conjunction of literals. As usual, a literal is either  $p$  or  $\neg p$  for any  $p \in Prop$ . Moreover, formulas of  $\mathcal{L}$  satisfy the following conditions:

1. formulas do not contain duplicates of conjunctive clauses;
2. conjunctive clauses of a formula do not contain duplicates of literals; and
3. literals in a conjunctive clause of a formula only occur in some fixed order.

A Dribble agent has both a belief base and a goal base which are finite subsets of  $\mathcal{L}$ . The agent's beliefs and goals are expressed in a language  $\mathcal{L}_{BG}$ . The syntax of  $\mathcal{L}_{BG}$  is defined as follows:

$$\beta \leftarrow \mathbf{B}\alpha \mid \mathbf{G}\alpha \mid \neg\beta \mid \beta_1 \wedge \beta_2 \text{ where } \alpha \in \mathcal{L}.$$

The meaning of  $\mathbf{B}\alpha$  is that  $\alpha$  can be propositionally derived from the belief base of an agent, and  $\mathbf{G}\alpha$  means that  $\alpha$  is the consequence of some single goal in the goal base of an agent. For convenience, a formula  $\beta$  is of  $\mathcal{L}_B$  ( $\mathcal{L}_G$ ) iff it does not contain any subformula of the form  $\mathbf{G}\alpha$  ( $\mathbf{B}\alpha$ , respectively).

A formula of  $\mathcal{L}_{BG}$  is interpreted by a pair of a belief base and a goal base  $\langle \delta, \gamma \rangle$ , in which both  $\delta$  and  $\gamma$  are finite subsets of formulas of  $\mathcal{L}$ . The truth of a formula  $\beta$  is defined inductively as follows.

- $\langle \delta, \gamma \rangle \models_{BG} \mathbf{B}\alpha \Leftrightarrow \delta \models_{Prop} \alpha$
- $\langle \delta, \gamma \rangle \models_{BG} \mathbf{G}\alpha \Leftrightarrow \exists g \in \gamma : g \models_{Prop} \alpha$
- $\langle \delta, \gamma \rangle \models_{BG} \neg\varphi \Leftrightarrow \langle \delta, \gamma \rangle \not\models_{BG} \varphi$
- $\langle \delta, \gamma \rangle \models_{BG} \beta \wedge \beta' \Leftrightarrow \langle \delta, \gamma \rangle \models_{BG} \beta \text{ and } \langle \delta, \gamma \rangle \models_{BG} \beta'$

## 2.2 Plans

A Dribble plan consists of basic actions and abstract plans composed by sequence and conditional choice operators. The sequence operator, ‘;’, takes two plans,  $\pi_1, \pi_2$ , as arguments and states that  $\pi_1$  should be performed before  $\pi_2$ . The conditional choice operator allows branching and generates plans of the form ‘if  $\phi$  then  $\pi_1$  else  $\pi_2$ ’. The syntax of plans is defined as follows:

$$\pi \leftarrow a \mid b \mid \text{if } \beta \text{ then } \pi'_1 \text{ else } \pi'_2 \mid \pi'_1; \pi'_2$$

where  $a$  is an abstract plan,  $b$  is a basic action and  $\beta \in \mathcal{L}_B$ .

We depart from [18] in that we do not have an empty plan (denoted by  $E$  in [18]) as a special kind of plan which can occur as part of other plans. Below, we will use  $E$  as a marker for an empty plan base, but not as a plan expression, to avoid introducing rewriting rules for  $E; E$  to  $E$  and  $\pi_1; E; \pi_2$  to  $\pi_1; \pi_2$ , etc.

We define length of a plan  $\pi$ ,  $len(\pi)$ , inductively as follows:

$$\begin{aligned} len(a) &= 1 \\ len(b) &= 1 \\ len(\text{if } \beta \text{ then } \pi'_1 \text{ else } \pi'_2) &= len(\pi'_1) + len(\pi'_2) + 4 \\ len(\pi'; \pi) &= len(\pi') + len(\pi) \end{aligned}$$

Notice that in the case of the if-then-else statement, the length is the sum of lengths of the plans  $\pi'_1$  and  $\pi'_2$  together with the number of extra symbols of the statement, i.e. *if*, *then*, *else* and  $\beta$ .

Since in reality, agents can hold a plan up to some fixed length, we make an assumption that all plans have length smaller than a certain preset number. Restricting the length of plans also makes the set of plans finite. This is necessary for the axiomatisation of the logic later in the paper.

In the rest of this paper, we denote by  $Plans$  the set of all plans whose lengths are smaller than  $len_{MAX}$ , where  $len_{MAX}$  is a natural number.

$$Plans = \{\pi \mid len(\pi) \leq len_{MAX}\}$$

## 2.3 Dribble agents

Writing a Dribble agent means writing a number of goal rules and practical reasoning rules. The syntax of goal rules (PG) and practical reasoning (PR) rules is given below.

- PG rules:  $\beta \rightarrow \pi$  where  $\beta \in \mathcal{L}_{BG}$  and  $\pi \in Plans$
- PR rules:  $\pi_1 \mid \beta \rightarrow \pi_2$  where  $\beta \in \mathcal{L}_B$  and  $\pi_1, \pi_2 \in Plans$ , and  $\pi_2$  may be empty.

One writes a PG rule to intend that an agent with an empty plan base will generate a plan  $\pi$  if its current belief and goal bases satisfy the condition encoded in  $\beta$ . If the agent has certain goals in its goal base, it will generate a plan based on its beliefs to hopefully achieve those goals. A PR rule proposes a possible revision  $\pi_2$  to (the prefix of) a plan  $\pi_1$  which is applicable if the belief base satisfies the condition encoded in  $\beta$ . That is, if the agent has certain beliefs which imply that the current plan will be unable to achieve the intended goal(s) or that the plan is redundant and can be simplified, it can modify the plan. Note that  $\pi_2$  can be empty, allowing the agent to drop part or all of a plan.

We have slightly modified the meaning of PR rules given in [18]. In Dribble, these rules apply to complete plans ( $\pi_1$  is the agent's plan in its entirety, not a plan prefix, for example a name for an abstract plan). In contrast we allow  $\pi_1$  to be a prefix of the agent's plan, which is replaced by  $\pi_2$  followed by the continuation of the original plan. We could have written PR rules as  $\pi'_1; \pi \mid \beta \rightarrow \pi_2; \pi$  where  $\pi$  is a plan variable. In cases where  $\pi_1$  matches the entire plan, our PR rules are equivalent to those in [18]. We believe that our generalisation is justified programmatically, and it presents an interesting challenge for logical formalisation, in particular model-checking. To enforce our assumption about the length of plans, we require that Dribble agents consist of PG and PR rules which do not produce plans of length more than  $len_{MAX}$ .

A Dribble agent only has one intention at a time, i.e., its plan base contains at most one plan and it can apply a goal rule only when its plan is empty, and is strongly committed to its goals, i.e., an agent drops a goal only when it believes that the goal has been achieved.

A **Dribble agent** is a tuple  $\langle \delta, \gamma, \Gamma, \Delta \rangle$  in which  $\Gamma$  is a set of goal rules,  $\Delta$  is set of practical reasoning rules,  $\delta$  and  $\gamma$  are the initial belief base and goal base and both satisfy the following conditions:

1.  $\delta$  is consistent
2.  $\forall \alpha \in \gamma, \delta \not\models_{Prop} \alpha$
3.  $\forall \alpha \in \gamma, \alpha$  is consistent

that is, the agent's beliefs are consistent, it does not have as a goal any formula it already believes to be the case, and each of its goals is consistent (though its goals may be inconsistent with each other as they can be achieved at different times).

## 2.4 Operational semantics

In this section, we describe how a Dribble agent operates. A Dribble program  $P$  is a pair  $(\Gamma, \Delta)$  of PG rules and PR rules.

A configuration of an agent is a tuple  $\langle \delta, \gamma, \{\pi\} \rangle$  where  $\delta, \gamma$  and  $\pi$  are the agent's current belief base, goal base and plan base (where  $\pi$  is the current plan, possibly partially executed), respectively. In what follows, we will omit the set brackets around the plan for readability, as in  $\langle \delta, \gamma, \pi \rangle$ . The plan base can also be empty, which we will write as  $\langle \delta, \gamma, \emptyset \rangle$ . The *initial configuration* of an agent is  $\langle \delta_0, \gamma_0, \emptyset \rangle$ .

We specify the operational semantics of a Dribble agent as a set of transition rules. Each transition corresponds to a single execution step and takes the system from one configuration/state to another. In the cases corresponding to applying PG and PR rules,

we have additional conditions to guarantee that we do not produce a plan of length more than  $len_{MAX}$ . Notice that transitions from a configuration in which the plan base begins with an abstract plan are included in application of PR rules.

Application of a goal rule

$$\frac{\varphi \rightarrow \pi \in \Gamma \quad len(\pi) \leq len_{MAX} \quad \langle \delta, \gamma \rangle \models_{BG} \varphi}{\langle \delta, \gamma, \emptyset \rangle \xrightarrow{apply(\varphi \rightarrow \pi)} \langle \delta, \gamma, \pi \rangle}$$

Application of a plan revision rule

$$\frac{\pi_1 \mid \beta \rightarrow \pi_2 \in \Delta \quad len(\pi_2; \pi) \leq len_{MAX} \quad \langle \delta, \gamma \rangle \models_{BG} \beta}{\langle \delta, \gamma, \pi_1; \pi \rangle \xrightarrow{apply(\pi_1 \mid \beta \rightarrow \pi_2)} \langle \delta, \gamma, \pi_2; \pi \rangle}$$

Basic action execution

$$\frac{\mathcal{T}(b, \delta) = \delta' \quad \gamma' = \gamma \setminus \{g \in \gamma \mid \delta' \models_{Prop} g\}}{\langle \delta, \gamma, b; \pi \rangle \xrightarrow{execute(b)} \langle \delta', \gamma', \pi \rangle}$$

where  $\mathcal{T}$  is a belief update function which takes an action and a belief base and returns the resulting belief base.  $\mathcal{T}$  is a partial function since an action may not be applicable in some situations.

Conditional statement

$$\frac{\langle \delta, \gamma \rangle \models_{BG} \beta}{\langle \delta, \gamma, \text{if } \beta \text{ then } \pi_1 \text{ else } \pi_2; \pi \rangle \xrightarrow{execute(if)} \langle \delta, \gamma, \pi_1; \pi \rangle}$$

$$\frac{\langle \delta, \gamma \rangle \not\models_{BG} \beta}{\langle \delta, \gamma, \text{if } \beta \text{ then } \pi_1 \text{ else } \pi_2; \pi \rangle \xrightarrow{execute(if)} \langle \delta, \gamma, \pi_2; \pi \rangle}$$

Note that in the last three rules,  $\pi$  may be absent (or be an empty string), in which case for example executing  $b; \pi$  will result in  $\langle \delta', \gamma', \emptyset \rangle$ .

For technical reasons, if in a configuration  $\langle \delta, \gamma, \pi \rangle$  no transition rule is applicable, we assume that there is a special ‘stall’ transition to the same configuration:  $\langle \delta, \gamma, \pi \rangle \xrightarrow{stall} \langle \delta, \gamma, \pi \rangle$ .

A computation tree  $CT(c_0, P)$  for a Dribble agent with a program  $P = (\Gamma, \Delta)$  is a tree with root  $c_0$  where each node is a configuration, such that for each node  $c$  and each child  $c'$  of  $c$ ,  $c \rightarrow c'$  is a transition in the transition system for  $P$ . The meaning of a Dribble agent  $\langle \delta_0, \gamma_0, P \rangle$  is a tree  $CT(\langle \delta_0, \gamma_0, \emptyset \rangle, P)$ .

### 3 A logic of Dribble programs

In this section, we introduce a logic which allows us to formalize the properties of Dribble agent programs. Formulas of the logic will be used as input to the model-checker. In addition, we give a complete axiomatisation of the models of the logic. Axiomatisation is, of course not necessary for model-checking, but it helps us to understand the logic and its models; for example, the axioms may be more intuitive or clearer than the semantic conditions on models.

The language of our logic  $\mathcal{L}_D$  is based on Computation Tree Logic (CTL) [7] which is a logic for reasoning about branching time. The syntax of CTL is as follows:

$$\mathcal{L}_{CTL} : \varphi \leftarrow p \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid EX\varphi \mid E(\varphi U\psi) \mid A(\varphi U\psi) \quad \text{where } p \in \mathcal{L}$$

The meaning of the temporal operators is as follows:  $EX\varphi$  means there is a successor state which satisfies  $\varphi$ ;  $E(\varphi U\psi)$  means that there is a branch where  $\varphi$  holds until  $\psi$  becomes true;  $A(\varphi U\psi)$  means that on all branches,  $\varphi$  holds until  $\psi$ .

### 3.1 Syntax

$\mathcal{L}_D$  extends  $\mathcal{L}_{CTL}$  with belief, goal and plan operators (**Bs**, **Gs** and **P**).

$$\mathcal{L}_D : \varphi \leftarrow \mathbf{Bs}\delta \mid \mathbf{Gs}\gamma \mid \mathbf{P}\pi \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid EX\varphi \mid E(\varphi U\psi) \mid A(\varphi U\psi)$$

where  $\delta, \gamma \subseteq \mathcal{L}$ ;  $\pi \in Plans \cup \{E\}$ . **Bs** and **Gs** describe the belief base and goal base of the agent. Note that these operators apply to *sets* of formulas. **P** is used to describe the plan base of the agent. If the agent's plan is  $\pi$ , this is expressed as  $\mathbf{P}\pi$ , and if the plan base is empty, this is expressed as  $\mathbf{P}E$ .

We will use the usual abbreviation:

$$\begin{aligned} AX\varphi &= \neg EX\neg\varphi \text{ (in all successor states, } \varphi) \\ AF\varphi &= A(\top U\varphi) \text{ (on all branches, in some future state, } \varphi) \\ EF\varphi &= E(\top U\varphi) \text{ (there exists a branch, where in some future state, } \varphi) \\ AG\varphi &= \neg EF\neg\varphi \text{ (on all branches, in all states, } \varphi) \\ EG\varphi &= \neg AF\neg\varphi \text{ (there is a branch, where in all states, } \varphi). \end{aligned}$$

### 3.2 Semantics

In this section we define models for the logic. We show in section 4 that they correspond exactly to the computation trees for Dribble agents generated by the operational semantics.

Given a Dribble agent program  $P = (I, \Delta)$ , a Dribble model of **P** is a triple  $M_P = (S, R, V)$  in which:

- $S$  is a nonempty set of states
- $R \subseteq S \times S$  satisfies the properties below
- $V = (V_b, V_g, V_p)$  a collection of three functions,  $V_b(s) : S \rightarrow 2^{\mathcal{L}}$ ,  $V_g(s) : S \rightarrow 2^{\mathcal{L}}$  and  $V_p(s) : S \rightarrow 2^{Plans}$  satisfying the following conditions, for all  $s \in S$ :
  1.  $V_b(s)$  and  $V_g(s)$  are finite subsets of propositional formulas
  2.  $V_b(s)$  is consistent
  3.  $\forall \alpha \in V_g(s) : V_b(s) \not\models_{Prop} \alpha$
  4.  $\forall \alpha \in V_g(s) : \alpha$  is (propositionally) consistent
  5.  $V_p(s)$  is a singleton or an empty set.

To simplify the definition, we use the following conventions:  $\forall g \in \Gamma$  of the form  $\varphi \rightarrow \pi$  then  $guard(g) = \varphi$ ,  $body(g) = \pi$ , i.e.  $guard(g)$  specifies the mental condition for which situation is a good idea to execute the plan, and  $body(g)$  is the plan generated after firing the goal rule. Also,  $\forall r \in \Delta$  of the form  $\pi_1 \mid \beta \rightarrow \pi_2$  then  $head(r) = \pi_1$ ,  $guard(r) = \beta$  and  $body(r) = \pi_2$ . If  $V_p(s) = \{\pi\}$ , we will write  $V_p(s) = \pi$  for readability.

Further requirements for  $R$  are listed below.

**EPG:** For all  $s \in S$ , if  $V_p(s) = \emptyset$  and there exists  $g \in \Gamma$  such that  $\langle V_b(s), V_g(s) \rangle \models_{BG} guard(g)$  then there is  $s' \in S$  such that  $(s, s') \in R$  with  $V_b(s') = V_b(s)$ ,  $V_g(s') = V_g(s)$  and  $V_p(s') = body(g)$

**APG:** For all  $(s, s') \in R$  such that  $V_p(s) = \emptyset$ , then  $V_b(s') = V_b(s)$ ,  $V_g(s') = V_g(s)$  and there is  $g \in \Gamma$  such that  $\langle V_b(s), V_g(s) \rangle \models_{BG} guard(g)$  and  $V_p(s') = body(g)$

**EBA:** For all  $s \in S$ , if  $V_p(s) = b; \pi$  then there is  $s' \in S$  such that  $(s, s') \in R$  with  $V_b(s') = T(b, V_b(s))$ ,  $V_g(s') = V_g(s) \setminus \{g \in V_g(s) \mid V_b(s') \models_{Prop} g\}$  and  $V_p(s') = \pi$

**EIF:** For all  $s \in S$ , if  $V_p(s) = \pi_{if}; \pi$ , where  $\pi_{if} = \text{if } \beta \text{ then } \pi_1 \text{ else } \pi_2$ , then there is  $s' \in S$  such that  $(s, s') \in R$  with  $V_b(s') = V_b(s)$ ,  $V_g(s') = V_g(s)$  and

$$V_p(s') = \begin{cases} \pi_1; \pi & \text{if } \langle V_b(s), V_g(s) \rangle \models_{BG} \beta \\ \pi_2; \pi & \text{otherwise} \end{cases}$$

**EPR:** For all  $s \in S$ , if  $V_p(s) = \pi_1; \pi$  and there exists  $r \in \Delta$  such that  $head(r) = \pi_1$  and  $\langle V_b(s), V_g(s) \rangle \models_{BG} guard(r)$  then there is  $s' \in S$  such that  $(s, s') \in R$  with  $V_b(s') = V_b(s)$ ,  $V_g(s') = V_g(s)$  and  $V_p(s') = body(r); \pi$

**ABAvPR:** For all  $(s, s') \in R$ , such that  $V_p(s) = b; \pi'; \pi$ , where  $\pi'$  might be empty, then either of the following is true:

1.  $V_b(s') = T(b, V_b(s))$ ,  $V_g(s') = V_g(s) \setminus \{g \in V_g(s) \mid V_b(s') \models_{Prop} g\}$  and  $V_p(s') = \pi'; \pi$
2.  $V_b(s') = V_b(s)$ ,  $V_g(s') = V_g(s)$ , and there is  $r \in \Delta$  such that  $\langle V_b(s), V_g(s) \rangle \models_{BG} guard(r)$ ,  $head(r) = b; \pi'$  and  $V_p(s') = body(r); \pi$

**AIFvPR:** For all  $(s, s') \in R$  such that  $V_p(s) = \pi_{if}; \pi'; \pi$  ( $\pi'$  might be empty), then either of the following is true:

1.  $V_b(s') = V_b(s)$ ,  $V_g(s') = V_g(s)$ , and

$$V_p(s') = \begin{cases} \pi_1; \pi'; \pi & \text{if } \langle V_b(s), V_g(s) \rangle \models_{BG} \beta \\ \pi_2; \pi'; \pi & \text{otherwise} \end{cases}$$

2.  $V_b(s') = V_b(s)$ ,  $V_g(s') = V_g(s)$ , and there is  $r \in \Delta$  such that  $\langle V_b(s), V_g(s) \rangle \models_{BG} guard(r)$ ,  $head(r) = \pi_{if}; \pi'$  and  $V_p(s') = body(r); \pi$

**APR:** For all  $(s, s') \in R$  such that  $V_p(s) = a; \pi'; \pi$  ( $\pi'$  might be empty), then  $V_b(s') = V_b(s)$ ,  $V_g(s') = V_g(s)$ , and there is  $r \in \Delta$  such that  $\langle V_b(s), V_g(s) \rangle \models_{BG} guard(r)$ ,  $head(r) = a; \pi'$  and  $V_p(s') = body(r); \pi$

For any state  $s$  such that there are no transitions  $R(s, s')$  required by the conditions above, we stipulate  $R(s, s)$ . This is required to make the transition relation serial. No other transitions apart from those required by the conditions above exist in the model.



Given a model  $M_P$  and a state  $s$  of  $M_P$ , the truth of a  $\mathcal{L}_D$  formula is defined inductively as follows:

- $M_P, s \models \mathbf{Bs}\delta \Leftrightarrow V_b(s) = \delta$
- $M_P, s \models \mathbf{Gs}\gamma \Leftrightarrow V_g(s) = \gamma$
- $M_P, s \models \mathbf{P}\pi \Leftrightarrow V_p(s) = \pi$
- $M_P, s \models \mathbf{PE} \Leftrightarrow V_p(s) = \emptyset$
- $M_P, s \models \neg\varphi \Leftrightarrow M_P, s \not\models \varphi$
- $M_P, s \models \varphi_1 \wedge \varphi_2 \Leftrightarrow M_P, s \models \varphi_1$  and  $M_P, s \models \varphi_2$
- $M_P, s \models EX\varphi \Leftrightarrow \exists s' : (s, s') \in R : M_P, s' \models \varphi$
- $M_P, s \models E(\varphi U \psi) \Leftrightarrow \exists \text{ path } (s_0, s_1, \dots, s_n) \text{ such that:}$   
 $s_0 = s; n \geq 0; (s_i, s_{i+1}) \in R \forall 0 \leq i < n$  and  $M_P, s_n \models \varphi$ , and for all  $i < n$ ,  
 $M_P, s_i \models \psi$
- $M_P, s \models A(\varphi U \psi) \Leftrightarrow \forall \text{ paths } (s_0, s_1, \dots) \text{ such that:}$   
 $s_0 = s$  and  $\forall i \geq 0 (s_i, s_{i+1}) \in R$ , exists  $n \geq 0: M_P, s_n \models \varphi$ , and for all  $i < n$ ,  
 $M_P, s_i \models \psi$

Note that the formulas of CTL are evaluated in state  $s$  on a tree corresponding to an unravelling of  $M_P$  with the root  $s$ . Without loss of generality, we can assume that each model of  $P$  is a tree with the root which intuitively corresponds to the initial configuration of the agent.

### 3.3 Axiomatization

We will refer to the axiom system below as the Dribble logic of a program  $P$ ,  $DL_P$ . To simplify the axioms, we use  $guard(g)$ ,  $body(g)$ ,  $head(r)$ ,  $guard(r)$  and  $body(r)$  with the same meanings as in the model. Finally, we use  $\pi_{if}$  for *if  $\beta$  then  $\pi_1$  else  $\pi_2$* .

**CL** classical propositional logic

**CTL** axioms of CTL

**A1a**  $\bigvee_{\delta \in \mathcal{L}} \mathbf{Bs}\delta$

**A1b**  $\mathbf{Bs}\delta \rightarrow \neg \mathbf{Bs}\delta', \forall \delta' \neq \delta$   
 An agent has only one belief base.

**A2a**  $\bigvee_{\gamma \in \mathcal{L}} \mathbf{Gs}\gamma$

**A2b**  $\mathbf{Gs}\gamma \rightarrow \neg \mathbf{Gs}\gamma', \forall \gamma' \neq \gamma$   
 An agent has only one goal base.

**A3a**  $\bigvee_{\pi \in Plans \cup \{E\}} \mathbf{P}\pi$

**A3b**  $\mathbf{P}\pi \rightarrow \neg \mathbf{P}\pi', \text{ where } \pi, \pi' \in Plans \cup \{E\}, \forall \pi' \neq \pi$   
 An agent has only one plan.

**A4**  $\neg \mathbf{Bs}\delta, \forall \delta$  such that  $\delta \models_{Prop} \perp$   
 Belief base is consistent.

**A5**  $\mathbf{Bs}\delta \rightarrow \neg \mathbf{Gs}\gamma$  for all  $\gamma$  such that  $\exists g \in \gamma : \delta \models_{Prop} g$   
 All goals in goal base are not consequences of belief base.

**A6**  $\neg \mathbf{Gs}\gamma$  for all  $\gamma$  such that  $\exists g \in \gamma : g \models_{Prop} \perp$   
 Each goal in goal base is consistent.

**EPG**  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{PE} \rightarrow EX(\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}\pi)$  if  $\exists g \in \Gamma$  such that  $\langle \delta, \gamma \rangle \models_{BG} guard(g)$  and  $\pi = body(g)$

In a state  $s$  where some planning goal rule is applicable, i.e. the current plan is empty, there exists a next state  $s'$  where its plan is the one generated by firing the planning goal rule.

**APG**  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{PE} \rightarrow AX(\bigvee_{g \in \Gamma'} (\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}\pi_g))$  where  $\Gamma'$  is a set of planning goal rules  $g$  that satisfies the following two conditions:  $\langle \delta, \gamma \rangle \models_{BG} guard(g)$  and  $\pi_g = body(g)$ , provided  $\Gamma' \neq \emptyset$ .

In a state  $s$  where the current plan is empty, all possible next states from  $s$  are only reachable by applying some PG rule, i.e. its plan is generated by firing the PG rule.

**EBA**  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}(b; \pi) \rightarrow EX(\mathbf{Bs}\delta' \wedge \mathbf{Gs}\gamma' \wedge \mathbf{P}\pi)$  where  $\delta' = T(b, \delta)$  and  $\gamma' = \gamma \setminus \{g \mid \delta' \models_{Prop} g\}$

In a state  $s$  where a basic action is applicable, there exists a next state  $s'$  in which the basic action is removed from its plan, and the belief base is updated according to the basic action (the goal base, therefore, also has to be changed in order to maintain the disjointness with the belief base).

**EIF**  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}(\pi_{if}; \pi) \rightarrow EX(\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}(\pi_i; \pi))$

where

$$\pi_i = \begin{cases} \pi_1 & \text{if } \langle V_b(s), V_g(s) \rangle \models_{BG} \beta \\ \pi_2 & \text{otherwise} \end{cases}$$

In a state  $s$  where the current plan begins with a conditional plan, there exists a next state  $s'$  in which the conditional plan is replaced by one of its two sub plans depending on whether its condition is derivable or not from the belief base in  $s$ , respectively.

**EPR**  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}(\pi_1; \pi) \rightarrow EX(\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}(\pi_2; \pi))$   
if  $\exists r \in \Delta$  such that  $\langle \delta, \gamma \rangle \models_{BG} guard(r)$ ,  $head(r) = \pi_1$  and  $\pi_2 = body(r)$

In a state  $s$  where a plan revision rule is applicable, i.e. the head of the rule is the beginning of the current state and the guard of the rule is derivable from the current belief and goal base, there exists a next state  $s'$  in which the beginning of the plan in  $s$  is replaced by the body of the rule.

**ABAvPR**  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}(b; \pi'; \pi) \rightarrow AX((\mathbf{Bs}\delta' \wedge \mathbf{Gs}\gamma' \wedge \mathbf{P}(\pi'; \pi)) \vee \bigvee_{r \in \Delta'} (\mathbf{Bs}\delta' \wedge \mathbf{Gs}\gamma' \wedge \mathbf{P}(\pi''; \pi)))$  where  $\Delta'$  is a set of plan revision rules  $r$  that satisfies the following three conditions:  $head(r) = b; \pi'$ ,  $\langle \delta, \gamma \rangle \models_{BG} guard(r)$  and  $body(r) = \pi''$ , provided  $\Delta' \neq \emptyset$  or  $\mathcal{T}(\perp, \delta)$  is defined.

In a state  $s$  where a basic action  $b$  is the first element of the plan, we can only transit to another state by executing the action or applying an applicable practical reasoning rule.

**AIFvPR**  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}(\pi_{if}; \pi'; \pi) \rightarrow AX((\mathbf{Bs}\delta' \wedge \mathbf{Gs}\gamma' \wedge \mathbf{P}(\pi_i; \pi'; \pi)) \vee \bigvee_{r \in \Delta'} (\mathbf{Bs}\delta' \wedge \mathbf{Gs}\gamma' \wedge \mathbf{P}(\pi''; \pi)))$  where  $\Delta'$  is a set of plan revision rules  $r$  that satisfy three following conditions:  $head(r) = \pi_{if}; \pi'$ ,  $\langle \delta, \gamma \rangle \models_{BG} guard(r)$  and  $body(r) = \pi''$ ; and

$$\pi_i = \begin{cases} \pi_1 & \text{if } \langle V_b(s), V_g(s) \rangle \models_{BG} \beta \\ \pi_2 & \text{otherwise} \end{cases}$$

In a state  $s$  where an if-then-else statement is the first element of the plan, we can only transit to another state by executing the if-then-else statement or applying an applicable practical reasoning rule.

**APR**  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}(a; \pi_1; \pi) \rightarrow AX \bigvee_{r \in \Delta'} (\mathbf{Bs}\delta' \wedge \mathbf{Gs}\gamma' \wedge \mathbf{P}(\pi_2; \pi))$  where  $\Delta'$  is a set of plan revision rules  $r$  that satisfy three following conditions:  $head(r) = a; \pi_1$ ,  $\langle \delta, \gamma \rangle \models_{BG} guard(r)$  and  $body(r) = \pi_2$ ; provided  $\Delta' \neq \emptyset$ .

In a state  $s$  where an abstract plan is the first element of the plan, we can only transit to another state by applying a practical reasoning rule.

**Stall**  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}\pi \rightarrow AX(\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}\pi)$  where  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}\pi$  describes a configuration from which no normal transitions are available.

We have the following result.

**Theorem 1.**  $DL_P$  is sound and complete with respect to the class of models of the program  $P$ .

*Proof.* The proof of soundness is straightforward and is omitted. In the rest of this section, we show the completeness of  $DL_P$ . Most of the proof is from that of  $CTL$  [13].

Let  $BGP = 2^{\mathcal{L}} \times 2^{\mathcal{L}} \times (Plans \cup \{E\})$ .  $BGP$  intuitively corresponds to the set of all possible configurations. Note that this is a finite set.

Given a consistent formula  $\varphi_0$ , we construct the generalised Fischer-Ladner closure of  $\varphi_0$ ,  $FL(\varphi_0)$ , as the least set  $H$  of formulas containing  $\varphi_0$  such that:

1.  $\mathbf{Bs}\delta \in H$  for all  $\delta \subseteq \mathcal{L}$
2.  $\mathbf{Gs}\gamma \in H$  for all  $\gamma \subseteq \mathcal{L}$
3.  $\mathbf{P}\pi \in H$  for all  $\pi \in Plans \cup \{E\}$
4.  $EX(\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}\psi)$  for all  $(\delta, \gamma, \pi) \in BGP$
5.  $EX(\bigvee_{(\delta, \gamma, \pi) \in BGP'} (\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{P}\psi))$  for all  $BGP' \subseteq BGP$
6.  $\neg\varphi \in H$ , then  $\varphi \in H$
7.  $\varphi \wedge \psi \in H$ , then  $\varphi, \psi \in H$
8.  $E(\varphi U \psi) \in H$ , then  $\varphi, EXE(\varphi U \psi) \in H$
9.  $A(\varphi U \psi) \in H$ , then  $\varphi, AXA(\varphi U \psi) \in H$
10.  $EX\varphi \in H$ , then  $\varphi \in H$
11.  $AX\varphi \in H$ , then  $\varphi \in H$
12.  $\varphi \in H$  and  $\varphi$  is not of the form  $\neg\psi$ , then  $\neg\varphi \in H$

It is obvious that  $FL(\varphi_0)$  is finite. As usual, we define a subset  $s$  of  $FL(\varphi_0)$  that is maximally consistent if  $s$  is consistent and for all  $\varphi, \neg\varphi \in FL(\varphi_0)$ , either  $\varphi$  or  $\neg\varphi$  is in  $s$ . Repeat the construction of a model  $M$  for  $\varphi_0$  as in [13] based on the set of maximally consistent sets of  $FL(\varphi_0)$ , with the condition that the assignments are as follows:

- $V_b(s) = \delta$  for any  $\delta$  such that  $\mathbf{Bs}\delta \in s$
- $V_g(s) = \gamma$  for any  $\gamma$  such that  $\mathbf{Gs}\gamma \in s$
- $V_p(s) = \pi$  for any  $\pi$  such that  $\mathbf{P}\pi \in s$  (and  $V_p(s) = \emptyset$  if  $\mathbf{P}E \in s$ ).

The above definition is well-defined because axioms **A1x**, **A2x** and **A3x** guarantee that there are exactly one  $\mathbf{Bs}\delta \in s$ ,  $\mathbf{Gs}\gamma \in s$  and  $\mathbf{P}\pi \in s$ . Our remaining task is to show that  $M$  is, in fact, a model of  $P$ .

**EPG:** Assume that  $V_b(s) = \delta$ ,  $V_g(s) = \gamma$  and  $V_p(s) = \emptyset$ , then we have  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{PE} \in s$ . Furthermore, assume that there is  $g \in \Gamma$  such that  $\langle \delta, \gamma \rangle \models_{BG} \text{guard}(g)$ . By axiom **EPG** and modus ponens (MP),  $EX(\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{Pbody}(g)) \in s$ . According to the construction of  $M$ , there is  $s'$  such that  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{Pbody}(g) \in s'$  and  $(s, s') \in R$ . It is obvious that  $V_b(s') = \delta$ ,  $V_g(s') = \gamma$  and  $V_p(s') = \text{body}(g)$ .

**APG:** Assume that  $V_b(s) = \delta$ ,  $V_g(s) = \gamma$  and  $V_p(s) = \emptyset$ , then we have  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{PE} \in s$ . Let  $\Gamma'$  is the set of PG rules  $g \in \Gamma$  such that  $\langle \delta, \gamma \rangle \models_{BG} \text{guard}(g)$ . By axiom **APG** and modus ponens (MP),  $AX(\bigvee_{g \in \Gamma'} (\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{Pbody}(g))) \in s$ .

According to the construction of  $M$ , for any  $s'$  such that  $(s, s') \in R$ ,

$$\bigvee_{g \in \Gamma'} (\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{Pbody}(g)) \in s$$

Then, there exists  $g \in \Gamma'$  such that  $\mathbf{Bs}\delta \wedge \mathbf{Gs}\gamma \wedge \mathbf{Pbody}(g) \in s'$ . It is obvious that  $V_b(s') = \delta$ ,  $V_g(s') = \gamma$  and  $V_p(s') = \text{body}(g)$ .

The proof of other conditions on  $R$  is similar to the two cases above and are omitted. This shows that  $M$  is a model of the program  $P$ .

## 4 Verification

We can express properties of Dribble programs using CTL operators in the usual way. For example, a safety property that ‘nothing bad will happen’ can be expressed as  $AG\neg\phi$ , where  $\phi$  is a description of the ‘bad’ situation. Similarly, a liveness property that ‘something good will happen’, for example the agent will achieve its goal, can be expressed as  $EF\phi$ , where  $\phi$  is a description of the ‘good’ situation. It is essential however that we know that we are verifying the properties with respect to the computation trees which precisely correspond to the operational semantics of the agent. We prove in the next section that models of the logic correspond to computation trees, and hence that a CTL formula is true at the root of a Dribble model for  $P$  if, and only if, the corresponding property holds for the initial configuration  $c_0$  of its computation tree  $CT(c_0, P)$ .

### 4.1 Correspondence theorem

We say that a state  $s$  of some model  $M_P$  is corresponding to a configuration  $c \in CT(c_0, P)$  (notation:  $s \sim c$ ) iff  $c = \langle V_b(s), V_g(s), V_p(s) \rangle$ .

Consider a Dribble agent  $\langle \delta_0, \gamma_0, P \rangle$  and its computational tree  $CT(\langle \delta_0, \gamma_0, \emptyset \rangle, P)$ . We claim that it is isomorphic to the Dribble model of  $P$  with the root  $s_0$  such that  $s_0 \sim \langle \delta_0, \gamma_0, \emptyset \rangle$  and for every state  $s$ , all children of  $s$  are distinct. The last condition is required for isomorphism; there may be Dribble models for  $P$  where there are duplicated transitions to identical states. Such duplication of identical successors does not affect the truth of  $DL_P$  formulas.

**Theorem 2.**  *$CT(c_0, P)$  is isomorphic to the Dribble model  $M_P$  of  $P$  with the root  $s_0$  such that  $s_0 \sim c_0$  satisfying the condition that for every state  $s$ , all children of  $s$  are distinct.*

*Proof.* We are going to show that  $\sim$  defines a bijection between the states of  $CT(c_0, P)$  and  $M_P$ . We prove the theorem by induction on the distance from the root of the tree.

**Base case:** Assume that  $(s_0, s) \in R$  in  $M_P$ . We will show that there exists a unique  $c_0$  with  $c_0 \rightarrow c$  in  $CT(c_0, P)$  and  $s \sim c$ . The other direction (if  $c_0 \rightarrow c$  then there is a unique  $s$  such that  $R(s_0, s)$  and  $s \sim c$ ) is similar.

Case 1: Assume that  $V_p(s_0) = \emptyset$ , by **APG**, there is  $g \in \Gamma$  such that

$$\langle V_b(s_0), V_g(s) \rangle \models_{BG} \text{guard}(g)$$

Furthermore,  $V_b(s) = V_b(s_0)$ ,  $V_g(s) = V_g(s_0)$  and  $V_p(s) = \text{body}(g)$ . Let  $c = \langle V_b(s), V_g(s), V_p(s) \rangle$ . By the operational semantics,  $c_0 \rightarrow_{\text{apply}(g)} c$ .

Case 2: Assume that  $V_p(s_0) = b; \pi'; \pi$  ( $\pi'$  might be empty). As  $(s_0, s) \in R$ , **ABAvPR** implies that there are two cases to consider. In the first case,  $V_b(s) = T(b, V_b(s_0))$ ,  $V_g(s) = V_g(s_0) \setminus \{\alpha \in V_g(s_0) \mid V_b(s) \models_{Prop} \alpha\}$  and  $V_p(s) = \pi'; \pi$ . Simply let  $c = \langle V_b(s), V_g(s), \pi' \pi \rangle$ , we have that  $c_0 \rightarrow_{\text{execute}(b)} c$ . In the second case, we have  $V_b(s) = V_b(s_0)$ ,  $V_g(s) = V_g(s_0)$  and there is  $r \in \Delta$  such that  $\text{head}(r) = b; \pi'$  and

$$\langle V_b(s_0), V_g(s_0) \rangle \models_{BG} \text{guard}(r)$$

and  $V_p(s) = \text{body}(r); \pi$ . Let  $c = \langle V_b(s), V_g(s), \text{body}(r); \pi \rangle$ , then we have  $c_0 \rightarrow_{\text{apply}(r)} c$ .

For the other cases of  $V_p(s_0)$ , the proof is done in a similar manner by using the suitable conditions of  $R$ .

**Induction step:** Assume that the path from  $s_0$  to  $s$  has length  $n > 1$ . That means there are  $s_1, \dots, s_n = s$  in  $M_P$  such that  $(s_i, s_{i+1}) \in R$  for all  $i > 0$ . By the induction hypothesis, there are  $c_1, \dots, c_{n-1}$  such that  $s_i \sim c_i$  for all  $i = 0, \dots, n-1$  and  $c_i \rightarrow_{x_i} c_{i+1}$  by some

$$x_i \in \{\text{execute}(b), \text{execute}(if), \text{apply}(g), \text{apply}(r)\}$$

By repeating the proof of the base case, we have that there is  $c_n$  such that  $s_n \sim c_n$  and  $c_{n-1} \rightarrow_x c_n$  by some

$$x \in \{\text{execute}(b), \text{execute}(if), \text{apply}(g), \text{apply}(r)\}.$$

## 4.2 Automated verification

In this section, we show how to encode  $DL_P$  models for a standard CTL model checker to allow the automated verification of properties of Dribble programs. For the examples reported here, we have used the MOCHA model checker [3], due to the ease with which we can specify Dribble programs in *reactive modules*, the description language used by MOCHA.<sup>1</sup>

<sup>1</sup> Note that model checkers such as MCMAS [15] intended for the verification of multi-agent systems are not appropriate, as they assume that epistemic modalities are defined in terms of accessibility relations, rather than syntactically as in  $\mathcal{L}_D$ .

States of the  $DL_P$  models correspond to an assignment of values to state variables in the model checker. In particular, the agent's mental state is encoded as a collection of state variables. The agent's goals and beliefs are encoded as boolean variables with the appropriate initial values. The agent's plan is encoded as an array of *steps* of length  $len_{MAX} + 1$ . *step* is an enumeration type which includes all basic actions and abstract plans declared in the agent's program, and a set of special *if-tokens*. Each if-token  $(\beta, u, v)$  corresponds to an if-then-else construct appearing one of the agent's plans, and encodes the belief(s) tested,  $\beta$  and the lengths of the 'then' and 'else' branches (denoted by  $u$  and  $v$  respectively). All elements of the plan array are initially assigned the value *null*

The execution of plans and the application of goal and practical reasoning rules are encoded as a MOCHA *atom* which describes the initial condition and transition relation for the variables corresponding to the agent's belief and plan bases. A basic action is performed if the corresponding step token is the first element in the plan array. Executing the action updates the agent's beliefs appropriately and advances the plan, i.e., for each plan element  $i > 0$ , the step at location  $i$  is moved to location  $i - 1$ . If the first element of the plan array is an if-token, a test is performed on the appropriate belief(s) and the plan advanced accordingly. For example, if the first element of the plan array is  $(\beta, u, v)$  and  $\beta$  is false, the plan is advanced  $u + 1$  steps. Goal rules can be applied when the agent has no plan, i.e., when the first element of the plan array contains 'null', and the rule's mental condition holds (i.e., the agent has the appropriate beliefs and goals). Firing the rule writes the plan which forms the body of the goal rule into the plan array. Practical reasoning rules can be applied when the plan to be revised matches a prefix of the plan array and the rule's belief condition is believed by the agent. Firing the rule writes the plan which forms the body of the rule into the plan array and appends the suffix of the original plan (if any). In the case in which the first element of the plan array is an abstract action, application of the appropriate practical reasoning rule inserts the plan corresponding to the abstract action at the beginning of the plan array. An additional atom encodes the agent's commitment strategy, and drops any goals the agent has come to believe as a result of executing a basic action.

The evolution of the system's state is described by an initial round followed by an infinite sequence of update rounds. State variables are initialised to their initial values in the initial round and new values are assigned to the variables in the subsequent update rounds. At each update round, MOCHA non-deterministically chooses between executing the next step in the plan (if any) and firing any applicable goal and practical reasoning rules.

### 4.3 Example

As an illustration, we show how to prove properties of a simple agent program written in Dribble. Consider the following Dribble program for a simple 'vacuum cleaner' agent. The agent's environment consists of two rooms, room1 and room2, and its goal is to clean both rooms. The agent has actions which allow it to move between rooms and to clean a room, and goal rules which allow it to select an appropriate plan to clean a room. To clean a room the agent's battery must be charged and cleaning discharges the battery. The agent has a PR rule which revises a plan which is not executable because

the battery has become discharged. The agent's beliefs and goals are expressed using the following propositional variables:  $c1$ ,  $c2$  which mean that room1 and room2 are clean,  $r1$ ,  $r2$  which mean that the agent is in room1 and room2, and  $b$  means that the agent's battery is charged. The agent has the following five basic actions:<sup>2</sup>

- $mR$  for 'move right'. It is applicable if the agent's belief base  $\delta$  is such that  $\delta \models_{prop} r1$ , for example  $\delta = \{b \wedge \neg c1 \wedge \neg c2 \wedge r1 \wedge \neg r2\}$ . For this particular  $\delta$ ,  $\mathcal{T}(mR, \delta) = \{b \wedge \neg c1 \wedge \neg c2 \wedge \neg r1 \wedge r2\}$ , that is, the agent no longer believes that it is in  $r1$  but believes that it is in  $r2$ . In general,  $\mathcal{T}(mR, \delta)$  is defined as follows: for every  $\alpha \in \delta$ , in every disjunct in  $\alpha$  containing  $r1$ , replace  $r1$  with  $\neg r1$  and for every disjunct containing  $\neg r2$ , replace  $\neg r2$  with  $r2$ .
- $mL$  for 'move left'. The belief update function is defined analogously.
- $cR$  for 'clean room'. It is applicable if the agent's belief base  $\delta$  is such that  $\delta \models_{prop} b$ . If the action is executed in room1, it makes  $c1$  true, similarly for executing  $cR$  in room2. In both cases  $b$  becomes false.
- $cB$  for 'charge battery'. It is only applicable in  $r2$  (intuitively because this is where the charger is) and it makes  $b$  true.

The agent's program consists of the following two goal rules:

$$\begin{aligned} g_1 = c1 &\rightarrow \text{if } r1 \text{ then } cR \text{ else } mL; cR \\ g_2 = c2 &\rightarrow \text{if } r2 \text{ then } cR \text{ else } mR; cR \end{aligned}$$

and the PR rule:

$$r_1 = cB; \pi \mid r1 \rightarrow mR; cB; \pi$$

We would like to verify that, starting in a state in which the agent believes that it is in room1,  $r1$ , and its battery is charged,  $b$ , the above program results in the agent achieving its goals  $c1$  and  $c2$ . This can be achieved by verifying, for the corresponding model, that  $AF(Bc1 \wedge Bc2)$  is true in the initial configuration where  $b$ ,  $\neg c1$ ,  $\neg c2$ ,  $r1$  and  $\neg r2$  are true.

## 5 Related work

A logic for proving properties of Dribble agents is presented in [18], based on dynamic logic rather than on CTL. However, no axiomatisation of the logic or automated verification procedure is given. In [1, 2], Alechina et al. introduced a logic which is also based on dynamic logic for verification of agent programs written in a sublanguage of 3APL, but this work does not consider rules to revise plans. In [10] Dastani et al. prove a correspondence between an APL-like agent programming language and a specification language based on CTL, and show that any agent implemented in the language satisfies some desirable properties, e.g., relating to commitment strategies. In contrast, our aim in this paper is to verify whether a given property holds for a particular agent program. In addition, the APL considered by Dastani et al. does not include practical reasoning

<sup>2</sup> For simplicity, we omit abstract actions.

rules, and hence their results are confined to agents which are unable to revise their plans.

A strand of work on model-checking properties of agent programming languages is represented by [4] and continued by [11, 12] who use Java Path Finder (JPF) which is a model checker for Java programs. This approach requires writing a Java interpreter for the language using Agent Infrastructure Layer (AIL). After that, verification proceeds as for a standard Java program, without exploiting any features specific to agent programming languages.

A model-checking approach to automated verification of ConGolog programs was described in [8]. The paper proposes a very expressive logic which includes first-order language for specifying properties of programs and defines a model-checking algorithm for the logic. Due to the first-order sublanguage, the algorithm is not guaranteed to terminate. While ConGolog is a very expressive language, it differs from the APL and Dribble family of languages in that it lacks an explicit mechanism for revising plans.

MetateM [14] is another language for agent programming which is based on executable temporal statements. Although it is easy to verify automatically properties of agents written in MetateM, the language is very different from agent programming languages such as 3APL, AgentSpeak and Dribble, where plan constructs are based on conventional imperative languages (e.g. plans with branching constructs and loops).

## 6 Conclusion

This paper describes a temporal logic which allows us to axiomatise the set of transition systems generated by the operational semantics of a Dribble program, and formulate properties of a Dribble agent, such as that the agent is guaranteed to achieve its goals or is not going to violate some safety restrictions. One of the interesting properties of Dribble are practical reasoning or plan rewriting rules, and we believe that they pose interesting challenges for logical formalisation; in particular, we had to introduce explicit ‘plan operators’ in the language to model those rules. We show how to encode the models of the logic as input to a standard model-checker, which gives an automatic procedure for verifying properties of a Dribble program.

## References

1. N. Alechina, M. Dastani, B. Logan, and J.-J. C. Meyer. A logic of agent programs. In *Proc. of the 22nd National Conf. on Artificial Intelligence (AAAI 2007)*, pages 795–800. AAAI Press, 2007.
2. N. Alechina, M. Dastani, B. Logan, and J.-J. C. Meyer. Reasoning about agent deliberation. In *Proc. of the 11th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR’08)*, pages 16–26. AAAI, 2008.
3. R. Alur, T. A. Henzinger, F. Y. C. Mang, S. Qadeer, S. K. Rajamani, and S. Tasiran. MOCHA: Modularity in model checking. In *Computer Aided Verification*, pages 521–525, 1998.
4. R. Bordini, M. Fisher, W. Visser, and M. Wooldridge. State-space reduction techniques in agent verification. In *Proceedings of the 3rd Int. Joint Conf. on Autonomous Agents and Multi-Agent Systems (AAMAS-2004)*, pages 896–903, New York, NY, 2004. ACM Press.



5. R. H. Bordini, M. Dastani, J. Dix, and A. E. Fallah-Seghrouchni, editors. *Multi-Agent Programming: Languages, Platforms and Applications*, volume 15 of *Multiagent Systems, Artificial Societies, and Simulated Organizations*. Springer, 2005.
6. R. H. Bordini, J. F. Hübner, and R. Vieira. **Jason** and the Golden Fleece of agent-oriented programming. In *Multi-Agent Programming: Languages, Platforms and Applications*, chapter 1. Springer-Verlag, 2005.
7. E. M. Clarke and E. A. Emerson. Design and synthesis of synchronization skeletons using branching time temporal logic. In D. Kozen, editor, *IBM Logics of Programs Workshop*, number 131 in LNCS, pages 52–71. Springer-Verlag, 1981.
8. J. Claßen and G. Lakemeyer. A logic for non-terminating Golog programs. In *Proceedings of the 11th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR'08)*, pages 589–599. AAAI, 2008.
9. M. Dastani, M. B. van Riemsdijk, F. Dignum, and J.-J. C. Meyer. A programming language for cognitive agents goal directed 3apl. In *Programming Multi-Agent Systems*, LNCS, pages 111–130. Springer, 2004.
10. M. Dastani, M. B. van Riemsdijk, and J.-J. C. Meyer. A grounded specification language for agent programs. In *Proceedings of the Sixth International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS'07)*, pages 578–585, New York, , USA, 2007. ACM Press.
11. L. A. Dennis, B. Farwer, R. H. Bordini, and M. Fisher. A flexible framework for verifying agent programs. In *7th Int. Joint Conf. on Autonomous Agents and Multiagent Systems (AAMAS 2008)*, pages 1303–1306. IFAAMAS, 2008.
12. L. A. Dennis and M. Fisher. Programming verifiable heterogeneous agent systems. In *6th Int. Workshop on Programming in Multi-Agent Systems (ProMAS'08)*, pages 27–42, 2008.
13. E. A. Emerson and J. Y. Halpern. Decision procedures and expressiveness in the temporal logic of branching time. *J. Comput. Syst. Sci.*, 30(1):1–24, 1985.
14. M. Fisher. Metatem: The story so far. In *Programming Multi-Agent Systems, 3rd Int. Workshop, (ProMAS 2005)*, volume 3862 of *Lecture Notes in Computer Science*, pages 3–22. Springer, 2006.
15. A. Lomuscio and F. Raimondi. MCMAS: A model checker for multi-agent systems. In H. Hermanns and J. Palsberg, editors, *Proceedings of the 12th International Conference on Tools and Algorithms for the Construction and Analysis of Systems, (TACAS 2006)*, volume 3920 of *Lecture Notes in Computer Science*, pages 450–454. Springer, 2006.
16. D. Morley and K. Myers. The SPARK agent framework. In *Proc. of the 3rd Int. Joint Conf. on Autonomous Agents and Multiagent Systems (AAMAS'04)*, pages 714–721. IEEE Computer Society, 2004.
17. J. Thangarajah, J. Harland, D. Morley, and N. Yorke-Smith. Aborting tasks in bdi agents. In *Proc. of the 6th Int. Joint Conf. on Autonomous Agents and Multi Agent Systems (AAMAS'07)*, pages 8–15, 2007.
18. B. van Riemsdijk, W. van der Hoek, and J.-J. C. Meyer. Agent programming in dribble: from beliefs to goals using plans. In *Proc. of the 2nd Int. Joint Conf. on Autonomous Agents and Multiagent Systems (AAMAS'03)*, pages 393–400. ACM Press, 2003.

## Author Index

- Adam, Carole, 1  
Alechina, Natasha, 162
- Baldoni, Matteo, V  
Bentahar, Jamal, V, 66
- Chesani, Federico, 147  
Chopra, Amit K., 114  
Colombetti, Marco, 33
- Dastani, Mehdi M., 98  
de Boer, Mathijs, 17  
de Lima, Tiago, 17  
Di Noia, Tommaso, 82  
Di Sciascio, Eugenio, 82  
Dignum, Frank, 98  
Doan, Thu Trang, 162  
Donini, Francesco M., 82
- Fornara, Nicoletta, 33
- Gomrokchi, Maziar, 66
- Herzig, Andreas, 1, 17
- Khan, Shakil M., 49
- Khosravifar, Babak, 66
- Lespérance, Yves, 49  
Lloyd, John, V  
Logan, Brian, 162  
Longin, Dominique, 1  
Lorini, Emiliano, 17  
Louis, Vincent, 1
- Mello, Paola, 147  
Meyer, John-Jules, 98  
Montali, Marco, 147
- Ragone, Azzurra, 82
- Sakama, Chiaki, 130  
Sindlar, Michal P., 98  
Singh, Munindar P., 114
- Torrioni, Paolo, 147  
Tran Cao, Son, 130
- van Riemsdijk, M. Birna, V
- Wellman, Michael P., 82