

Peter McBurney  
Iyad Rahwan  
Simon Parsons (Eds.)

LNAI 6614

# Argumentation in Multi-Agent Systems

7th International Workshop, ArgMAS 2010  
Toronto, ON, Canada, May 2010  
Revised, Selected and Invited Papers

 Springer

Lecture Notes in Artificial Intelligence

6614

Edited by R. Goebel, J. Siekmann, and W. Wahlster

Subseries of Lecture Notes in Computer Science

Peter McBurney Iyad Rahwan  
Simon Parsons (Eds.)

# Argumentation in Multi-Agent Systems

7th International Workshop, ArgMAS 2010  
Toronto, ON, Canada, May 10, 2010  
Revised, Selected and Invited Papers



Springer

## Series Editors

Randy Goebel, University of Alberta, Edmonton, Canada  
Jörg Siekmann, University of Saarland, Saarbrücken, Germany  
Wolfgang Wahlster, DFKI and University of Saarland, Saarbrücken, Germany

## Volume Editors

Peter McBurney  
King's College London, Department of Informatics, Strand Building S6.04  
The Strand, London WC2R 2LS, UK  
E-mail: peter.mcburney@kcl.ac.uk

Iyad Rahwan  
Masdar Institute  
P.O. Box 54224, Abu Dhabi, United Arab Emirates  
E-mail: irahwan@acm.org

Simon Parsons  
City University of New York, Brooklyn College  
Department of Computer and Information Science  
2900 Bedford Avenue, Brooklyn, NY 11210, USA  
E-mail: parsons@sci.brooklyn.cuny.edu

ISSN 0302-9743 e-ISSN 1611-3349  
ISBN 978-3-642-21939-9 e-ISBN 978-3-642-21940-5  
DOI 10.1007/978-3-642-21940-5  
Springer Heidelberg Dordrecht London New York

Library of Congress Control Number: 2011929700

CR Subject Classification (1998): I.2, I.2.11, C.2.4, F.4.1, H.4, H.3

LNCS Sublibrary: SL 7 – Artificial Intelligence

© Springer-Verlag Berlin Heidelberg 2011

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer. Violations are liable to prosecution under the German Copyright Law.

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

*Typesetting:* Camera-ready by author, data conversion by Scientific Publishing Services, Chennai, India

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

# Preface

This volume contains revised versions of the papers presented at the seventh edition of the International Workshop on Argumentation in Multi-Agent Systems, (ArgMAS 2010), held in Toronto, Canada, in association with the 9th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS 2010) in May 2010. Previous ArgMAS workshops have been held in New York City, USA (2004), Utrecht, The Netherlands (2005), Hakodate, Japan (2006), Honolulu, USA (2007), Estoril, Portugal (2008) and Budapest, Hungary (2009). The event is now a regular feature on the international calendar for researchers in computational argument and dialectics in multi-agent systems. Toronto will be remembered by many as the ash-cloud workshop, as travel plans for many people were disrupted by ash from the Icelandic volcano, Eyjafjallajökull.

A brief word to explain these topics is in order. Different agents within a multi-agent system (MAS) potentially have differential access to information and different capabilities, different beliefs, different preferences and desires, and different goals. A key aspect of the scientific and engineering study of multi-agent systems therefore has been the development of methods and procedures for identifying, assessing, reconciling, arbitrating between, managing, and mitigating such differences. Market mechanisms and voting procedures are two methods for dealing with these differences. Argumentation is another. Argumentation can be understood as the formal interaction of different arguments for and against some conclusion (e.g., a proposition, an action intention, a preference, etc.). An agent may use argumentation techniques to perform individual reasoning for itself alone, in order to resolve conflicting evidence or to decide between conflicting goals it may have. Two or more agents may also jointly use dialectical argumentation to identify, express and reconcile differences between themselves, by means of interactions such as negotiation, persuasion, inquiry and joint deliberation.

In recent years, formal theories of argument and argument interaction have been proposed and studied, and this has led to the study of computational models of argument. The ArgMAS series of workshops has focused on computational argumentation within the context of agent reasoning and multi-agent systems. The ArgMAS workshops are of interest to anyone studying or applying: default reasoning in autonomous agents; single-agent reasoning and planning under uncertainty; strategic single-agent reasoning in the context of potential competitor actions; and the rational resolution of the different beliefs and intentions of multiple agents within multi-agent systems. There are close links between these topics and other topics within the discipline of autonomous agents and multi-agent systems, particularly: agent communications languages and protocols; game theory; AI planning; logic programming; and human-agent interaction.

The papers in this volume were selected for inclusion in the ArgMAS 2010 workshop following a peer-review process undertaken by anonymous reviewers,

resulting in 14 papers being selected for inclusion in the workshop. We thank all authors who made submissions to ArgMAS 2010, and we thank the members of the Program Committee listed here for their efforts in reviewing the papers submitted. We also thank the two reviewers of the paper submitted by two of the co-editors who undertook their reviews anonymously through a process of indirection, arranged and decided by the third co-editor. As for the 2009 workshop, we tasked official pre-chosen respondents to provide short, prepared critiques to a number of the papers presented at the workshop. This innovation was borrowed from conferences in philosophy, where it is standard, and we found that it works very well. The comments of respondents, who each knew of their assignment ahead of time and so could make a careful reading of their assigned paper, better focused the discussions at the workshop, and led to improvements in the quality of the revised papers later published here. This volume also contains a paper from the invited keynote speaker at the workshop, prominent argumentation-theorist and philosopher David Hitchcock of McMaster University, Hamilton, Ontario, Canada. His talk explored some of the philosophical issues behind decisions over actions, and led to a lively debate at the workshop. We were honored by Professor Hitchcock's participation, and we thank him for giving the keynote address.

As in collections of papers at previous ArgMAS workshops, we have also invited several papers from the main AAMAS Conference of relevance to argumentation in multi-agent systems. There are three invited papers here: a paper by David Grossi entitled "Argumentation in the View of Modal Logic"; a paper by Nabila Hadidi, Yannis Dimopoulos, and Pavlos Moraitis, entitled "Argumentative Alternating Offers"; and a paper by Matthias Thimm and Alejandro J. García entitled, "On Strategic Argument Selection in Structured Argumentation Systems." Apart from Professor Hitchcock's invited paper, papers in this volume are listed alphabetically by first author within three topic areas: Practical Reasoning and Reasoning About Action; Applications; and Theoretical Aspects of Argumentation.

We hope that you enjoy reading this collection.

February 2011

Peter McBurney  
Iyad Rahwan  
Simon Parsons

# Organization

## Program Chairs

Peter McBurney	King's College London, UK
Iyad Rahwan	Masdar Institute, Abu Dhabi, UAE
	MIT, USA
Simon Parsons	Brooklyn College, City University of New York, USA

## ArgMAS Steering Committee

Antonis Kakas	University of Cyprus, Cyprus
Nicolas Maudet	Université Paris Dauphine, France
Peter McBurney	King's College, London, UK
Pavlos Moraitis	Paris Descartes University, France
Simon Parsons	Brooklyn College, City University of New York, USA
Iyad Rahwan	Masdar Institute, Abu Dhabi, UAE
	MIT, USA
Chris Reed	University of Dundee, UK

## Program Committee

Leila Amgoud	IRIT, France
Katie Atkinson	University of Liverpool, UK
Jamal Bentahar	Laval University, Canada
Elizabeth Black	Oxford University, UK
Guido Boella	Università di Torino, Italy
Carlos Chesnevar	Universitat de Lleida, Spain
Frank Dignum	Utrecht University, The Netherlands
Yannis Dimopoulos	University of Cyprus, Cyprus
Sylvie Doutre	IRIT, Toulouse, France
Rogier van Eijk	Utrecht University, The Netherlands
Anthony Hunter	University College, London, UK
Antonis Kakas	University of Cyprus, Cyprus
Nikos Karacapilidis	University of Patras, Greece
Nicolas Maudet	Université Paris Dauphine, France
Peter McBurney	University of Liverpool, UK
Jarred McGinnis	Royal Holloway, University of London, UK
Sanjay Modgil	King's College London, UK

VIII Organization

Pavlos Moraitis	Paris Descartes University, France
Tim Norman	University of Aberdeen, UK
Nir Oren	King's College London, UK
Fabio Paglieri	ISTC-CNR, Rome, Italy
Simon Parsons	City University of New York, USA
Enric Plaza	Spanish Scientific Research Council, Spain
Henri Prade	IRIT, Toulouse, France
Henry Prakken	Utrecht University, The Netherlands
Iyad Rahwan	Masdar Institute, UAE
	MIT, USA
Chris Reed	University of Dundee, UK
Michael Rovatsos	University of Edinburgh, UK
Guillermo Simari	Universidad Nacional del Sur, Argentina
Francesca Toni	Imperial College, London, UK
Leon van der Torre	University of Luxembourg, Luxembourg
Paolo Torroni	Università di Bologna, Italy
Gerard Vreeswijk	Utrecht University, The Netherlands
Douglas Walton	University of Winnipeg, Canada
Simon Wells	University of Dundee, UK
Michael Wooldridge	University of Liverpool, UK



# Table of Contents

## Part I: Practical Reasoning and Argument about Action

Instrumental Rationality . . . . .	1
<i>David Hitchcock</i>	
Agreeing What to Do . . . . .	12
<i>Elizabeth Black and Katie Atkinson</i>	
A Formal Argumentation Framework for Deliberation Dialogues . . . . .	31
<i>Eric M. Kok, John-Jules Ch. Meyer, Henry Prakken, and Gerard A.W. Vreeswijk</i>	
Empirical Argumentation: Integrating Induction and Argumentation in MAS . . . . .	49
<i>Santiago Ontañón and Enric Plaza</i>	
Arguing about Preferences and Decisions . . . . .	68
<i>T.L. van der Weide, F. Dignum, J.-J.Ch. Meyer, H. Prakken, and G.A.W. Vreeswijk</i>	

## Part II: Applications

On the Benefits of Argumentation-Derived Evidence in Learning Policies . . . . .	86
<i>Chukwuemeka David Emele, Timothy J. Norman, Frank Guerin, and Simon Parsons</i>	
Argumentative Alternating Offers . . . . .	105
<i>Nabila Hadidi, Yannis Dimopoulos, and Pavlos Moraitis</i>	
On a Computational Argumentation Framework for Agent Societies . . . .	123
<i>Stella Heras, Vicente Botti, and Vicente Julián</i>	
Towards a Dialectical Approach for Conversational Agents in Selling Situations . . . . .	141
<i>Maxime Morge, Sameh Abdel-Naby, and Bruno Beaufils</i>	
Reasoning about Trust Using Argumentation: A position paper . . . . .	159
<i>Simon Parsons, Peter McBurney, and Elizabeth Sklar</i>	

**Part III: Theoretical Aspects**

An Argument-Based Multi-agent System for Information Integration ... <i>Marcela Capobianco and Guillermo R. Simari</i>	171
Argumentation in the View of Modal Logic ..... <i>Davide Grossi</i>	190
Towards Pragmatic Argumentative Agents within a Fuzzy Description Logic Framework ..... <i>Ioan Alfred Letia and Adrian Groza</i>	209
Dynamic Argumentation in Abstract Dialogue Frameworks ..... <i>M. Julieta Marcos, Marcelo A. Falappa, and Guillermo R. Simari</i>	228
Argumentation System Allowing Suspend/Resume of an Argumentation Line ..... <i>Kenichi Okuno and Kazuko Takahashi</i>	248
Computing Argumentation in Polynomial Number of BDD Operations: A Preliminary Report ..... <i>Yuqing Tang, Timothy J. Norman, and Simon Parsons</i>	268
On Strategic Argument Selection in Structured Argumentation Systems ..... <i>Matthias Thimm and Alejandro J. García</i>	286
Preference-Based Argumentation Capturing Prioritized Logic Programming ..... <i>Toshiko Wakaki</i>	306
<b>Author Index</b> .....	327

# Instrumental Rationality

David Hitchcock

Department of Philosophy, McMaster University, Hamilton, Ontario, Canada L8S 4K1  
hitchckd@mcmaster.ca

**Abstract.** Comprehensive reasoning from end to means requires an initiating intention to bring about some goal, along with five premisses: a specified means would immediately contribute to realization of the goal, the goal is achievable, the means is permissible, no alternative means is preferable, and the side effects do not outweigh the benefits of achieving the goal. Its conclusion is a decision to bring about the means. The scheme can be reiterated until an implementable means is reached. In a particular context, resource limitations may warrant truncation of the reasoning.

**ACM Category:** I.2.11 Multiagent systems. **General terms:** Theory.

**Keywords:** rationality, instrumental rationality, means-end reasoning, reasoning schemes.

## 1 Introduction

Instrumental rationality is rationality in the selection of means, or instruments, for achieving a definite goal. The goal is some state of affairs to be brought about at some future time through the agency of some person or group of persons, who need not be identical with the person or persons reasoning from end to means. A presupposition of such reasoning is that the intended end does not already obtain, and will not come about without some effort on the part of one or more agents to realize it. The means selected may be a single action by a single person, such as leaving one's home at a certain time in order to keep an appointment. But it may also be a plan, more or less completely specified at first, such as the plan of the declarer in a game of contract bridge to draw trumps first in an attempt to make the contract. Or it may be a policy, such as a policy of working out on a regular basis in order to maintain one's fitness. The goal may be a personal goal of the agent, as in the examples just mentioned. It may also be a broad social goal, like the initial target proposed by James Hansen *et al.* [1] of reducing the concentration of carbon dioxide in the Earth's atmosphere to at most 350 parts per million. The goal may be difficult to realize, with severe restrictions on available means and unreliable or incomplete information available to the reasoner or reasoners about the relevant initial conditions and causal relationships.

As pointed out in [2], reasoning from a goal in mind to a chosen means is just one form of reasoning about what is to be done, a genus often called "practical reasoning". Means-end reasoning should be distinguished, for example, from deciding to act in a certain way on the basis that the action has in itself a certain character, apart from its

consequences, as when someone notices that a store clerk has neglected to charge them for an item and decides to bring the omission to the clerk's attention, on the ground that doing so is the honourable thing to do in the situation. Here mentioning the omission is not a means to behaving honourably, but is an instance of such behaviour in concrete circumstances. The distinction between such reasoning and means-end reasoning may be difficult to draw, since as Anscombe [3] has pointed out one and the same action can have a variety of descriptions, some of which may incorporate the (expected or actual) achievement of a goal; means-end reasoning is however distinctive in involving reference to a causal chain from the selected means to the intended goal. Another form of practical reasoning is reasoning from general prescriptions or proscriptions to a conclusion about what must or cannot be done in some particular situation, as when one decides to keep silent about confidential information that one's audience has no right to know. Still other forms of practical reasoning concern the determination of what goal is to be pursued. In some cases, the goal is an intermediate goal, itself reached by a process of means-end reasoning, as in the proposal to reduce atmospheric carbon dioxide to at most 350 parts per million: *"If humanity wishes to preserve a planet similar to that on which civilization developed and to which life on Earth is adapted, paleoclimate evidence and ongoing climate change suggest that CO<sub>2</sub> will need to be reduced from its current 385 ppm to at most 350 ppm, but likely less than that."* [1] In other cases, the goal may be a final goal, not a means to achieving some further end; Richardson [4] has argued persuasively that it is possible to reason in various ways about such final ends. Another form of practical reasoning is deciding what to do on the basis of a number of relevant but separately inconclusive considerations, as when one chooses whether to spend a free evening watching a movie or reading a novel or going out for a drink with some friends. So-called "pragmatic argumentation" for or against some policy on the basis of its consequences [5] involves yet another form of practical reasoning. Still another form of practical reasoning is that envisaged by standard causal decision theory [6], in which one calculates the expected utility of one's options with a view to choosing the one with the highest expected utility.

Any decision-making about what is to be done may need to take into account a variety of types of factors: goals, prescriptions, prohibitions, valuable and "disvaluable" features, likes and dislikes, and so forth. Hence there may be considerable overlap, and even identity, between the sets of "critical questions" associated with reasoning schemes or argumentation schemes for two different types of practical reasoning. Nevertheless, it is useful to consider means-end reasoning separately from other forms of practical reasoning, because of its specific characteristic of starting from an intention to pursue a definite goal.

It is often taken to be obvious what instrumental rationality is. Habermas [7, p. 9] remarks simply that from the perspective of "cognitive-instrumental rationality" goal-directed actions are rational to the extent that their claims to effectiveness can be defended against criticism. Larry Laudan, advocating an instrumental conception of scientific rationality, writes: *"The theory of instrumental rationality simply insists that, once one has settled on one's ... desired ends, then the issue of the appropriate methods of appraisal to use depends on what strategies conduce to the realization of the selected end"* [8, p. 318].

Effectiveness of the means in securing the selected end is however often a difficult matter to determine in advance. Further, an agent may be simultaneously pursuing several goals at once, for example in conversational interaction [9]. Further, effectiveness is not always the only factor that needs to be kept in mind. As Perelman and Olbrechts-Tyteca point out ([10, p. 278]), everyday reasoning can rarely eliminate all considerations of value other than those that relate to the end in view. Hence there is more than Habermas and Laudan acknowledge to be said about instrumental rationality.

In what follows, I review the factors that may need to be taken into account when someone reasons from a concrete end in view to a means adopted with a view to achieving it, and as a result of that review propose a comprehensive scheme for means-end reasoning, whose implementation in particular domains or circumstances may be truncated, for example because of resource constraints. I focus on solo reasoning by a single agent, on the ground that such reasoning is simpler than that involved in a deliberation dialogue where two or more agents seek to arrive at an agreement on what is to be done in the pursuit of one or more antecedently agreed goals. Solo means-end reasoning is also simpler than justification of one's choice of means to a rational critic. One can of course represent solo means-end reasoning as a kind of dialogue with oneself, in which one alternately takes the role of a proponent and of a rational critic. But this representation only occasionally corresponds to the way in which solo means-end reasoning actually proceeds, and there seems to be no theoretical gain from shoe-horning solo means-end reasoning into an implicitly dialogical format. In fact, there is a theoretical risk in this approach of taking recognition that some means will achieve an agent's intended goal as establishing a presumption that the agent should perform it (cf. [11, p. 12]) — an assumption that Christian Kock [12] has cogently refuted.

## 2 Selection of the Goal

Means-end reasoning begins with the adoption as one's aim of one or more concrete ends in view. The standard belief-desire model of how reasoning issues in action, a model that comes from Aristotle (Nicomachean Ethics III.3 [13]) and Hume (Treatise II.3.3 [14]), treats the mental state of having a goal in mind as a desire. So does the more sophisticated belief-desire-intention (BDI) model due to Bratman [15]. Certainly one wants to achieve whatever one has decided to pursue as a goal. But there is more to having something as one's goal than wanting it to come about, as Searle has noted [16]. One can want something that one recognizes to be impossible, such as personal immortality on Earth, so that one makes no effort to pursue it as a goal, while nevertheless still wishing that it might come about. One can quite rationally have two desires that one recognizes cannot both be satisfied, such as the proverbial desire to have one's cake and eat it too, but one cannot rationally pursue as a goal the satisfaction of both desires once one has recognized that both cannot be satisfied. The starting-point of solo means-end reasoning might better be described as an intention to bring about an end, rather than a desire. It is not a judgment that one has the end as one's goal, and its verbal expression (to oneself or someone else) is not a statement that one has the end as one's goal. The speech act corresponding to the intention that initiates means-end reasoning would be some sort of directive, expressible linguistically by a first-person-singular imperative of the sort grammaticalized in some languages, for example classical Greek.

This proposed alternative to belief-desire and belief-desire-intention models of means-end reasoning was articulated independently of the belief-goal model proposed by Castelfranchi and Paglieri [17], to which it is similar in some respects. Castelfranchi and Paglieri conceive of a goal as “an *anticipatory internal representation* of a state of the world that has the *potential for* and the function of (eventually) *constraining/governing the behaviour of an agent towards its realization*” [17, p. 239] (italics in original). This conception is broader than the conception of a goal assumed in the present paper, in that for them a goal is not necessarily actively adopted as a constraint on action; it may merely have the potential for such constraint. In the present paper, a goal is conceived as something adopted as a concrete end in view and as actually constraining at least the agent’s thinking about what is to be done.

Intentions to pursue something as a goal are subject to rational criticism. The goal may be unattainable, so that attempts to pursue it are a waste of time and resources. Once achieved, it may turn out to be quite different than one imagined it to be, or just much less to one’s liking than one had supposed—an eventuality warned against in the saying, “Be careful what you wish for, ’cause you just might get it”, echoed in the title of cautionary lyrics by the American rapper Eminem [18]. If the goal is an intermediate goal, it can be criticized on the ground that it is ineffective for its intended purpose. It can also be criticized because it does not in fact realize the values that motivate its pursuit. Atkinson and Bench-Capon [19] have proposed to distinguish the goal pursued from the values realized by its implementation, as a way of providing for multi-agent agreement on a course of action despite differences in value preferences. Values in their approach are prized features of states of affairs, as opposed to concrete states of affairs like the examples in [17]: marrying a certain person, cooking liver Venetian style, becoming a Catholic priest, completing a dissertation, submitting an article to a journal. A distinction between goals and values is useful in solo means-end reasoning, as a way of opening up a mental space for reformulation of the goal if it seems difficult to achieve, by adopting a different goal that realizes the same value. In fact, a goal can be pursued in order to realize simultaneously a number of values. For instance, in the repressive regime in the Soviet Union from late 1982 to early 1984, a young university student was determined to lose his virginity before marriage as a form of resistance to the regime’s ideological pressures (in this case, pressure to have sex only within marriage), as well as of gaining self-respect and respect in the eyes of his peers (and sexual satisfaction); achievement of the goal would thus realize simultaneously, in his view, political, psychological and social values. Objections that achieving the goal would not in fact realize one or other of these values would count as a criticism of the intention to pursue the goal, a criticism that could be countered by taking realization of the remaining values as sufficient grounds.

The fact that adopted goals are subject to rational criticism opens up the question of the ultimate touchstone of practical reasoning, including means-end reasoning. In reasoning and argument about what is the case, the ultimate touchstone, if one adopts an epistemological rather than a purely dialectical or rhetorical perspective, is what is the case. Ideally, one’s reasons should be known to be true, and each conclusion in one’s chain of reasoning should be known to follow from the reasons offered in its immediate support, where following means that it is impossible for the reasons to be true while

the conclusion is untrue. Less stringent epistemic criteria of premiss adequacy and inference adequacy get their rationale from their aptness at tracking what is the case; for example, justified beliefs or beliefs acquired by a generally reliable process are likely to be true, and instances of inductively strong or *ceteris paribus* forms of argument tend to have a true conclusion if they have true premisses.

Is there an analogous touchstone for reasoning about what is to be done? From a purely dialectical perspective, the touchstone is acceptance by one's interlocutor of whatever starting-points and rules of inference are used to generate a conclusion about what is to be done. From a purely rhetorical perspective, the touchstone is adherence by one's intended audience to the starting-points and rules of inference. An epistemological perspective looks for some factor other than agreement or adherence. A plausible candidate is what Pollock [20] calls a "situation-liking", a feeling produced by an agent's situation as the agent believes it to be, of which humans are introspectively aware, and from which Pollock proposes to construct a cardinal measure of how much an agent likes a token situation. This cardinal measure, which has some similarities to measures of a person's utilities on the basis of their qualified preferences, can be fed into standard decision-theoretic calculations of the sort described by Weirich [6]. Pollock's proposal for the architecture of a rational agent, complex as it is, suffers from being solipsistic, asocial and amoral [21]. It might profitably be supplemented by the account of the common morality of humanity developed by Bernard Gert [22]. Gert construes morality as an informal institution for reducing the harm that human beings suffer. He defines an evil or harm as something that all rational persons avoid unless they have an adequate reason not to, and a good or benefit as something that no rational person will give up or avoid without an adequate reason [22, p. 91]. On this basis, and taking into account the types of treatment that count as punishment and the types of conditions that count as maladies for medical purposes, the basic personal evils or harms are death, pain, disability, loss of freedom and loss of pleasure; and the basic personal goods are consciousness, ability, freedom and pleasure [22, p. 108]. Gert's list of basic personal harms and basic personal benefits can be regarded as common inputs for rational human beings to the situation-likings (and situation-dislikings) that Pollock takes as fundamental to practical reasoning.

### 3 Consideration of Possible Means

However the adoption as a goal of some concrete end in view is to be critiqued or justified, and whatever the ultimate touchstone for any such critique or justification, the goal is just the starting-point of means-end reasoning. The next stage is the consideration of possible means of achieving the goal (or goals, if the reasoner aims to pursue more than one goal at once).

Two constraints on the search for effective means ought to be noted at the outset.

First, the search takes time and resources, which must be weighed against the benefits of finding some theoretically optimal path to one's goal, as compared to other desirable results from using the time and resources in a different way. Aristotle tells us that, "*if it [the end-DH] seems to be produced by several means, they [those who deliberate] consider by which it is most easily and best produced*" (Nicomachean Ethics

III.3.1112b16-17 [13]). His description has the merit of recognizing more than one criterion for choosing among possible sufficient means, not just ease or efficiency but what we might translate as “fineness”. But the cost of discovering the most efficient and finest path to one’s goal may be greater than the payoff in extra efficiency or beauty, as is commonly recognized in work on agent reasoning in computer science. As Perelman and Olbrechts Tyteca point out, “*If the value of the means is to be enhanced by the end, the means must obviously be effective; but this does not mean that it has to be the best. The determination of the best means is a technical problem, which requires various data to be brought into play and all kinds of argumentation to be used*” [10, p. 277].

Second, there are often ethical, legal or institutional constraints on acceptable means. For example, researchers designing a study to determine the effectiveness of an educational or therapeutic intervention must make sure that the design respects ethical guidelines for research using human subjects. The declarer in a game of contract bridge who works out a strategy that maximizes the chance of making the contract does so within the framework of the rules of the game, such as the rule that each player must follow suit if possible. And so on. Constraints of these sorts usually operate in the background of a person’s thinking, in the sense that the person considers only means of achieving the goal that are allowed by the constraints. Nevertheless, their operation should be acknowledged in a comprehensive account of instrumental rationality.

Perhaps the simplest case of selecting a means for achieving a goal is the case where exactly one means is required. This case seems to be the only type of means-end reasoning where something akin to the strictness of formal deductive validity comes into play. Kant expresses the underlying principle as follows: “*Whoever wills the end, also wills (insofar as reason has decisive influence on his actions) the means that are indispensably necessary to it that are in his control*” [23, p. 34, Ak4:417]. Kant maintains that this principle is an analytic necessary truth, that there is a kind of volitional inconsistency in the combination of setting out to achieve some goal, recognizing that some action in one’s power is required for the achievement of that goal, but nevertheless not proceeding to perform the required action. John Searle, despite his claim that “*there is no plausible logic of practical reason*” [16, p. 246], concedes that in one special sense Kant’s claim is correct: “*If I intend an end E, and I know that in order to achieve E I must intentionally do M, then I am committed to intending to do M*” [16, p. 266]. Searle’s formulation qualifies Kant’s claim in three ways, each needed to block counter-examples. The agent does not merely desire the end but intends it. The means is not just necessary for achieving the end but is known by the agent to be necessary. And for achievement of the end it is necessary that the agent intends to bring M about, not just that M occur.

The scope of Kant’s principle is however rather narrow, since we rarely know that we must intend to do something in order to achieve some intended goal. And the principle is a two-edged sword. One can use it either to justify implementing the necessary means or to justify abandoning or modifying one’s goal. In general, as Perelman and Olbrechts-Tyteca note, the end justifies the means only sometimes: “*the use of the means may be blameworthy in itself or have disastrous consequences outweighing the end one wishes to secure*” [10, p. 276]. In the case of a necessary means, it may also turn out that the goal will not be achieved even if one implements the means, because of other factors



beyond one's control; in that case, the reasonable thing to do is to abandon or modify the goal rather than to implement the means (unless there is some independent reason for implementing it).

If one determines that a means to one's goal is not only necessary but sufficient, that the means is achievable and permissible, that it is not in itself undesirable, that it brings with it no overriding unwelcome side-effects, and that it does not impede the pursuit of one's other goals, then one's course is clear: One should adopt the means as one's intermediate goal, and as a plan of action if one can implement it directly.

A slightly more complicated situation arises when achievement of the goal requires implementation of one of a number of means, which might for example be specifications of some generic means. Here consideration of the ease of achieving each means, its permissibility or impermissibility, its intrinsically desirable or undesirable features, the desirability or undesirability of its side-effects, and its effect on the possibility of achieving one's other goals may come into play in selecting among these disjunctively necessary means. It seems difficult to propose an algorithm or quasi-algorithm for taking such considerations into account. Walton's selection premiss in his necessary condition schema for practical reasoning is perhaps the best one can do by way of a general statement: "*I have selected one member Bi [of the set of means, at least one of which is necessary for achieving my goal–DH] as an acceptable, or the most acceptable, necessary condition for A [my goal–DH].*" [24, p. 323]; cf. [25].

A different situation arises when there are several ways of achieving the goal, each of them sufficient. It is this situation that Aristotle envisages when he describes a deliberator as selecting the easiest and best means. As indicated by previous remarks in this paper, however, it is not necessarily rational to select the easiest and best of a number of means that are each sufficient to achieve one's goal. The easiest and finest way to bring about an intended end might have foreseeable consequences whose disadvantages outweigh the benefits of achieving the goal. Or all the available means might violate a moral, legal or institutional constraint. The time and resources required to achieve the goal might not be worth it. Again, perhaps the best one can do in the way of a general statement about how to select among a number of sufficient means for achieving one's goal is Walton's selection premiss in his sufficient condition schema for practical reasoning: "*I have selected one member Bi [of the set of means, each of which is by itself sufficient for achieving my goal–DH] as an acceptable, or the most acceptable, sufficient condition for A [my goal–DH].*" [24, p. 323]; cf. [25].

In many cases, the information available does not permit identification of either a necessary or a sufficient means for achieving one's goal. One may know only the probable consequences of the options open to us, especially if those consequences depend on the actions and reactions of other agents. Perhaps less importantly, one's information about causal connections or initial conditions may be incomplete, inaccurate or even inconsistent. One may have to settle for an action that only makes it probable that one will achieve one's goal. Indeed, in some situations the most rational decision is to do something that has only a slim chance of achieving it, if it is the only possible way.

Whether a means under consideration is a necessary, sufficient, probable or even merely possible way of achieving one's goal, a number of considerations can make one hesitate before proceeding to bring about the means in question: conflicting goals,

alternative means, practical difficulties, side-effects. These considerations are well captured in the premisses and critical questions of Walton's necessary condition and sufficient condition schemata for practical reasoning [24, pp. 323-324].

Provision needs to be made, however, for the sort of backwards chaining that Aristotle describes, from an ultimate goal through intermediate means to a means that is in one's power (or in the power of an agent on whose behalf one is reasoning): "... if it [the end-DH] is achieved by one <means-DH> only they consider how it will be achieved by this and by what means this will be achieved, till they come to the first cause, which in the order of discovery is last . . ." (Nicomachean Ethics III.3.1112b17-19 [13]).

The conclusion of means-end reasoning is not a judgment that something is the case, or even a judgment that something ought to be brought about. It is a decision to bring something about, as Aristotle already recognized, or a recommendation that someone else bring it about. Its verbal expression would be some sort of directive rather than an assertive.

## 4 Conclusion

If we put together the considerations raised in the preceding discussion, we get the following rather complicated scheme for solo reasoning from a goal in mind to a selected means:

*Initiating intention of an agent A:* to bring about some goal G (where G is described as some future state of affairs, possibly but not necessarily including a reference to A)

*Immediate means premiss:* Means M1 would immediately contribute to bringing about goal G (where M1 is describable as a present or future state of affairs and may or may not be an action of A).

*Achievability premiss:* M1 is achievable as the result of a causal sequence initiated by some policy P of some agent (where the agent may or may not be A) in the present circumstances C (where achievability may be a matter of possibility or probability rather than something guaranteed).

*Permissibility premiss:* M1 violates no applicable moral, legal or institutional rule without adequate justification for the violation.

*Alternative means premiss:* No other permissible means that would immediately contribute to bringing about goal G is preferable to M1 in the light of the sum total of considerations relevant to choosing a means of bringing about an end, such as the probability in the circumstances that the means will bring about the end, the economy of time and resources involved in producing the means, the value on balance of the side effects of the means, and the intrinsic merits and demerits of the means.

*Side effects premiss:* The side effects of M1, including its effect on the achievement of other goals of A, do not outweigh the expected benefit of achieving G (where the expected benefit is a function of the values promoted by G, the degree to which achieving G will promote each of those values, and the probability that G will occur as the result of M1).

*Concluding decision:* to bring about M1.

If M1 is not a policy that an agent can immediately implement in circumstances C, then the scheme would need to be applied again, with M1 as the new goal and M2 as the hoped-for new means. Application of the scheme should continue until a means is selected that is within the power of the relevant agent.

The alternative means premiss is schematic, and would need to be fleshed out for a given practical context in a given domain. In a situation where neither of two mutually exclusive means that would contribute to achievement of the goal is preferable to the other, there is no basis for choosing one means over the other. It would be equally rational to choose either.

The scheme needs supplementation with a scheme for selection of goals, including refinement or replacement of a goal that turns out to be unachievable in an acceptable way. Castelfranchi and Paglieri [17] make some helpful suggestions in this direction, with a general characterization of belief-based goal selection, a characterization that could serve as inspiration for critical questions in various form of practical reasoning. The approach of Atkinson and Bench-Capon [19] of distinguishing goals from the values they promote could also be useful in this context.

There is also a need to supplement the generic scheme for means-end reasoning with a general framework for updating one's plans in the light of new information, as for example when the play of cards in a game of contract bridge reveals more information to the declarer about the opponents' hands.

It may not make sense to deploy the full scheme in a given situation where one has a goal in mind and needs to work out a means of achieving it. The cost of deploying the full scheme may not be worth any extra benefits so obtained. But, as pointed out by Fabio Paglieri in his review of an earlier version of this paper, such cost-benefit considerations do not diminish the analytical value of the scheme, since even simplified heuristics for decision making can be seen as abridged or modified versions of it. For instance, focusing one's attention only on a few options simply means applying the alternative means premiss to a limited sub-set of potential means or considerations relevant to the choice of such means. Adopting a satisficing perspective, as proposed by Herbert Simon [26], requires a modified version of the alternative means premiss: that no other satisficing means has been discovered that is preferable to the satisficing means M1. More generally, economies in decision making would likely involve neglecting or simplifying the alternative means premiss and/or the side effects premiss, since these are the most costly premisses in this scheme. In particular contexts, it may make sense to treat the issues of alternative means and side effects as the subject of critical questions, answers to which might overturn a presumption in favour of some means of achieving one's goal but would not be required to establish the presumption in the first place. These possible changes suggest some continuity between the present proposal of a general and idealized scheme for means-end reasoning and various bounded rationality models of the same phenomenon. Unfortunately, the present author is constrained by resource limitations to leave to others the work of exploring this continuity and the implications of the proposed scheme for work in computer science.

An abstract and high-level reasoning scheme for solo means-end reasoning like the one just proposed is perhaps not of much direct use as a guide to real-life

decision-making. It may be of most use as a guide for the formulation of lower-level domain-specific reasoning schemes. And no doubt it is subject to counter-examples that can be an occasion for further refinement.

## Acknowledgements

I would like to acknowledge helpful comments on my presentation of earlier versions of this paper at the Seventh International Workshop on Argumentation in Multi-Agent Systems (ArgMAS 2010) and at the Seventh Conference on Argumentation of the International Society for the Study of Argumentation (ISSA). At the latter conference, comments by Michael Gilbert, Hanns Hohmann, Erik Krabbe, Alain Létourneau and Christoph Lumer were particularly helpful. I would also like to acknowledge helpful reviews of the paper by Fabio Paglieri for ISSA and by two anonymous reviewers for ArgMAS 2010.

## References

1. Hansen, J., Sato, M., Kharecha, P., Beerling, D., Berner, R., Masson-Delmotte, V., Pagani, M., Raymo, M., Royer, D.L., Zachos, J.C.: Target atmospheric CO<sub>2</sub>: Where should humanity aim? *The Open Atmosphere*. *Science Journal* 2, 217–231 (2008); doi:10.2174/1874282300802010217
2. Girle, R., Hitchcock, D., McBurney, P., Verheij, B.: Decision support for Practical Reasoning. In: Reed, C., Norman, T.J. (eds.) *Argumentation Machines: New Frontiers in Argumentation and Computation*, vol. 9, pp. 55–83. Kluwer Academic Publishers, Dordrecht (2004)
3. Anscombe, E.: *Intention*, 2nd edn. Blackwell, Oxford (1963)
4. Richardson, H.S.: *Practical Reasoning about Final Ends*. Cambridge University Press, Cambridge (1994)
5. Ihnen, C.: Evaluating pragmatic argumentation: A pragma-dialectical perspective. In: Van Eemeren, F.H., Garssen, B., Godden, D., Mitchell, G. (eds.) *Proceedings of the 7th International Conference on Argumentation of the International Society for the Study of Argumentation, SICSAT, Amsterdam, The Netherlands (forthcoming)*
6. Weirich, P.: Causal decision theory. In: Zalta, E.N. (ed.) *The Stanford Encyclopedia of Philosophy* (Fall 2010), <http://plato.stanford.edu/archives/fall2010/entries/decision-causal/>
7. Habermas, J.: *The Theory of Communicative Action*, vol. 1: Reason and the Rationalization of Society, Heinemann, London, UK (1984); Translation by T. McCarthy of: *Theorie des Kommunikativen Handelns, Band I, Handlungsrationalität und gesellschaftliche Rationalisierung*, Suhrkamp, Frankfurt, Germany (1981)
8. Laudan, L.: Aim-less epistemology? *Studies in the History and Philosophy of Science* 21, 315–322 (1990)
9. Tracy, K., Coupland, N.: Multiple goals in discourse: an overview of issues. *Journal of Language and Social Psychology* 9, 1–13 (1990)
10. Perelman, C., Olbrechts-Tyteca, L.: *The New Rhetoric: A Treatise on Argumentation*. University of Notre Dame Press, Notre Dame (1969); Translation of *la Nouvelle Rhétorique. Traité de l'argumentation*. Presses Universitaires de France, Paris (1958)
11. Walton, D.N.: *Argument Schemes for Presumptive Reasoning*. Lawrence Erlbaum, Mahwah (1996)

12. Kock, C.: Is practical reasoning presumptive? *Informal Logic* 27, 91–108 (2007)
13. Aristotle: *The Complete Works of Aristotle: The Revised Oxford Translation*. vol. 1, 2, edited by J. Barnes. Princeton University Press, Princeton (1984)
14. Hume, D.: *A Treatise of Human Nature: Being an Attempt to Introduce the Experimental Method of Reasoning into Moral Subjects (1739-1740)*, <http://socserv.mcmaster.ca/econ/ugcm/3113/hume/treat.html> (retrieved)
15. Bratman, M.E.: *Intention, Plans, and Practical Reason*. Harvard University Press, Cambridge (1987); Reprinted in 1999 by CSLI, Stanford
16. Searle, J.R.: *Rationality in Action*. MIT Press, Cambridge (2001)
17. Castelfranchi, C., Paglieri, F.: The role of beliefs in goal dynamics: prolegomena to a constructive theory of intentions. *Synthese* 155, 237–263 (2007)
18. Eminem: Careful what you wish for lyrics (2008), [http://www.metrolyrics.com/careful\\_what\\_you\\_wish\\_for\\_lyrics\\_eminem.html](http://www.metrolyrics.com/careful_what_you_wish_for_lyrics_eminem.html) (retrieved)
19. Atkinson, K., Bench-Capon, T.J.M.: Practical reasoning as presumptive argumentation using action based alternating transition systems. *Artificial Intelligence* 171, 855–874 (2007); doi:10.1016/j.artint.2007.04.009
20. Pollock, J.L.: *Cognitive Carpentry: A Blueprint for How to Build a Person*. MIT Press, Cambridge (1995)
21. Hitchcock, D.: Pollock on practical reasoning. *Informal Logic* 22, 247–256 (2002)
22. Gert, B.: *Morality: Its Nature and Justification*. Revised edn. Oxford University Press, New York (2005)
23. Kant, I.: *Groundwork for the Metaphysics of Morals*. Yale University Press, New Haven (2002); Translator and Editor: A. W. Wood. Originally Published (1785, 1786)
24. Walton, D., Reed, C., Macagno, F.: *Argumentation Schemes*. Cambridge University Press, Cambridge (2008)
25. Walton, D.N.: *Practical Reasoning: Goal-Driven, Knowledge-Based, Action-Guiding Argumentation*. Rowman and Littlefield, Savage (1990)
26. Simon, H.A.: Rational choice and the structure of the environment. *Psychological Review* 63, 129–138 (1956)

# Agreeing What to Do

Elizabeth Black<sup>1</sup> and Katie Atkinson<sup>2</sup>

<sup>1</sup> Department of Engineering Science, University of Oxford, UK  
lizblack@robots.ox.ac.uk

<sup>2</sup> Department of Computer Science, University of Liverpool, UK  
katie@liverpool.ac.uk

**Abstract.** When deliberating about what to do, an autonomous agent must generate and consider the relative pros and cons of the different options. The situation becomes even more complicated when an agent is involved in a joint deliberation, as each agent will have its own preferred outcome which may change as new information is received from the other agents involved in the deliberation. We present an argumentation-based dialogue system that allows agents to come to an agreement on how to act in order to achieve a joint goal. The dialogue strategy that we define ensures that any agreement reached is acceptable to each agent, but does not necessarily demand that the agents resolve or share their differing preferences. We give properties of our system and discuss possible extensions.

**ACM Category:** I.2.11 Multiagent systems. **General terms:** Theory.

**Keywords:** dialogue, argumentation, agreement, strategy, deliberation, action.

## 1 Introduction

When agents engage in dialogues their behaviour is influenced by a number of factors including the type of dialogue taking place (e.g. negotiation or inquiry), the agents' own interests within the dialogue, and the other parties participating in the dialogue. Some of these aspects have been recognised in Walton and Krabbe's characterisation of dialogue types [1]. Some types of dialogue are more adversarial than others. For example, in a persuasion dialogue an agent may try to force its opponent to contradict itself, thus weakening the opponent's position. In a deliberation dialogue, however, the agents are more co-operative as they each share the same goal to establish agreement, although individually they may wish to influence the outcome in their own favour.

We present a dialogue system for deliberation that allows agents to reason and argue about what to do to achieve some joint goal but does not require them to pool their knowledge, nor does it require them to aggregate their preferences. Few existing dialogue systems address the problem of deliberation ([2,3] are notable exceptions). Ours is the first system for deliberation that provides a dialogue strategy that allows agents to come to an agreement about how to act that each is happy with, despite the fact that they may have different preferences and thus may each be agreeing for different reasons; it couples a dialectical setting with formal methods for argument evaluation and allows strategic manoeuvring in order to influence the dialogue outcome. We present an analysis of when agreement can and cannot be reached with our system; this provides an essential foundation to allow us to explore mechanisms that allow agents to come to an agreement in situations where the system presented here may fail.

We assume that agents are co-operative in that they do not mislead one another and will come to an agreement wherever possible; however, each agent aims to satisfy its own preferences. For the sake of simplicity, here we present a two party dialogue; however, the assumed co-operative setting means that many of the difficult issues which normally arise with multi party dialogues (e.g. [4]) are avoided here. We believe it to be straightforward to extend the system to allow multiple participants, for example following the approach taken in [5].

We describe the setting envisaged through a characteristic scenario. Consider a situation where a group of colleagues is attending a conference and they would all like to go out for dinner together. Inevitably, a deliberation takes place where options are proposed and critiqued and each individual will have his own preferences that he wishes to be satisfied by the group's decision. It is likely that there will be a range of different options proposed that are based on criteria such as: the type of cuisine desired; the proximity of the restaurant; the expense involved; the restaurant's capacity; etc.

To start the dialogue one party may put forward a particular proposal, reflecting his own preferences, say going to a French restaurant in the town centre. Such an argument may be attacked on numerous grounds, such as it being a taxi ride away, or it being expensive. If expense is a particular consideration for some members of the party, then alternative options would have to be proposed, each of which may have its own merits and disadvantages, and may need to consider the preferences already expressed. We can see that in such a scenario the agents, whilst each having their own preferred options, are committed to finding an outcome that everyone can agree to.

We present a formal argumentation-based dialogue system to handle joint deliberation. In section 2 we present the reasoning mechanism through which agents can construct and propose arguments about action. In section 3 we define the dialogue system and give an example dialogue. In section 4 we present an analysis of our system and in section 5 we discuss important extensions. In section 6 we discuss related work, and we conclude the paper in section 7.

## 2 Practical Arguments

We now describe the model of argumentation that we use to allow agents to reason about how to act. Our account is based upon a popular approach to argument characterisation, whereby argumentation schemes and critical questions are used as presumptive justification for generating arguments and attacks between them [6]. Arguments are generated by an agent instantiating a *scheme for practical reasoning* which makes explicit the following elements: the initial circumstances where action is required; the action to be taken; the new circumstances that arise through acting; the goal to be achieved; and the social value promoted by realising the goal in this way. The scheme is associated with a set of characteristic critical questions (CQs) that can be used to identify challenges to proposals for action that instantiate the scheme. An unfavourable answer to a CQ will identify a potential flaw in the argument. Since the scheme makes use of what are termed as 'values', this caters for arguments based on subjective preferences as well as more objective facts. Such values represent qualitative social interests that an agent wishes (or does not wish) to uphold by realising the goal stated [7].

To enable the practical argument scheme and critical questions approach to be precisely formalised for use in automated systems, in [8] it was defined in terms of an Action-based Alternating Transition System (AATS) [9], which is a structure for modelling game-like multi-agent systems where the agents can perform actions in order to attempt to control the system in some way. Whilst the formalisms given in [8,9] are intended to represent the overall behaviour of a multi-agent system and the effects of joint actions performed by the agents, we are interested in representing the knowledge of individual agents within a system. Hence, we use an adaptation of their formalisms (first presented in [5]) to define a *Value-based Transition System* (VATS) as follows.

**Definition 1. A Value-based Transition System (VATS)**, for an agent  $x$ , denoted  $S^x$ , is a 9-tuple  $\langle Q^x, q_0^x, Ac^x, Av^x, \rho^x, \tau^x, \Phi^x, \pi^x, \delta^x \rangle$  s.t.:

$Q^x$  is a finite set of states;

$q_0^x \in Q^x$  is the designated initial state;

$Ac^x$  is a finite set of actions;

$Av^x$  is a finite set of values;

$\rho^x : Ac^x \mapsto 2^{Q^x}$  is an action precondition function, which for each action  $a \in Ac^x$  defines the set of states  $\rho(a)$  from which  $a$  may be executed;

$\tau^x : Q^x \times Ac^x \mapsto Q^x$  is a partial system transition function, which defines the state  $\tau^x(q, a)$  that would result by the performance of  $a$  from state  $q$ —n.b. as this function is partial, not all actions are possible in all states (cf. the precondition function above);

$\Phi^x$  is a finite set of atomic propositions;

$\pi^x : Q^x \mapsto 2^{\Phi^x}$  is an interpretation function, which gives the set of primitive propositions satisfied in each state: if  $p \in \pi^x(q)$ , then this means that the propositional variable  $p$  is satisfied (equivalently, true) in state  $q$ ; and

$\delta^x : Q^x \times Q^x \times Av^x \mapsto \{+, -, =\}$  is a valuation function, which defines the status (promoted (+), demoted (-), or neutral (=)) of a value  $v \in Av^x$  ascribed by the agent to the transition between two states:  $\delta^x(q, q', v)$  labels the transition between  $q$  and  $q'$  with respect to the value  $v \in Av^x$ .

Note,  $Q^x = \emptyset \leftrightarrow Ac^x = \emptyset \leftrightarrow Av^x = \emptyset \leftrightarrow \Phi^x = \emptyset$ .

Given its VATS, an agent can now instantiate the practical reasoning argument scheme in order to construct arguments for (or against) actions to achieve a particular goal because they promote (or demote) a particular value.

**Definition 2. An argument** constructed by an agent  $x$  from its VATS  $S^x$  is a 4-tuple  $A = \langle a, p, v, s \rangle$  s.t.:  $q_x = q_0^x$ ;  $a \in Ac^x$ ;  $\tau^x(q_x, a) = q_y$ ;  $p \in \pi^x(q_y)$ ;  $v \in Av^x$ ;  $\delta^x(q_x, q_y, v) = s$  where  $s \in \{+, -\}$ .

We define the functions:  $Act(A) = a$ ;  $Goal(A) = p$ ;  $Val(A) = v$ ;  $Sign(A) = s$ .

If  $Sign(A) = +$  (–resp.), then we say  $A$  is an argument **for** (**against** resp.) action  $a$ .

We denote the **set of all arguments an agent  $x$  can construct from  $S^x$**  as  $Args^x$ ; we let  $Args_p^x = \{A \in Args^x \mid Goal(A) = p\}$ .

The set of **values** for a set of arguments  $\mathcal{X}$  is defined as  $Vals(\mathcal{X}) = \{v \mid A \in \mathcal{X} \text{ and } Val(A) = v\}$ .

If we take a particular argument for an action, it is possible to generate attacks on that argument by posing the various CQs related to the practical reasoning argument scheme.



In [8], details are given of how the reasoning with the argument scheme and posing CQs is split into three stages: *problem formulation*, where the agents decide on the facts and values relevant to the particular situation under consideration; *epistemic reasoning*, where the agents determine the current situation with respect to the structure formed at the previous stage; and *action selection*, where the agents develop, and evaluate, arguments and counter arguments about what to do. Here, we assume that the agents' problem formulation and epistemic reasoning are sound and that there is no dispute between them relating to these stages; hence, we do not consider the CQs that arise in these stages. That leaves CQ5-CQ11 for consideration (as numbered in [8]):

**CQ5:** Are there alternative ways of realising the same consequences?

**CQ6:** Are there alternative ways of realising the same goal?

**CQ7:** Are there alternative ways of promoting the same value?

**CQ8:** Does doing the action have a side effect which demotes the value?

**CQ9:** Does doing the action have a side effect which demotes some other value?

**CQ10:** Does doing the action promote some other value?

**CQ11:** Does doing the action preclude some other action which would promote some other value?

We do not consider CQ5 or CQ11 further, as the focus of the dialogue is to agree to an action that achieves the *goal*; hence, the incidental consequences (CQ5) and other potentially precluded actions (CQ11) are of no interest. We focus instead on CQ6-CQ10; agents participating in a deliberation dialogue use these CQs to identify attacks on proposed arguments for action. These CQs generate a set of arguments for and against different actions to achieve a particular goal, where each argument is associated with a motivating value. To evaluate the status of these arguments we use a Value Based Argumentation Framework (VAF), introduced in [7]. A VAF is an extension of the argumentation frameworks (AF) of Dung [10]. In an AF an argument is admissible with respect to a set of arguments  $S$  if all of its attackers are attacked by some argument in  $S$ , and no argument in  $S$  attacks an argument in  $S$ . In a VAF an argument succeeds in defeating an argument it attacks only if its value is ranked as high, or higher, than the value of the argument attacked; a particular ordering of the values is characterised as an *audience*. Arguments in a VAF are admissible with respect to an audience  $A$  and a set of arguments  $S$  if they are admissible with respect to  $S$  in the AF which results from removing all the attacks which are unsuccessful given the audience  $A$ . A maximal admissible set of a VAF is known as a *preferred extension*.

Although VAFs are commonly defined abstractly, here we give an instantiation in which we define the attack relation between the arguments. Condition 1 of the following attack relation allows for CQ8 and CQ9; condition 2 allows for CQ10; condition 3 allows for CQ6 and CQ7. Note that attacks generated by condition 1 are not symmetrical, whilst those generated by conditions 2 and 3 are.

**Definition 3.** An **instantiated value-based argumentation framework (iVAF)** is defined by a tuple  $\langle \mathcal{X}, \mathcal{A} \rangle$  s.t.  $\mathcal{X}$  is a finite set of arguments and  $\mathcal{A} \subset \mathcal{X} \times \mathcal{X}$  is the **attack relation**. A pair  $(A_i, A_j) \in \mathcal{A}$  is referred to as “ $A_i$  attacks  $A_j$ ” or “ $A_j$  is attacked by  $A_i$ ”. For two arguments  $A_i = \langle a, p, v, s \rangle$ ,  $A_j = \langle a', p', v', s' \rangle \in \mathcal{X}$ ,  $(A_i, A_j) \in \mathcal{A}$  iff  $p = p'$  and either:

1.  $a = a'$ ,  $s = -$  and  $s' = +$ ; or
2.  $a = a'$ ,  $v \neq v'$  and  $s = s' = +$ ; or
3.  $a \neq a'$  and  $s = s' = +$ .

An **audience** for an agent  $x$  over the values  $V$  is a binary relation  $\mathcal{R}^x \subset V \times V$  that defines a total order over  $V$ . We say that an argument  $A_i$  is **preferred** to the argument  $A_j$  in the audience  $\mathcal{R}^x$ , denoted  $A_i \succ_x A_j$ , iff  $(\text{Val}(A_i), \text{Val}(A_j)) \in \mathcal{R}^x$ . If  $\mathcal{R}^x$  is an audience over the values  $V$  for the iVAF  $\langle \mathcal{X}, \mathcal{A} \rangle$ , then  $\text{Vals}(\mathcal{X}) \subseteq V$ .

We use the term audience here to be consistent with the literature, it does not refer to the preference of a *set* of agents; rather, we define it to represent a particular agent's preference over a set of values.

Given an iVAF and a particular agent's audience, we can determine acceptability of an argument as follows. Note that if an attack is symmetric, then an attack only succeeds in defeat if the attacker is more preferred than the argument being attacked; however, as in [7], if an attack is asymmetric, then an attack succeeds in defeat if the attacker is at least as preferred as the argument being attacked.

**Definition 4.** Let  $\mathcal{R}^x$  be an audience and let  $\langle \mathcal{X}, \mathcal{A} \rangle$  be an iVAF.

For  $(A_i, A_j) \in \mathcal{A}$  s.t.  $(A_j, A_i) \notin \mathcal{A}$ ,  $A_i$  **defeats**  $A_j$  under  $\mathcal{R}^x$  if  $A_j \not\succeq_x A_i$ .

For  $(A_i, A_j) \in \mathcal{A}$  s.t.  $(A_j, A_i) \in \mathcal{A}$ ,  $A_i$  **defeats**  $A_j$  under  $\mathcal{R}^x$  if  $A_i \succ_x A_j$ .

An argument  $A_i \in \mathcal{X}$  is **acceptable w.r.t**  $S$  under  $\mathcal{R}^x$  ( $S \subseteq \mathcal{X}$ ) if: for every  $A_j \in \mathcal{X}$  that defeats  $A_i$  under  $\mathcal{R}^x$ , there is some  $A_k \in S$  that defeats  $A_j$  under  $\mathcal{R}^x$ .

A subset  $S$  of  $\mathcal{X}$  is **conflict-free** under  $\mathcal{R}^x$  if no argument  $A_i \in S$  defeats another argument  $A_j \in S$  under  $\mathcal{R}^x$ .

A subset  $S$  of  $\mathcal{X}$  is **admissible** under  $\mathcal{R}^x$  if:  $S$  is conflict-free in  $\mathcal{R}^x$  and every  $A \in S$  is acceptable w.r.t  $S$  under  $\mathcal{R}^x$ .

A subset  $S$  of  $\mathcal{X}$  is a **preferred extension** under  $\mathcal{R}^x$  if it is a maximal admissible set under  $\mathcal{R}^x$ .

An argument  $A$  is **acceptable** in the iVAF  $\langle \mathcal{X}, \mathcal{A} \rangle$  under audience  $\mathcal{R}^x$  if there is some preferred extension containing it.

We have now defined a mechanism with which an agent can determine attacks between arguments for and against actions, and can then use an ordering over the values that motivate such arguments (its audience) in order to determine their acceptability. In the next section we define our dialogue system.

### 3 Dialogue System

The communicative acts in a dialogue are called *moves*. We assume that there are always exactly two agents (*participants*) taking part in a dialogue, each with its own identifier taken from the set  $\mathcal{I} = \{1, 2\}$ . Each participant takes it in turn to make a move to the other participant. We refer to participants using the variables  $x$  and  $\bar{x}$  such that:  $x$  is 1 if and only if  $\bar{x}$  is 2;  $x$  is 2 if and only if  $\bar{x}$  is 1.

A move in our system is of the form  $\langle \text{Agent}, \text{Act}, \text{Content} \rangle$ . *Agent* is the identifier of the agent generating the move, *Act* is the type of move, and the *Content* gives

**Table 1.** Format for moves used in deliberation dialogues:  $\gamma$  is a goal;  $a$  is an action;  $A$  is an argument;  $x \in \{1, 2\}$  is an agent identifier

Move	Format
<i>open</i>	$\langle x, \text{open}, \gamma \rangle$
<i>assert</i>	$\langle x, \text{assert}, A \rangle$
<i>agree</i>	$\langle x, \text{agree}, a \rangle$
<i>close</i>	$\langle x, \text{close}, \gamma \rangle$

the details of the move. The format for moves used in deliberation dialogues is shown in Table 1, and the set of all moves meeting the format defined in Table 1 is denoted  $\mathcal{M}$ . Note that the system allows for other types of dialogues to be generated and these might require the addition of extra moves. Also,  $\text{Sender} : \mathcal{M} \mapsto \mathcal{I}$  is a function such that  $\text{Sender}(\langle \text{Agent}, \text{Act}, \text{Content} \rangle) = \text{Agent}$ .

We now informally explain the different types of move: an *open* move  $\langle x, \text{open}, \gamma \rangle$  opens a dialogue to agree on an action to achieve the goal  $\gamma$ ; an *assert* move  $\langle x, \text{assert}, A \rangle$  asserts an argument  $A$  for or against an action to achieve a goal that is the topic of the dialogue; an *agree* move  $\langle x, \text{agree}, a \rangle$  indicates that  $x$  agrees to performing action  $a$  to achieve the topic; a *close* move  $\langle x, \text{close}, \gamma \rangle$  indicates that  $x$  wishes to end the dialogue.

A dialogue is simply a sequence of moves, each of which is made from one participant to the other. As a dialogue progresses over time, we denote each timepoint by a natural number. Each move is indexed by the timepoint when the move was made. Exactly one move is made at each timepoint.

**Definition 5.** A **dialogue**, denoted  $D^t$ , is a sequence of moves  $[m_1, \dots, m_t]$  involving two participants in  $\mathcal{I} = \{1, 2\}$ , where  $t \in \mathbb{N}$  and the following conditions hold:

1.  $m_1$  is a move of the form  $\langle x, \text{open}, \gamma \rangle$  where  $x \in \mathcal{I}$
2.  $\text{Sender}(m_s) \in \mathcal{I}$  for  $1 \leq s \leq t$
3.  $\text{Sender}(m_s) \neq \text{Sender}(m_{s+1})$  for  $1 \leq s < t$

The **topic** of the dialogue  $D^t$  is returned by  $\text{Topic}(D^t) = \gamma$ . The set of all dialogues is denoted  $\mathcal{D}$ .

The first move of a dialogue  $D^t$  must always be an open move (condition 1 of the previous definition), every move of the dialogue must be made by a participant (condition 2), and the agents take it in turns to send moves (condition 3). In order to terminate a dialogue, either: two close moves must appear one immediately after the other in the sequence (a *matched-close*); or two moves agreeing to the same action must appear one immediately after the other in the sequence (an *agreed-close*).

**Definition 6.** Let  $D^t$  be a dialogue s.t.  $\text{Topic}(D^t) = \gamma$ . We say that  $m_s$  ( $1 < s \leq t$ ), is

- a **matched-close for**  $D^t$  iff  $m_{s-1} = \langle x, \text{close}, \gamma \rangle$  and  $m_s = \langle \bar{x}, \text{close}, \gamma \rangle$ .
- an **agreed-close for**  $D^t$  iff  $m_{s-1} = \langle x, \text{agree}, a \rangle$  and  $m_s = \langle \bar{x}, \text{agree}, a \rangle$ .

We say  $D^t$  has a **failed outcome** iff  $m_t$  is a *matched-close*, whereas we say  $D^t$  has a **successful outcome** of  $a$  iff  $m_t = \langle x, \text{agree}, a \rangle$  is an *agreed-close*.

So a *matched-close* or an *agreed-close* will terminate a dialogue  $D^t$  but only if  $D^t$  has not already terminated.

**Definition 7.** Let  $D^t$  be a dialogue.  $D^t$  **terminates at**  $t$  iff  $m_t$  is a matched-close or an agreed-close for  $D^t$  and  $\neg\exists s$  s.t.  $s < t$ ,  $D^t$  **extends**  $D^s$  (i.e. the first  $s$  moves of  $D^t$  are the same as the sequence  $D^s$ ) and  $D^s$  terminates at  $s$ .

We shortly give the particular protocol and strategy functions that allow agents to generate deliberation dialogues. First, we introduce some subsidiary definitions. At any point in a dialogue, an agent  $x$  can construct an iVAF from the union of the arguments it can construct from its VATS and the arguments that have been asserted by the other agent; we call this  $x$ 's *dialogue iVAF*.

**Definition 8.** A **dialogue iVAF** for an agent  $x$  participating in a dialogue  $D^t$  is denoted  $dVAF(x, D^t)$ . If  $D^t$  is the sequence of moves  $= [m_1, \dots, m_t]$ , then  $dVAF(x, D^t)$  is the iVAF  $\langle \mathcal{X}, \mathcal{A} \rangle$  where  $\mathcal{X} = \text{Args}_{\text{Topic}(D^t)}^x \cup \{A \mid \exists m_k = \langle \bar{x}, \text{assert}, A \rangle (1 \leq k \leq t)\}$ .

An action is *agreeable* to an agent  $x$  if and only if there is some argument for that action that is acceptable in  $x$ 's dialogue iVAF under the audience that represents  $x$ 's preference over values. Note that the set of actions that are agreeable to an agent may change over the course of the dialogue.

**Definition 9.** An action  $a$  is **agreeable** in the iVAF  $\langle \mathcal{X}, \mathcal{A} \rangle$  under the audience  $\mathcal{R}^x$  iff  $\exists A = \langle a, \gamma, v, + \rangle \in \mathcal{X}$  s.t.  $A$  is acceptable in  $\langle \mathcal{X}, \mathcal{A} \rangle$  under  $\mathcal{R}^x$ . We denote the **set of all actions that are agreeable to an agent  $x$  participating in a dialogue  $D^t$**  as  $\text{AgActs}(x, D^t)$ , s.t.  $a \in \text{AgActs}(x, D^t)$  iff  $a$  is agreeable in  $dVAF(x, D^t)$  under  $\mathcal{R}^x$ .

A protocol is a function that returns the set of moves that are permissible for an agent to make at each point in a particular type of dialogue. Here we give a deliberation protocol. It takes the dialogue that the agents are participating in and the identifier of the agent whose turn it is to move, and returns the set of permissible moves.

**Definition 10.** The **deliberation protocol** for agent  $x$  is a function  $\text{Protocol}_x : \mathcal{D} \mapsto \wp(\mathcal{M})$ . Let  $D^t$  be a dialogue ( $1 \leq t$ ) with participants  $\{1, 2\}$  s.t.  $\text{Sender}(m_t) = \bar{x}$  and  $\text{Topic}(D^t) = \gamma$ .

$$\text{Protocol}_x(D^t) = P_x^{\text{ass}}(D^t) \cup P_x^{\text{ag}}(D^t) \cup \{\langle x, \text{close}, \gamma \rangle\}$$

where the following are sets of moves and  $x' \in \{1, 2\}$ .

$$P_x^{\text{ass}}(D^t) = \{\langle x, \text{assert}, A \rangle \mid \text{Goal}(A) = \gamma$$

**and**

$$\neg\exists m_{t'} = \langle x', \text{assert}, A \rangle (1 < t' \leq t)\}$$

$$P_x^{\text{ag}}(D^t) = \{\langle x, \text{agree}, a \rangle \mid \text{either}$$

$$(1) m_t = \langle \bar{x}, \text{agree}, a \rangle\}$$

**else**

$$(2) (\exists m_{t'} = \langle \bar{x}, \text{assert}, \langle a, \gamma, v, + \rangle \rangle (1 < t' \leq t)$$

**and**

$$(\text{if } \exists m_{t''} = \langle x, \text{agree}, a \rangle$$

**then } \exists A, m\_{t'''} = \langle x, \text{assert}, A \rangle
$$(t'' < t''' \leq t)))\}$$**

$$\text{Strategy}_x(D^t) = \begin{cases} \text{Pick}(S_x^{\text{ag}})(D^t) & \text{iff } S_x^{\text{ag}}(D^t) \neq \emptyset \\ \text{Pick}(S_x^{\text{prop}})(D^t) & \text{iff } S_x^{\text{ag}}(D^t) = \emptyset \text{ and } S_x^{\text{prop}}(D^t) \neq \emptyset \\ \text{Pick}(S_x^{\text{att}})(D^t) & \text{iff } S_x^{\text{ag}}(D^t) = S_x^{\text{prop}}(D^t) = \emptyset \text{ and } S_x^{\text{att}}(D^t) \neq \emptyset \\ \langle x, \text{close}, \text{Topic}(D^t) \rangle & \text{iff } S_x^{\text{ag}}(D^t) = S_x^{\text{prop}}(D^t) = S_x^{\text{att}}(D^t) = \emptyset \end{cases}$$

where the choices for the moves are given by the following subsidiary functions ( $x' \in \{x, \bar{x}\}$ ,  $\text{Topic}(D^t) = \gamma$ ):

$$\begin{aligned} S_x^{\text{ag}}(D^t) &= \{\langle x, \text{agree}, a \rangle \in P_x^{\text{ag}}(D^t) \mid a \in \text{AgActs}(x, D^t)\} \\ S_x^{\text{prop}}(D^t) &= \{\langle x, \text{assert}, A \rangle \in P_x^{\text{ass}}(D^t) \mid A \in \text{Args}_\gamma^x, \text{Act}(A) = a, \text{Sign}(A) = + \text{ and} \\ &\quad a \in \text{AgActs}(x, D^t)\} \\ S_x^{\text{att}}(D^t) &= \{\langle x, \text{assert}, A \rangle \in P_x^{\text{ass}}(D^t) \mid A \in \text{Args}_\gamma^x, \text{Act}(A) = a, \text{Sign}(A) = -, \\ &\quad a \notin \text{AgActs}(x, D^t) \text{ and } \exists m_{t'} = \langle x', \text{assert}, A' \rangle \\ &\quad (1 \leq t' \leq t) \text{ s.t. } \text{Act}(A') = a \text{ and } \text{Sign}(A') = +\} \end{aligned}$$

**Fig. 1.** The **strategy** function uniquely selects a move according to the following preference ordering (starting with the most preferred): an agree move (ag), a proposing assert move (prop), an attacking assert move (att), a close move (close)

The protocol states that it is permissible to assert an argument as long as that argument has not previously been asserted in the dialogue. An agent can agree to an action that has been agreed to by the other agent in the preceding move (condition 1 of  $P_x^{\text{ag}}$ ); otherwise an agent  $x$  can agree to an action that has been proposed by the other participant (condition 2 of  $P_x^{\text{ag}}$ ) as long as if  $x$  has previously agreed to that action, then  $x$  has since then asserted some new argument. This is because we want to avoid the situation where an agent keeps repeatedly agreeing to an action that the other agent will not agree to: if an agent makes a move agreeing to an action and the other agent does not wish to also agree to that action, then the first agent must introduce some new argument that may convince the second agent to agree before being able to repeat its agree move. Agents may always make a close move. Note, it is straightforward to check conformance with the protocol as it only refers to public elements of the dialogue.

We now define a *basic deliberation strategy*. It takes the dialogue  $D^t$  and returns exactly one of the permissible moves. Note, this strategy makes use of a function  $\text{Pick} : \wp(\mathcal{M}) \mapsto \mathcal{M}$ . We do not define  $\text{Pick}$  here but leave it as a parameter of our strategy (in its simplest form  $\text{Pick}$  may return an arbitrary move from the input set); hence our system could generate more than one dialogue depending on the definition of the  $\text{Pick}$  function. In future work, we plan to design particular  $\text{Pick}$  functions; for example, taking into account an agent's perception of the other participant (more in section 5).

**Definition 11.** The **basic strategy** for an agent  $x$  is a function  $\text{Strategy}_x : \mathcal{D} \mapsto \mathcal{M}$  given in Figure 1.

A *well-formed deliberation dialogue* is a dialogue that has been generated by two agents each following the basic strategy.

**Definition 12.** A **well-formed deliberation dialogue** is a dialogue  $D^t$  s.t.  $\forall t' (1 \leq t' \leq t)$ ,  $\text{Sender}(m^{t'}) = x$  iff  $\text{Strategy}_x(D^{t'-1}) = m_{t'}$

We now present a simple example. There are two participating agents ( $\{1, 2\}$ ) who have the joint goal to go out for dinner together ( $din$ ).  $Ac^1 \cup Ac^2 = \{it, ch\}$  ( $it$ : go to an Italian restaurant;  $ch$ : go to a Chinese restaurant) and  $Av^1 \cup Av^2 = \{d, e1, e2, c\}$  ( $d$ : distance to travel;  $e1$ : agent 1's enjoyment;  $e2$ : agent 2's enjoyment;  $c$ : cost). The agents' audiences are as follows.

$$\begin{aligned} d &\succ_1 e1 \succ_1 c \succ_1 e2 \\ c &\succ_2 e2 \succ_2 e1 \succ_2 d \end{aligned}$$

Agent 1 starts the dialogue.

$$m_1 = \langle 1, open, din \rangle$$

The agents' dialogue iVAFs at this opening stage in the dialogue can be seen in Figs. 2 and 3, where the nodes represent arguments and are labelled with the action that they are for (or the negation of the action that they are against) and the value that they are motivated by. The arcs represent the attack relation between arguments, and a double circle round a node means that the argument it represents is acceptable to that agent.

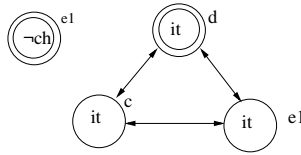


Fig. 2. Agent 1's dialogue iVAF at  $t = 1$ ,  $dVAF(1, D^1)$

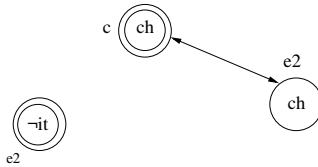


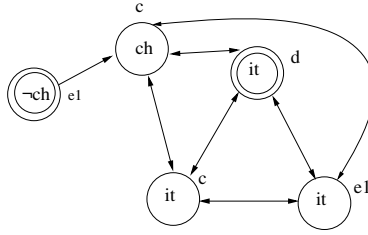
Fig. 3. Agent 2's dialogue iVAF at  $t = 1$ ,  $dVAF(2, D^1)$

At this point in the dialogue, there is only one argument for an action that is acceptable to 2 ( $\langle ch, din, c, + \rangle$ ), hence  $ch$  is the only action that is agreeable to 2. 2 must therefore assert an argument that it can construct for going to the Chinese restaurant. There are two such arguments that the Pick function could select ( $\langle ch, din, c, + \rangle$ ,  $\langle ch, din, e2, + \rangle$ ). Let us assume that  $\langle ch, din, c, + \rangle$  is selected.

$$m_2 = \langle 2, assert, \langle ch, din, c, + \rangle \rangle$$

This new argument is added to 1's dialogue iVAF, to give  $dVAF(1, D^2)$  (Fig. 4).

Although agent 2 has proposed going to the Chinese restaurant, this action is not agreeable to agent 1 at this point in the dialogue (as there is no argument for this

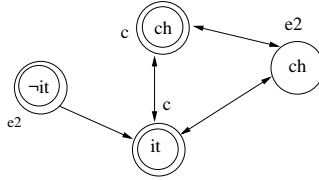


**Fig. 4.** Agent 1's dialogue iVAF at  $t = 2$ ,  $dVAF(1, D^2)$

action that is acceptable in Fig. 4). There is, however, an argument for the action *it* ( $\langle it, din, d, + \rangle$ ) that is acceptable in 1's dialogue iVAF (Fig. 4), and so going to the Italian restaurant is agreeable to 1. Hence, 1 must make an assert move proposing an argument for the action *it*, and there are three such arguments that the Pick function can select from ( $\langle it, din, d, + \rangle$ ,  $\langle it, din, c, + \rangle$ ,  $\langle it, din, e1, + \rangle$ ). Let us assume that  $\langle it, din, c, + \rangle$  is selected.

$$m_3 = \langle 1, assert, \langle it, din, c, + \rangle \rangle$$

This new argument is added to 2's dialogue iVAF, to give  $dVAF(2, D^3)$  (Fig. 5).



**Fig. 5.** Agent 2's dialogue iVAF at  $t = 3$ ,  $dVAF(2, D^3)$

Going to the Italian restaurant is now agreeable to agent 2 since the new argument introduced promotes the value ranked most highly for agent 2, i.e. cost, and so this argument is acceptable. So, 2 agrees to this action.

$$m_4 = \langle 2, agree, it \rangle$$

Going to the Italian restaurant is also agreeable to agent 1 (as the argument  $\langle it, din, d, + \rangle$  is acceptable in its dialogue iVAF, which is still the same as that shown in Fig. 4 as 2 has not asserted any new arguments), hence 1 also agrees to this action.

$$m_5 = \langle 1, agree, it \rangle$$

Note that the dialogue has terminated successfully and the agents are each happy to agree to go to the Italian restaurant; however, this action is agreeable to each agent for a different reason. Agent 1 is happy to go to the Italian restaurant as it promotes the value of distance to travel (the Italian restaurant is close by), whereas agent 2 is happy to go

to the Italian restaurant as it will promote the value of cost (as it is a cheap restaurant). The agents need not be aware of one another's audience in order to reach an agreement.

It is worth mentioning that, as we have left the Pick function unspecified, our strategy could have generated a longer dialogue if, for example, agent 1 had instead chosen to assert the argument  $\langle it, din, d, + \rangle$  at the move  $m_3$ . This illustrates how an agent's perception of the other participant may be useful: in the previous example agent 1 may make the assumption that, as agent 2 has previously asserted an argument that promotes cost, cost is something that agent 2 values; or an agent may use its perception of another agent's personality to guide argument selection [11].

Another point to note concerns the arguments generated by CQ10. Such arguments do not dispute that the action should be performed, but do dispute the reasons as to why, and so they are modelled as attacks despite being for the same action. Pinpointing this distinction here is important for two main reasons. Firstly, an advantage of the argumentation approach is that agents make explicit the reasons as to why they agree and disagree about the acceptability of arguments, and the acceptability may well turn on such reasons. Where there are two arguments proposed for the same action but each is based upon different values, an agent may only accept the argument based on one of the values. Hence such arguments are seen to be in conflict. Secondly, by participating in dialogues agents reveal what their value orderings are, as pointed out in [12]. If an agent will accept an argument for action based upon one particular value but not another, then this is potentially useful information for future dialogue interactions; if agreement is not reached about a particular action proposal, then dialogue participants will know the values an opposing agent cares about and this can guide the selection of further actions to propose, as we discuss later on in section 5.

A final related issue to note is that of accrual of arguments. If there are multiple arguments for an action and the values promoted are acceptable to the agents then some form of accrual might seem desirable. However, the complex issue of how best to accrue such arguments has not been fully resolved and this is not the focus here.

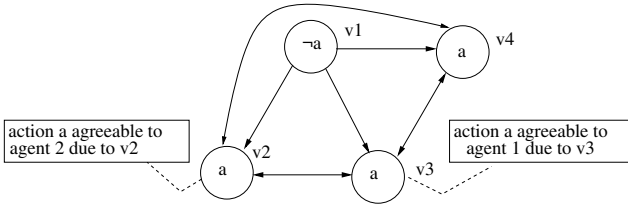
## 4 Properties

Certainly (assuming the cooperative agents do not abandon the dialogue for some reason), all dialogues generated by our system terminate. This is clear as we assume that the sets of actions and values available to an agent are finite, hence the set of arguments that an agent can construct is also finite. As the protocol does not allow the agents to keep asserting the same argument, or to keep agreeing to the same action unless a new argument has been asserted, either the dialogue will terminate successfully else the agents will run out of legal assert and agree moves and so each will make a close move.

**Proposition 1.** *If  $D^t$  is a well-formed deliberation dialogue, then  $\exists t' (t \leq t')$  s.t.  $D^{t'}$  is a well-formed deliberation dialogue that terminates at  $t'$  and  $D^{t'}$  extends  $D^t$ .*

It is also clear from the definition of the strategy (which only allows an action to be agreed to if that action is agreeable to the agent) that if the dialogue terminates with a successful outcome of action  $a$ , then  $a$  is agreeable to both agents.





**Fig. 6.** The joint iVAF

**Proposition 2.** If  $D^t$  is a well-formed deliberation dialogue that terminates successfully at  $t$  with outcome  $a$ , then  $a \in \text{AgActs}(x, D^t)$  and  $a \in \text{AgActs}(\bar{x}, D^t)$ .

Similarly, we can show that if there is an action that is agreeable to both agents when the dialogue terminates, then the dialogue will terminate successfully. In order to show this, however, we need a subsidiary lemma that states: if an agent makes a close move, then any arguments that it can construct that are for actions that it finds agreeable must have been asserted by one of the agents during the dialogue. This follows from the definition of the strategy, which only allows agents to make a close move once they have exhausted all possible assert moves.

**Lemma 1.** Let  $D^t$  be a well-formed deliberation dialogue with  $\text{Topic}(D^t) = \gamma$ , s.t.  $m_t = \langle x, \text{close}, \gamma \rangle$  and  $\text{dVAF}(x, D^t) = \langle \mathcal{X}, \mathcal{A} \rangle$ . If  $A = \langle a, \gamma, v, + \rangle \in \mathcal{X}$  and  $a \in \text{AgActs}(x, D^t)$ , then  $\exists m_{t'} = \langle x', \text{assert}, A, \rangle$  ( $1 < t' \leq t$ ,  $x' \in \{x, \bar{x}\}$ ).

Now we show that if there is an action that is agreeable to both agents when the dialogue terminates, then the dialogue will have a successful outcome.

**Proposition 3.** Let  $D^t$  be a well-formed deliberation dialogue that terminates at  $t$ . If  $a \in \text{AgActs}(x, D^t)$  and  $a \in \text{AgActs}(\bar{x}, D^t)$ , then  $D^t$  terminates successfully.

**Proof:** Assume that  $D^t$  terminates unsuccessfully at  $t$  and that  $\text{Sender}(m_t) = \bar{x}$ . From Lemma 1, there is at least one argument  $A$  for  $a$  that has been asserted by one of the agents. There are two cases. Case 1:  $x$  asserted  $A$ . Case 2:  $\bar{x}$  asserted  $A$ .

Case 1:  $x$  asserted  $A$ . Hence (from the protocol) it would have been legal for  $\bar{x}$  to make the move  $m_t = \langle \bar{x}, \text{agree}, a \rangle$  (in which case  $x$  would have had to replied with an agree, giving successful termination), unless  $\bar{x}$  had previously made a move  $m_{t'} = \langle \bar{x}, \text{agree}, a \rangle$  but had not made a move  $m_{t''} = \langle \bar{x}, \text{assert}, A \rangle$  with  $t' < t'' < t$ . However, if this were the case, then we would have  $\text{AgActs}(x, D^{t'}) = \text{AgActs}(x, D^t)$  (because no new arguments have been put forward by  $\bar{x}$  to change  $x$ 's dialogue iVAF), hence  $x$  would have had to respond to the move  $m_{t'}$  with an agree, terminating the dialogue successfully. Hence contradiction.

Case 2: works equivalently to case 1. Hence,  $D^t$  terminates successfully.  $\square$

We have shown then: all dialogues terminate; if a dialogue terminates successfully, then the outcome will be agreeable to both participants; if a dialogue terminates and there is some action that is agreeable to both agents, then the dialogue will have a successful outcome.

It would be desirable to show that if there is some action that is agreeable in the **joint iVAF**, which is the iVAF that can be constructed from the union of the agents' arguments (i.e. the iVAF  $\langle \mathcal{X}, \mathcal{A} \rangle$ , where  $\mathcal{X} = \text{Args}_\gamma^x \cup \text{Args}_\gamma^{\bar{x}}$  and  $\gamma$  is the topic of the dialogue), then the dialogue will terminate successfully. However, there are some cases where there is an action that is agreeable in the joint iVAF to each of the participants and yet still they may not reach an agreement. Consider the following example in which there is an action  $a$  that is agreeable to both the agents given the joint iVAF (see Fig. 6) and yet the dialogue generated here terminates unsuccessfully.

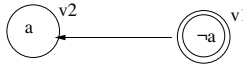
The participants ( $\{1, 2\}$ ) have the following audiences.

$$\begin{aligned} v3 \succ_1 v1 \succ_1 v4 \succ_1 v2 \\ v2 \succ_2 v1 \succ_2 v4 \succ_2 v3 \end{aligned}$$

Agent 1 starts the dialogue.

$$m_1 = \langle 1, \text{open}, p \rangle$$

The agents' dialogue iVAFs at this stage in the dialogue can be seen in Figs. 7 and 8.



**Fig. 7.** Agent 1's dialogue iVAF at  $t = 1$ ,  $dVAF(1, D^1)$



**Fig. 8.** Agent 2's dialogue iVAF at  $t = 1$ ,  $dVAF(2, D^1)$

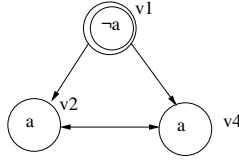
At this point in the dialogue, there is one action that is agreeable to agent 2 ( $a$ , as there is an argument *for*  $a$  that is acceptable in Fig. 8); hence (following the basic dialogue strategy), agent 2 must assert one of the arguments that it can construct for  $a$  (either  $\langle a, p, v3, + \rangle$  or  $\langle a, p, v4, + \rangle$ ). Recall, we have not specified the Pick function that has to choose between these two possible proposing assert moves. Let us assume that the Pick function makes an arbitrary choice to assert  $\langle a, p, v4, + \rangle$ .

$$m_2 = \langle 2, \text{assert}, \langle a, p, v4, + \rangle \rangle$$

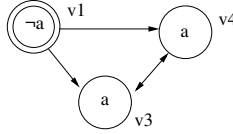
This new argument is added to agent 1's dialogue iVAF, to give  $dVAF(1, D^2)$  (Fig. 9).

From Fig. 9 we see that the only argument that is now acceptable to agent 1 is the argument *against*  $a$  ( $\langle a, p, v1, - \rangle$ ), hence there are no actions that are agreeable to agent 1. Thus agent 1 must make an attacking assert move.

$$m_3 = \langle 1, \text{assert}, \langle a, p, v1, - \rangle \rangle$$



**Fig. 9.** Agent 1's dialogue iVAF at  $t = 2$ ,  $dVAF(1, D^2)$



**Fig. 10.** Agent 2's dialogue iVAF at  $t = 3$ ,  $dVAF(2, D^3)$

This new argument is added to agent 2's dialogue iVAF, to give  $dVAF(2, D^3)$  (Fig. 10).

We see from Fig. 10 that the only argument that is now acceptable to agent 2 is the argument *against*  $a$  that 1 has just asserted ( $\langle a, p, v1, - \rangle$ ); hence,  $a$  is now no longer an agreeable action for agent 2. As there are now no actions that are agreeable to agent 2, it cannot make any proposing assert moves. It also cannot make any attacking assert moves, as the only argument that it can construct against an action has already been asserted by agent 1. Hence, agent 2 makes a close move.

$$m_4 = \langle 2, close, p \rangle$$

Thus, the dialogue iVAF for 1 is still the same as that which appears in Fig. 9. As there are no actions that are agreeable to agent 1, it cannot make any proposing assert moves. It cannot make any attacking assert moves, as the only argument that it can construct against an action has already been asserted. Hence, agent 1 also makes a close move.

$$m_5 = \langle 1, close, p \rangle$$

The dialogue has thus terminated unsuccessfully and the agents have not managed to reach an agreement as to how to achieve the goal  $p$ . However, we can see that if the Pick function instead selected the argument  $\langle a, p, v3, + \rangle$  for agent 2 to assert for the move  $m_2$ , then the resulting dialogue would have led to a successful outcome.

This example then illustrates a particular problem: the arguments exist that will enable the agents to reach an agreement (we can see this in the joint iVAF, Fig. 6 in which each agent finds  $a$  agreeable) and yet the particular arguments selected by the Pick function may not allow agreement to be reached. The choice of moves made in a deliberation dialogue affects the dialogue outcome; hence, strategic manoeuvring within the dialogue is possible in order to try to influence the dialogue outcome.

This evaluation helps us to understand the complex issues and difficulties involved in allowing agents with different preferences to agree how to act. We discuss possible responses to some of these difficulties in the next section.

## 5 Proposed Extensions

One way in which we could aim to avoid the problem illustrated in the previous example is by allowing agents to develop a model of which values they believe are important to the other participant. This model can then be used by the Pick function in order to select arguments that are more likely to lead to agreement (i.e. those that the agent believes promote or demote values that are highly preferred by the other participant). Consider the above example, if agent 2 believed that value  $v3$  was more preferred to agent 1 than value  $v4$ , then 2 would have instead asserted  $\langle a, p, v3, + \rangle$  for the move  $m_2$ , which would have led to a successful outcome.

Therefore, the first extension that we plan to investigate is to design a particular Pick function that takes into account what values the agent believes are important to the other participant. We also plan to develop a mechanism which allows the agent to build up its model of the other participant, based on the other participant's dialogue behaviour; for example, if an agent  $x$  asserts an argument for an action  $a$  because it promotes a particular value  $v$ , and the other participant  $\bar{x}$  does not then agree to  $a$ , agent  $x$  may have reason to believe that  $\bar{x}$  does not highly rank the value  $v$ .

Another problem that may be faced with our dialogue system is when it is not possible for the agents to come to an agreement no matter which arguments they choose to assert. The simplest example of this is when each agent can only construct one argument to achieve the topic  $p$ : agent 1 can construct  $\langle a1, p, v1, + \rangle$ ; agent 2 can construct  $\langle a2, p, v2, + \rangle$ . Now if agent 1's audience is such that it prefers  $v1$  to  $v2$  and agent 2's audience is such that it prefers  $v2$  to  $v1$ , then the agents will not be able to reach an agreement with the dialogue system that we have proposed here; this is despite the fact that both agents do share the goal of coming to some agreement on how to act to achieve  $p$ . The agents in this case have reached an impasse, where there is no way of finding an action that is agreeable to both agents given their individual preferences over the values.

The second extension that we propose to investigate aims to overcome such an impasse when agreement is nevertheless necessary. We plan to define a new type of dialogue (which could be embedded within the deliberation dialogue we have defined here) that allows the agents to discuss their preferences over the values and to suggest and agree to compromises that allow them to arrive at an agreement in the deliberation dialogue. For example, if agent 1's audience is  $v1 \succ_1 v2 \succ_1 v3$  and agent 2's audience is  $v3 \succ_2 v2 \succ_2 v1$ , then they may both be willing to switch their first and second most preferred values if this were to lead to an agreement (i.e. giving  $v2 \succ_1 v1 \succ_1 v3$  and  $v2 \succ_2 v3 \succ_2 v1$ ).

We would also like to extend our system to deal with the situation in which the other stages of practical reasoning (problem formulation and epistemic reasoning) may be flawed. In [5], an approach to dealing with epistemic reasoning was presented, that allowed an embedded inquiry subdialogue with which agents could jointly reason epistemically about the state of the world. Thus, the third extension that we propose is to develop a new type of dialogue that will allow agents to jointly reason about the elements of a VATS in order to consider possible flaws in the problem formulation stage.

## 6 Related Work

There is existing work in the literature on argumentation that bears some relation to what we have presented here, though the aims and contributions of these approaches are markedly different.

Our proposal follows the approach in [5][13] but the types of moves are different, and the protocol and strategy functions are substantially altered from those presented in either [5] or [13]. This alteration is necessary as neither of [5][13] allow agents to participate in deliberation dialogues. In [13], a dialogue system is presented for epistemic inquiry dialogues; it allows agents to jointly construct argument graphs (where the arguments refer only to beliefs) and to use a shared defeat relation to determine the acceptability of particular arguments.

The proposal of [5] is closer to that presented here, as both are concerned with how to act. However, the dialogue system in [5] does not allow deliberation dialogues as the outcome of any dialogue that it generates is predetermined by the union of the participating agents' knowledge. Rather, the dialogues of [5] are better categorised as a joint inference; they ensure that the agents assert all arguments that may be relevant to the question of how to act, after which a universal value ordering is applied to determine the outcome. As a shared universal value ordering is used in [5], there is an objective view of the "best" outcome (being that which you would get if you pooled the agents' knowledge and applied the shared ordering); this is in contrast to the dialogue system we present here, where the "best" outcome is subjective and depends on the point of view of a particular agent. As the agents presented here each have their own distinct audience, they must come to an explicit agreement about how to act (hence the introduction of an agree move) despite the fact that their internal views of argument acceptability may conflict. Also, here we define the attack relation (in the iVAF), which takes account of the relevant CQs, whilst in [5] the attack relation is only informally discussed.

Deliberation dialogues have only been considered in detail by the authors of [2][3]. Unlike in our work, in [2] the evaluation of arguments is not done in terms of argumentation frameworks, and strategies for reaching agreement are not considered; and in [3] the focus is on goal selection and planning.

In [12] issues concerning audiences in argumentation frameworks are addressed where the concern is to find particular audiences (if they exist) for which some arguments are acceptable and others are not. Also considered is how preferences over values emerge through a dialogue; this is demonstrated by considering how two agents can make moves within a dialogue where both are dealing with the same joint graph. However, the graph can be seen as a static structure within which agents are playing moves, i.e. putting forward acceptable arguments, rather than constructing a graph that is not complete at the outset, as in the approach we have presented.

There is also some work that considers how Dungian argumentation frameworks associated with individual agents can be merged together [14]. The merging is done not through taking the union of the individual frameworks, but through the application of criteria that determine when arguments and attacks between them can be merged into a larger graph. The main goal of the work is to characterise the sets of arguments acceptable by the whole group of agents using notions of joint acceptability, which include voting methods. In our work we are not interested in merging individual agent's graphs

*per se*; rather, an agent develops its own individual graph and uses this to determine if it finds an action agreeable. In [14] no dialogical interactions are considered, and it is also explicitly noted that consideration has not been given to how the merging approach can be applied to value-based argument systems.

Prakken [15] considers how agents can come to a public agreement despite their internal views of argument acceptability conflicting, allowing them to make explicit attack and surrender moves. However, Prakken does not explicitly consider value-based arguments, nor does he discuss particular strategies.

Strategic argumentation has been considered in other work. For example, in [16] a dialogue game for persuasion is presented that is based upon one originally proposed in [1] but makes use of Dungian argumentation frameworks. Scope is provided for three strategic considerations which concern: revealing inconsistencies between an opponent's commitments and his beliefs; exploiting the opponent's reasoning so as to create such inconsistencies; and revealing blunders to be avoided in expanding the opponent's knowledge base. These strategies all concern reasoning about an opponent's beliefs, as opposed to reasoning about action proposals with subjective preferences, as done in our work, and the game in [16] is of an adversarial nature, whereas our setting is more co-operative.

One account that does consider strategies when reasoning with value-based arguments is given in [7], where the objective is to create obligations on the opponent to accept some argument based on his previously expressed preferences. The starting point for such an interaction is a fixed joint VAF, shared by the dialogue participants. In our approach the information is not centralised in this manner, the argument graphs are built up as the dialogue proceeds, we do not assume perfect knowledge of the other agent's graph and preferences, and our dialogues have a more co-operative nature.

A related new area that is starting to receive attention is the application of game theory to argumentation (e.g. [17]). This work has investigated situations under which rational agents will not have any incentive to lie about or hide arguments; although this is concerned mainly with protocol design, it appears likely that such work will have implications for strategy design.

A few works do explicitly consider the selection of dialogue targets, that is the selection of a particular previous move to respond to. In [15] a move is defined as relevant if its target would (if attacked) cause the status of the original move to change; properties of dialogues are considered where agents are restricted to making relevant moves. In [18] this is built on to consider other classes of move relevance and the space that agents then have for strategic manoeuvring. However, these works only investigate properties of the dialogue protocols; they do not consider particular strategies for such dialogues as we do here.

## 7 Concluding Remarks

We have presented a dialogue system for joint deliberation where the agents involved in the decision making may each have different preferences yet all want an agreement to be reached. We defined how arguments and critiques are generated and evaluated, and how this is done within the context of a dialogue. A key aspect concerns how

agents' individual reasoning fits within a more global context, without the requirement to completely merge all knowledge. We presented some properties of our system that show when agreement can be guaranteed, and have explored why an agreement may not be reached. Identifying such situations is crucial for conflict resolution and we have discussed how particular steps can be taken to try to reach agreement when this occurs. In future work we intend to give a fuller account of such resolution steps whereby reasoning about other agents' preferences is central.

Ours is the first work to provide a dialogue strategy that allows agents with different preferences to come to an agreement as to how to act. The system allows strategic manoeuvring in order to influence the dialogue outcome, thus laying the important foundations needed to understand how strategy design affects dialogue outcome when the preferences involved are subjective.

## References

1. Walton, D.N., Krabbe, E.C.W.: *Commitment in Dialogue: Basic Concepts of Interpersonal Reasoning*. SUNY Press, Albany (1995)
2. McBurney, P., Hitchcock, D., Parsons, S.: The eightfold way of deliberation dialogue. *International Journal of Intelligent Systems* 22(1), 95–132 (2007)
3. Tang, Y., Parsons, S.: Argumentation-based dialogues for deliberation. In: 4th Int. Joint Conf. on Autonomous Agents and Multi-Agent Systems, pp. 552–559 (2005)
4. Dignum, F., Vreeswijk, G.: Towards a testbed for multi-party dialogues. In: AAMAS Int. Workshop on Agent Communication Languages and Conversation Policies, pp. 63–71 (2003)
5. Black, E., Atkinson, K.: Dialogues that account for different perspectives in collaborative argumentation. In: 8th Int. Joint Conf. on Autonomous Agents and Multi-Agent Systems, pp. 867–874 (2009)
6. Walton, D.N.: *Argumentation Schemes for Presumptive Reasoning*. Lawrence Erlbaum Associates, Mahwah (1996)
7. Bench-Capon, T.J.M.: Agreeing to differ: Modelling persuasive dialogue between parties without a consensus about values. *Informal Logic* 22(3), 231–245 (2002)
8. Atkinson, K., Bench-Capon, T.J.M.: Practical reasoning as presumptive argumentation using action based alternating transition systems. *Artificial Intelligence* 171(10-15), 855–874 (2007)
9. Wooldridge, M., van der Hoek, W.: On obligations and normative ability: Towards a logical analysis of the social contract. *J. of Applied Logic* 3, 396–420 (2005)
10. Dung, P.M.: On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and  $n$ -person games. *Artificial Intelligence* 77, 321–357 (1995)
11. van der Weide, T., Dignum, F., Meyer, J.-J., Prakken, H., Vreeswijk, G.: Personality-based practical reasoning. In: Rahwan, I., Moraitis, P. (eds.) *ArgMAS 2008*. LNCS, vol. 5384, pp. 3–18. Springer, Heidelberg (2009)
12. Bench-Capon, T.J.M., Doutre, S., Dunne, P.E.: Audiences in argumentation frameworks. *Artificial Intelligence* 171(1), 42–71 (2007)
13. Black, E., Hunter, A.: An inquiry dialogue system. *Autonomous Agents and Multi-Agent Systems* 19(2), 173–209 (2009)
14. Coste-Marquis, S., Devred, C., Konieczny, S., Lagasque-Schiex, M.C., Marquis, P.: On the merging of Dung's argumentation systems. *Artificial Intelligence* 171(10-15), 730–753 (2007)

15. Prakken, H.: Coherence and flexibility in dialogue games for argumentation. *J. of Logic and Computation* 15, 1009–1040 (2005)
16. Devereux, J., Reed, C.: Strategic argumentation in rigorous persuasion dialogue. In: McBurney, P., Rahwan, I., Parsons, S., Maudet, N. (eds.) *ArgMAS 2009*. LNCS, vol. 6057, pp. 94–113. Springer, Heidelberg (2010)
17. Rahwan, I., Larson, K.: Mechanism design for abstract argumentation. In: *5th Int. Joint Conf. on Autonomous Agents and Multi-Agent Systems*, pp. 1031–1038 (2008)
18. Parsons, S., McBurney, P., Sklar, E., Wooldridge, M.: On the relevance of utterances in formal inter-agent dialogues. In: *6th Int. Joint Conf. on Autonomous Agents and Multi-Agent Systems*, pp. 1002–1009 (2007)



# A Formal Argumentation Framework for Deliberation Dialogues

Eric M. Kok, John-Jules Ch. Meyer, Henry Prakken,  
and Gerard A.W. Vreeswijk

Department of Information and Computing Sciences,  
Utrecht University,  
The Netherlands

**Abstract.** Agents engage in deliberation dialogues to collectively decide on a course of action. To solve conflicts of opinion that arise, they can question claims and supply arguments. Existing models fail to capture the interplay between the provided arguments as well as successively selecting a winner from the proposals. This paper introduces a general framework for agent deliberation dialogues that uses an explicit reply structure to produce coherent dialogues, guides in outcome selection and provide pointers for agent strategies.

**Category:** I.2.11 [Artificial Intelligence] Distributed Artificial Intelligence— Languages and structures, multi-agent systems.

**General Terms:** Design, Languages.

**Keywords:** Argumentation, multi-agent communication, deliberation dialogues.

## 1 Introduction

In multi-agent systems the agents need to work together in order to achieve their personal and mutual goals. Working together means communication and often these dialogues will be on finding consensus over some belief, action or goal. In the last decade frameworks and protocols for such dialogues have been designed using argumentation theory. Walton and Krabbe [14] give a classification of dialogues types based on their initial situation, main goals and participant aims. In a *persuasion dialogue* agents need to find resolution for some conflicting point of view. They will try to persuade the others by forwarding arguments. In *negotiation*, there is not a conflict on some claim, but rather a potential conflict on the division of resources. A deal needs to be made in which each agent tries to get their most preferred resource allocation. *Deliberation* dialogues in contrast, have a significant cooperative aspect. There is a need for action and the agents need to mutually reach a decision. Although agreement is pursued, individual interests also play part.

The literature on argumentation in multi-agent systems has mainly focused on persuasion and negotiation type dialogues. Few systems for deliberation have

so far been proposed. The most sophisticated work is that of McBurney et al. [5]. To accommodate deliberating agents, a language and protocol are given that allow for forwarding and discussing of proposals for action. The protocol that they use is liberal in the sense that very few restrictions are imposed on the agents. The modelling of conflicts on beliefs and interests of the agents is limited to the assessment of commitments. It is stated that a voting phase can be used to derive a winner.

It seems that the inquisitive nature of the deliberation process has been well captured in the existing literature. However, the conflicts that arise are left more or less indeterminate. In persuasion, on the other hand, dealing with conflicts is explicitly modelled. Frameworks for these dialogues allow to determine whether given arguments are justified and consequently point out a winner. Such conflicts can be modelled this way for deliberation as well. This can be used to control the deliberation process by maintaining focus on the topic and support the selection of a winning proposal.

For persuasion dialogues, Prakken [10] has proposed a framework that uses an explicit reply structure to capture the relation between arguments. This in turn is used to ensure coherent dialogues as well as to determine the dialogical status of the initial claim. Our framework will be based on this work, adjusting it for use with deliberation dialogues. This will give several advantages. First, proposals can be assigned a status, which can be used to ensure coherent dialogues. Second, the proposed actions can be classified to guide in the selection of a winner. Moreover, the framework will be general to allow for domain specific instantiations and to capture existing protocols in it.

## 2 The Deliberation Dialogue

A deliberation dialogue commences when the need for action arises. In other words, it needs to be decided upon what action should be taken. A group of people may need to decide where to go for dinner or some automotive company needs to plan what type of car to develop. Agents will need to conceive novel proposals for action and move them in the dialogue. These proposed actions can then be reasoned upon by the agents. If a proposal is unfavourable to the agent it can question it, while it can support the proposal if it seems advantageous. Agents can even express preferences on the proposals. All this is done to influence the dialogue outcome.

In a multi-agent system, deliberation dialogues are only a part of the full communication system. Other types of dialogue, such as argument-based mutual planning [12] or persuasion, can also be part of the system. Deliberation dialogues are thus part of a context. In particular, it commences when in the context the agents believe they mutually need to decide on some action to realize a common goal. Both the goal and need for action can originate from various sources in the context, such as an authority or an earlier dialogue. When the deliberation dialogue starts, agents have, at least in our framework, already agreed on them and can start generating and evaluating proposals.

Agents will have different personal interests and beliefs, because of which conflicts of opinion will come to light during the dialogue. These conflicts can be solved by embedding persuasion-style dialogues. Agents move arguments and question claims to convince other agents. A decision on the winning proposal may be reached through agreement, a voting system or through some authority. Depending on the domain however, both the supplied arguments and the expressed preferences can still be used.

While persuasion is always competitive, deliberation is partially a cooperative process as well. This is expressed in a mutual goal that every agent needs to respect once they accept to engage in deliberation. Support for their proposals needs to show how the action will achieve this common goal. Agents thus need to mediate between their personal opinions and the mutual objective.

As an example, consider a dialogue between three agents that need to find a place for dinner where they will all enjoy the food. They all have an incentive to work towards an agreement on the restaurant, but as the dialogue progresses, differences on beliefs will also need to be resolved.

- $a_1$ : We should go to the local pizzeria.
- $a_2$ : Why should we go there? I propose we go to the nearby bistro.
- $a_1$ : Well, the pizzeria serves tasty pizza's. Why should we go to the bistro?
- $a_2$ : The toppings at the pizzeria are very dull, while the bistro has the best steaks in town.
- $a_3$ : I agree on going to the bistro, because the seafood there is great.
- $a_1$ : The bistro doesn't even server steaks any more.
- $a_3$ : What makes you think the pizza toppings are so dull?
- $a_2$ : Because the menu hasn't been changed for a very long time. We could also just go to pub.
- $a_1$ : No, I don't want to go there.

### 3 A Formal Deliberation Framework

As explained, our framework will build on the argumentation framework for persuasion dialogues of Prakken, altering and extending it for use with deliberation dialogues. It models persuasion as a dialogue game in which agents make utterances in a communication language while being restricted by a protocol. The utterances, or moves, are targeted at earlier moves. Every reply is either an attacker surrender, forming an explicit dialogue reply structure. The moves contain claims and arguments in the topic language with an argumentation logic. Since it is a framework it allows for various instantiations of the languages and protocol. In the most basic form the protocol is very liberal, only disallowing agents to speak at the same time and requiring that moves are replies to earlier moves. The dialogue terminates when one of the agents cannot make a legal move. The protocol is defined such that there are no legal moves when there is agreement on the original claim.

The explicit reply structure is utilized in two ways. First, moves have a dialectic status. The idea is that a dialogue move is *in* if it is surrendered or else all

its attackers are *out*, and that it is *out* if it has an attacker that is *in*. Now the outcome of the persuasion dialogue can be determined based on the dialogical status of the original claim, viz. if at termination this claim is *in* the proponent is the winner. Second, the protocol may be extended with a relevance rule. This compels the agents to stay focussed on the dialogue topic, giving rise to more coherent dialogues.

To make the framework suitable for deliberation dialogues, several modifications are needed. First, multiple agents need to be supported, while the persuasion framework only covers one proponent and one opponent. Several notions, such as relevance, and protocol rules, such as for termination, need to be revised accordingly. Second, there are multiple proposals instead of a single claim to discuss. The communication language needs support for forwarding, rejecting and questioning them. Multiple proposals also means there are multiple dialogical trees to which the agents may contribute. Third, the dialogue outcome is no longer a direct result of the moves. A winning function is needed to select a single action from all actions that are proposed, or possible none if there is no acceptable option.

Now the formal specification for deliberation systems in our framework is introduced. This definition is taken from [10], with the appropriate additions and revisions.

**Definition 1 (Deliberation system).** A dialogue system for deliberation dialogues is defined by:

- A *topic language*  $L_t$  is a logical language closed under classical negation.
- An *argumentation logic*  $\mathcal{L}$  as defined in [11]. It is an instance of the Dung [4] argumentation model in which arguments can be formed using inference trees of strict and defeasible rules. Here, an argument will be written as  $A \Rightarrow p$  where  $A$  is a set of premises and sub-arguments,  $\Rightarrow$  is the top inference rule and  $p$  is the conclusion of the argument. Such an argument can be attacked by rebutting the conclusion or a sub-argument, by undermining some premise it uses or by undercutting one of the used inference rules.
- A *communication language*  $L_c$ , which is a set of locutions  $\mathcal{S}$  and two binary relations  $R_a$  and  $R_s$  of attacking and surrendering reply on  $\mathcal{S}$ . Every  $s \in \mathcal{S}$  is of the form  $p(l)$  where  $p$  is a performative and  $l \in L_t$ ,  $l \subseteq L_t$  or  $l$  is an argument in  $\mathcal{L}$ .  $R_a$  and  $R_s$  are disjunct and irreflexive. Locutions cannot attack one locution and surrender to another. Finally, every surrendering locution has an *attacking counterpart*, which is an attacking locution in  $L_c$ .
- The set  $\mathcal{A}$  of agents.
- The set of *moves*  $M$  defined as  $\mathbb{N} \times \mathcal{A} \times L_c \times \mathbb{N}$  where each element of a move  $m$  respectively is denoted by:
  - $\text{id}(m)$ , the move identifier,
  - $\text{player}(m)$ , the agent that played the move,
  - $\text{content}(m)$ , the speech act, or content, of the move,
  - $\text{target}(m)$ , the move target.
- The set of *dialogues*  $M^{\leq \infty}$  is the set of all sequences  $m_1, \dots, m_i, \dots$  from  $M$ , where each  $i^{\text{th}}$  element in the sequence has identifier  $i$  and for each  $m_i$

in the sequence it holds if  $\mathbf{target}(m_i) \neq 0$  then  $\mathbf{target}(m_i) = j$  for some  $m_j$  preceding  $m_i$  in  $d$ . The set of finite dialogues  $M^{<\infty}$  is the set of all those dialogues that are finite, where one such dialogue is denoted by  $d$ .

- A *dialogue purpose* to reach a decision on a single *course of action*, which is a  $P \in L_t$ .  $P$  is a proposition stating that some action should be done.
- A *deliberation context* consisting of the *mutual goal*  $g_d \in L_t$ .
- A *protocol*  $\mathcal{P}$  that specifies the legal moves at each point in the dialogue. Formally a protocol on  $M$  is a function that works on a non-empty set of *legal finite dialogues*  $D \subseteq M^{<\infty}$  and the mutual goal such that  $\mathcal{P} : D \times L_t \rightarrow Pow(M)$ . The elements of  $\mathcal{P}(d)$  are called the legal moves after  $d$ .  $\mathcal{P}$  must satisfy the condition that for all legal finite dialogue  $d$  and moves  $m$  it holds that  $d \in D$  and  $m \in \mathcal{P}(d)$  iff  $d, m \in D$ .
- A *turntaking function*  $\mathcal{T} : D \rightarrow \mathcal{A}$  mapping a legal finite deliberation dialogue to a single agent.
- A *deliberation outcome* specified by a function  $\mathcal{O} : D \times L_t \rightarrow L_t$ , mapping all legal finite dialogues and the mutual goal  $g_d$  to a single course of action  $\alpha$ .

This deliberation system specification gives rise to a dialogue game with an explicit reply structure. The types of locutions of  $L_c$  that are available to the agents are enumerated in Table 1, each with the appropriate attacking and surrendering replies. The attacking counterpart for each surrendering locution is displayed in the same row. The locutions that deal with proposals (propose, reject, why-propose and prefer) are taken from McBurney et al. while the ones dealing with persuasion (argue, why, retract, concede) are adopted from Prakken’s framework. Below the term *proposal move* is used when the  $\mathbf{content}(m) = \mathit{propose}(P)$ , *argue move* is used when the  $\mathbf{content}(m) = \mathit{argue}(A \Rightarrow p)$ , etc.

**Table 1.** The available speech acts in the communication language  $L_c$

speech act	attacks	surrenders
$\mathit{propose}(P)$	$\mathit{why-propose}(P)$ $\mathit{reject}(P)$	
$\mathit{reject}(P)$	$\mathit{why-reject}(P)$	
$\mathit{why-propose}(P)$	$\mathit{argue}(A \Rightarrow p)$	$\mathit{drop-propose}(P)$
$\mathit{why-reject}(P)$	$\mathit{argue}(A \Rightarrow \neg p)$	$\mathit{drop-reject}(P)$
$\mathit{drop-propose}(P)$		
$\mathit{drop-reject}(P)$		
$\mathit{prefer}(P, Q)$		
$\mathit{prefer-equal}(P, Q)$		
$\mathit{skip}$		
$\mathit{argue}(A \Rightarrow p)$	$\mathit{argue}(B \Rightarrow q)$ where $B \Rightarrow q$ defeats $A \Rightarrow p$ $\mathit{why}(q)$ where $q \in A$	$\mathit{concede}(p)$ $\mathit{concede}(q)$ where $q \in A$
$\mathit{why}(p)$	$\mathit{argue}(A \Rightarrow p)$	$\mathit{retract}(p)$
$\mathit{concede}(p)$		
$\mathit{retract}(p)$		

Argue moves have a well-formed argument in  $\mathcal{L}$  as content. If it attacks some other argue move it should defeat the argument contained in that targeted move following the defeat relation of  $\mathcal{L}$ . All other speech acts have some well-formed formula in  $L_t$  as content. Note that for every move  $m$  where  $\mathbf{content} = \mathit{propose}$ ,  $\mathit{prefer}$  or  $\mathit{prefer-equal}$  it holds that  $\mathbf{target}(m) = 0$  and for all other locutions  $\mathbf{target} \neq 0$ . Specific instantiations of our framework may use a different communication language with different speech acts, as long as the reply relation is defined.

Series of moves that agents make are called turns.

**Definition 2 (Turn).** A *turn*  $T$  in a deliberation dialogue is a maximal sequence of moves  $\langle m_i, \dots, m_j \rangle$  where the same player is to move. A complete deliberation dialogue  $d$  can be split up in the sequence of turns  $\langle T_1, \dots, T_k, \dots, T_n \rangle$  where  $k \in \mathbb{N}$  is the turn identifier. A turn thus only has moves from a single player, defined by  $\mathbf{player}(T)$ .

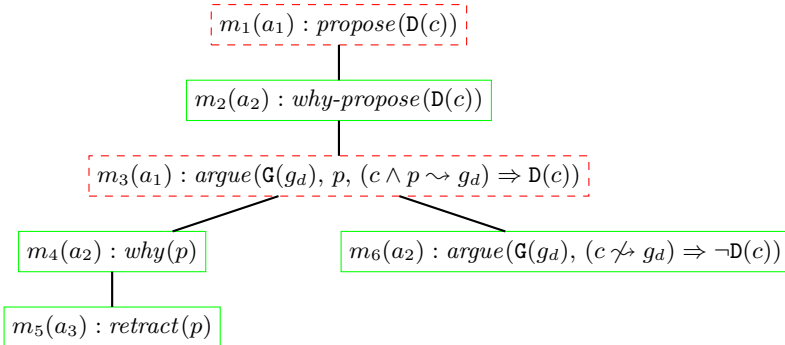
A deliberation dialogue may be represented a set of ordered directed trees.

**Definition 3 (Proposal tree).** For each proposal move  $m_i$  in dialogue  $d$  a *proposal tree*  $P$  is defined as follows:

1. The root of  $P$  is  $m_i$ .
2. For each move  $m_j$  that is a node in  $P$ , its children are all moves  $m_k$  in  $d$  such that  $\mathbf{target}(m_k) = m_j$ .

This is a tree since every move in  $d$  has a single target. Now, for any move  $m$  in proposal tree  $P$  we write  $\mathbf{proposal}(m) = m_i$ .

An example proposal tree is displayed in Fig. 1, which represents a dialogue between three agents. A proposal is moved, questioned and being supported with an argument that in turn had several replies. For each move  $m_i$  the number  $i$  is its identifier in the dialogue and between brackets the playing agent is noted. Moves in a dotted box are *out*, those in a solid box are *in*.



**Fig. 1.** A small example proposal tree

## 4 Dialogical Status of a Move

At every point in time, the dialogical status of a move can be evaluated. The use for this is twofold. First, it helps making dialogues coherent through the notion of move relevance. Secondly, the status of proposal moves can later be used during the selection of the final dialogue outcome.

Every move in a proposal tree is always either *in* or *out*. The distinction between attacking and surrendering replies is used here to make the status of moves concrete.

**Definition 4 (Move status).** A move  $m$  in a dialogue  $d$  is *in*, also called *warranted*, iff:

1.  $m$  is surrendered in  $d$  by every agent  $a \in \mathcal{A}$ ; or else,
2.  $m$  has no attacking replies in  $d$  that are *in*.

Otherwise it is *out*.

Although this definition is directly taken from [10], special attention here is required to the surrendering attacks. A move is not yet *out* until it is surrendered by every agent in the dialogue, not only by the agent that originally made the attacked move. Take for example the dialogue of Fig. 1. Although agent  $a_3$  moved a *retract*( $p$ ) in response to  $a_2$ 's *why*( $p$ ) this targeted move was still *in*. It is not until agent  $a_1$  replied with a *retract*( $p$ ) as well that the *why*( $p$ ) move is *in* again. A surrendering move is more a statement of no commitment. This idea is made concrete in the following definition of a surrendered move.

**Definition 5 (Surrendered move).** A move  $m$  is *surrendered* in a dialogue  $d$  by some agent  $a$  iff:

1.  $m$  is an argue move  $A \Rightarrow p$  and  $a$  has made a reply  $m'$  to  $m$  that has  $\text{content}(m') = \text{concede}(p)$ ; or else
2.  $a$  has made a surrendering reply to  $m$  in  $d$ .

Otherwise it is *out*.

The notion of relevance can now be formalised.

**Definition 6 (Relevance).** An attacking move  $m$  in a dialogue  $d$  is *relevant* iff it changes the move status of  $\text{proposal}(m)$ . A surrendering move is relevant iff its attacking counterpart is.

Depending on the domain a different notion of surrendered move or relevance may be useful. Prakken describes a notion of weak relevance that may be adopted. It is weaker in the sense that an agent can contribute multiple ways to change the proposal tree root and still be relevant. This is achieved by only requiring a move to create an additional way to change the status of a proposal. A protocol with weak relevance allows an agent to make multiple attacks per turn in a proposal tree as opposed to a single one if the earlier notion is used, which we below use the term *strong relevance* for.

**Definition 7 (Weak relevance).** An attacking move  $m$  in a dialogue  $d$  is *weakly relevant* iff it creates a new or removes an existing *winning part* in the proposal tree  $P$  associated with  $\text{proposal}(m)$  in  $d$ . A surrendering move is weakly relevant iff its attacking counterpart is. If the  $\text{proposal}(m)$  is *in*, a winning part  $w^P$  for this tree  $P$  is defined as follows:

1. First include the root of  $P$ ;
2. For each  $m$  of even depth, if  $m$  is surrendered by every agent  $a \in \mathcal{A}$ , include all its surrendering replies, otherwise include all its attacking replies;
3. For each  $m$  of even depth, include one attacking reply  $m'$  that is *in* in  $d$ ;

The idea of a winning part is that it is 'a reason' why the proposal is *in* at that moment. Since this is not unique, there may be alternative attacking replies, a move is already weakly relevant if it succeeds to create an additional winning part or removes a winning part. Take for example the dialogue of Fig. 1 again. After  $\text{argue}(\mathbf{G}(g_d), (c \rightsquigarrow g_d) \Rightarrow \mathbf{D}(c))$  was moved by agent  $a_1$  there are no more strongly relevant moves in this proposal tree, while there exists new weakly relevant moves, for example  $\text{argue}(s \Rightarrow g_d)$ . This results in a more liberal deliberation process.

## 5 Turntaking and Termination

We have still not made concrete how agents take turns and when the dialogue terminates.

**Definition 8 (Turntaking).** Agents take turns in sequence and end their turns explicitly with a skip move. Formally, for a dialogue  $d = \langle m_1, \dots, m_n \rangle$   $\mathcal{T}(d) = \text{player}(m_n)$  unless  $\text{content}(m_n) = \text{skip}$  in which case  $\mathcal{T}(d) = \text{player}(m_n) + 1$ .

Clearly, when there are no more legal moves besides the skip move, that is  $\mathcal{P}(d) = \{\text{skip}\}$ , the turn switches. Now, the dialogue terminates if all agents no longer make other moves than directly skipping.

**Definition 9 (Termination).** A dialogue  $d$  terminates on  $|\mathcal{A}| + 1$  consecutive skip moves.

The rationale behind the termination rule is that each agent should have the opportunity to make new moves when it stills want to. However, to prevent agents from endlessly skipping until some other agent makes a beneficial move or even a mistake, the number of skip moves is limited.

## 6 Protocol Rules

Now various protocol rules are discussed. Depending on the domain some might or might not be desirable. First, some rules that prevent agents from playing incoherent moves are added. More precisely, these rules require the agents to be relevant, not to overflow the dialogue.



1. Agents can only reply to moves of others. Formally, for every attacking or surrendering move  $m$  in a dialogue  $player(m) \neq player(target(m))$ .
2. Every attacking and surrendering move must be relevant.
3. A turn can contain at most one proposal move.
4. A proposal must be unique in the dialogue. Formally, for every proposal move  $m$  in  $d$  it holds that  $content(m) \notin \{p | p = content(n) \text{ of some proposal move } n \in d\}$ .

The first rule may be dropped for domains where a more liberal deliberation process is appropriate. This would allow agents to attack their own proposals as well. The relevance of the second rule may be strong or weak relevance. Note that in case of strong relevance there can be at most one attacking move per proposal tree.

Not only the dialogue should be coherent. The same holds for the agents' preference statements on the proposals. A protocol rule is added to ensure that an agent is consistent in his ordering.

5. An agent may only make a prefer move if the resulting option ordering maintains transitivity and antisymmetry. This is further explained below.

The last rules are used to ensure that arguments for (and against) a proposal explain how it (fails to) achieve the mutual goal.

6. Every argue move  $m$  with  $target(m) = m'$  and  $content(m') = why-propose(D(P))$  will contain an argument in  $\mathcal{L}$  with  $g_d$  as one of its premises and  $D(P)$  as conclusion.
7. Every argue move  $m$  with  $target(m) = m'$  and  $content(m') = why-reject(D(P))$  will contain an argument in  $\mathcal{L}$  with  $\neg g_d$  as one of its premises and  $\neg D(P)$  as conclusion.

The arguments that these protocol rules require are used to make sure that a proposal for action  $P$  will indeed (fail to) achieve the mutual goal  $g_d$ . Put differently, the proposed action needs to be *appropriate* in relation to our dialogue topic. The topic language and used logic therefore need support to express this. One option, used below, is to include an inference rule for the practical syllogism in our logic  $\mathcal{L}$ . Similar to [2] a practical reasoning rule will then be used that says 'if  $g_d$  is a goal and  $P$  will achieve  $g_d$  then  $P$  is an appropriate proposal for action'. Such arguments, below written as  $G(g_d), P \rightsquigarrow g_d \Rightarrow D(P)$ , can then be moved.

## 7 Dialogue Outcome

At any moment in time the outcome of the deliberation dialogue can be determined. As the outcome function dictates, this is a single course of action, or no action at all when there is a structural disagreement. To establish this, the options, which are the moved proposals, are first specified and then classified based on their status. This set of proposals is then considered over the agent preferences to determine a winner.

**Definition 10 (Options).** The dialogue *options* are defined by a function  $O : D \rightarrow Pow(L_i)$  mapping all legal dialogues to a subset of proposals. For any dialogue  $d$  the set of options is  $O(d) = \{o \mid o = \text{content}(m) \text{ for each proposal move } m \in d\}$  (below written simply as  $O$ ). In reverse,  $\text{move}(o)$  is used to refer to the move in which the option  $o$  was proposed.

The proposal moves that introduced the various options have a move status, which will be used to classify the options. Such a classification is any-time and can thus not only be used in selecting the dialogue outcome, but also during the dialogue by agent strategies.

**Definition 11 (Option status).** An option  $o \in O(d)$  for any dialogue  $d$  is:

- *justifiable* iff  $\text{move}(o)$  is *in*,
- *invalid* iff  $\text{player}(\text{move}(o))$  played a move  $m$  such that  $\text{target}(m) = \text{move}(o)$  and  $\text{content}(m) = \text{drop-propose}(o)$ ,
- otherwise it is *defensible*.

Justifiable options are proposals that were questioned but were successfully defended. None of the agents was able to build a warranted case against the proposal. Defensible options are proposals that were attacked by some move that is still warranted. These are thus options that might be reasonable alternatives albeit not being properly supported. Invalid options are those that were retracted by the proposing agent. From the perspective of the multi-agent system, the status of each option hints at its acceptability as dialogue outcome. To settle on one of the options they are first ordered according to some preference.

**Definition 12 (Option preference).** An *option preference* relation  $\preceq$  is a partial order of  $O$ . This is defined as  $o_i \prec o_j$  (strictly preferred) if  $o_i \preceq o_j$  but  $o_j \not\preceq o_i$  and we have  $o_i \approx o_j$  (equally preferred) if  $o_i \preceq o_j$  and  $o_j \preceq o_i$ .

A preliminary ordering on the options can be made. This captures the idea of preferring justifiable options over non-justifiable ones. This may be used during the selection of a dialogue outcome.

**Definition 13 (Preliminary ordering).** Using the set of all options a partition  $O = O_j \cup O_i \cup O_d$  is created such that

- $O_j = \{o \mid o \in O \text{ where } o \text{ is justifiable}\}$ ,
- $O_d = \{o \mid o \in O \text{ where } o \text{ is defensible}\}$ ,
- $O_i = \{o \mid o \in O \text{ where } o \text{ is invalid}\}$ .

Now  $\preceq_p$  is the total *preliminary ordering* over  $O$  such that:

- for every two options  $o_k, o_l \in O_j, O_d$  or  $O_i$  it holds that  $o_k \approx_p o_l$ ,
- for every  $o_j \in O_j$  and  $o_d \in O_d$  it holds that  $o_j \prec_p o_d$ ,
- for every  $o_d \in O_d$  and  $o_i \in O_i$  it holds that  $o_d \prec_p o_i$ .

Justifiable proposals are in principle preferred as dialogue outcome over defensible proposals, which in turn are preferred over invalid ones. However, justifiable options should not always be selected as winner over defensible ones. For one, the preferences as moved by the agents using prefer and prefer-equal moves may be taken into account.

**Definition 14 (Agent option ordering).** Every agent  $a$  has a partial *agent option ordering*  $\preceq_a$  over  $O$  such that for any two options  $o_i, o_j \in O$ :

- $o_i \prec_a o_j$  if the agent played some move  $m$  where  $\mathbf{content}(m) = \mathit{prefer}(o_j, o_i)$ ,
- $o_i \approx_a o_j$  if the agent played some move  $m$  where  $\mathbf{content}(m) = \mathit{prefer-equal}(o_j, o_i)$ .

The protocol forces an agent to be consistent in its preference utterances with relation to the strict ordering of options.

When the dialogue terminates, the deliberation dialogue outcome should be selected from the set of options. How this final selection is achieved is totally dependent on the domain and the purpose of the system. For example, there may be an agent authority that gets to choose the winner, an additional phase may be introduced in which agents vote on the outcome or a function may be used to aggregate all (preliminary and agent-specific) preference orderings. In any case we need to leave open the option for *mutual disagreement* [5].

Preference aggregation is extensively studied in the field of social choice theory and is out of the scope of the present paper. [9] It is interesting to note, though, that when maximum social welfare is desirable it may be good to incorporate the notion of our option status in the winner selection. The valuable information obtained during the deliberation dialogue can be used with a public calculus. This would decide on the outcome in a way similar to the use of public semantics and would not need to rely on agents considering these notions in their voting strategies. For single agents, this is already studied in [1]. How to make use of this is left as future research.

## 8 An Example

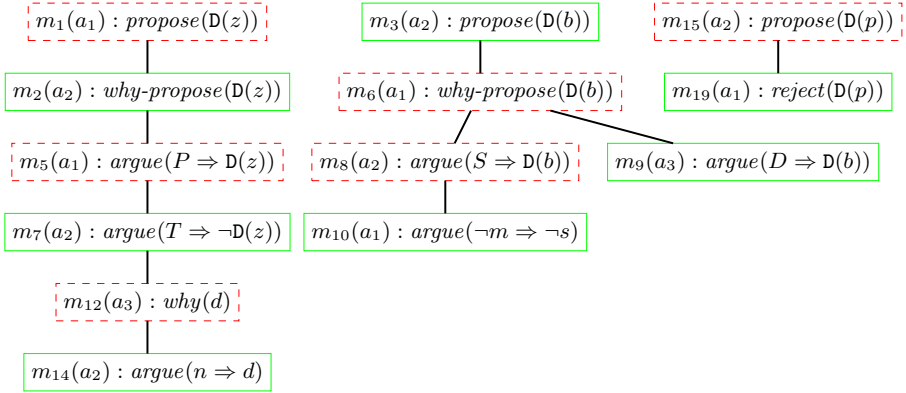
To further explain how the different notions work together, consider an example of three agents  $\mathcal{A} = \{a_1, a_2, a_3\}$  participating in a deliberation dialogue with mutual goal  $g_d$ . We will use all the protocol rules discussed above and adopt a weak form of move relevance. The turns are as follows:

- $T_1$  by  $a_1$   
 $m_1 : \mathit{propose}(\mathbb{D}(z))$  where  $z = \mathit{goToPizzeria}$
- $T_2$  by  $a_2$   
 $m_2 : \mathit{why-propose}(\mathbb{D}(z))$   
 $m_3 : \mathit{propose}(\mathbb{D}(b))$  where  $b = \mathit{goToBistro}$
- $T_3$  by  $a_3$   
 $m_4 : \mathit{skip}$

- $T_4$  by  $a_1$   
 $m_5 : argue(P \Rightarrow \mathbb{D}(z))$  where  
 $P = \{\mathbf{G}(enjoyFood), tastyPizza, goToPizzeria \wedge tastyPizza \rightsquigarrow enjoyFood\}$   
 $m_6 : why-propose(\mathbb{D}(b))$
- $T_5$  by  $a_2$   
 $m_7 : argue(T \Rightarrow \neg \mathbb{D}(z))$  where  
 $T = \{\mathbf{G}(enjoyFood), dullTopping, goToPizzeria \wedge dullTopping \rightsquigarrow \neg enjoyFood\}$   
 $m_8 : argue(S \Rightarrow \mathbb{D}(b))$  where  
 $S = \{\mathbf{G}(enjoyFood), bestSteaks, goToBistro \wedge bestSteaks \rightsquigarrow enjoyFood\}$
- $T_6$  by  $a_3$   
 $m_9 : argue(D \Rightarrow \mathbb{D}(b))$  where  
 $D = \{\mathbf{G}(enjoyFood), greatSeafood, goToBistro \wedge greatSeafood \rightsquigarrow enjoyFood\}$
- $T_7$  by  $a_1$   
 $m_{10} : argue(\neg m \Rightarrow \neg s)$  where  $m = steakOnMenu$
- $T_8$  by  $a_2$   
 $m_{11} : skip$
- $T_9$  by  $a_3$   
 $m_{12} : why(d)$  where  $d = dullTopping$
- $T_{10}$  by  $a_1$   
 $m_{13} : skip$
- $T_{11}$  by  $a_2$   
 $m_{14} : argue(n \Rightarrow d)$  where  $m = menuNeverChanged\}$   
 $m_{15} : propose(\mathbb{D}(p))$  where  $b = goToPub$   
 $m_{16} : prefer(b, p)$   $m_{17} : prefer(p, z)$
- $T_{12}$  by  $a_3$   
 $m_{18} : prefer(b, p)$
- $T_{13}$  by  $a_1$   
 $m_{19} : reject(p)$   
 $m_{20} : prefer(z, b)$   
 $m_{21} : prefer-equal(b, p)$
- $T_{14}$  by  $a_2$   
 $m_{22} : skip$
- $T_{15}$  by  $a_3$   
 $m_{23} : skip$
- $T_{16}$  by  $a_1$   
 $m_{24} : skip$
- $T_{17}$  by  $a_2$   
 $m_{25} : skip$

At that point, the proposal trees of the dialogue will look as represented Fig. 2. To see how the dialogical status and protocol rules affected the agents, consider turn  $T_5$ , in which agent  $a_2$  tries to refute the proposal for  $do(goToPizzeria)$  as made by agent  $a_1$  and support its own proposal for  $do(goToBistro)$ .

To somehow attack proposal  $\mathbb{D}(goToPizzeria)$  the agent needs to find a point of attack, which should always be a relevant move. Within this proposal branch, the only points of attack are to attack  $m_5$  or to move another reply



**Fig. 2.** The proposal trees of the example

to  $m_1$ . A relevant move to  $m_5$  can be both an argue (rebuttal, undercutter or underminer) or a why move. Since the proposal move  $m_1$  was already questioned with a *why-propose* the only remaining valid reply there is to move a *reject*( $D(\text{goToPizzeria})$ ). The agent chooses to rebut the conclusion of  $m_5$  with some argument  $T \Rightarrow \neg D(\text{goToPizzeria})$ .

Within the same turn, the agent also decides to give support to its own proposal  $D(\text{goToBistro})$ . To make this proposal *in*, it will have to find a relevant attack move. In this case the only legal attacking move is to forward an argument with conclusion  $D(\text{goToBistro})$  in reply to  $m_6$ , which it does in the form  $S \Rightarrow D(\text{goToBistro})$ .

Weak relevance is displayed in turn  $T_6$  where agent  $a_3$  make the move *argue*( $D \Rightarrow D(\text{goToBistro})$ ). Although at that point a winning part for the proposal tree of  $D(\text{goToBistro})$  already existed, specifically  $\{m_3, m_6, m_8\}$ , a new winning part  $\{m_3, m_6, m_9\}$  is created. If instead strong relevance is used, then move  $m_9$  is not relevant and thus illegal. In turn, the move  $m_{10}$  by agent  $a_1$  is only weakly relevant because it removed one winning part in the proposal tree without changing the status of **proposal**( $m_8$ ).

The dialogue terminates after turn  $T_{25}$ , when agent  $a_2$  was the first to skip twice in a continuous series of skips. The proposal moves of *goToPizzeria* and *goToPub* are *out* so those options are defensible. The proposal move of *goToBistro* on the other hand is *in* and so this option is justifiable. Partitioning the options set  $O$  according to the option status results in  $O_j = \{\text{goToBistro}\}$  and  $O_d = \{\text{goToPizzeria}, \text{goToPub}\}$ . This gives a preliminary ordering  $\text{goToPizzeria} \approx_p \text{goToPub} \prec_p \text{goToBistro}$ . The agent orderings follow directly from the prefer moves they made. The agent option ordering for  $a_1$  is  $\text{goToBistro} \approx_{a_1} \text{goToPub} \prec_{a_1} \text{goToPizzeria}$ , while that of  $a_2$  is  $\text{goToPizzeria} \prec_{a_2} \text{goToPub} \prec_{a_2} \text{goToBistro}$  and the ordering of  $a_3$  is  $\text{goToBistro} \prec_{a_3} \text{goToPub}$ .

## 9 Basic Fairness and Efficiency Requirements

McBurney et al. [6] have proposed a set of 13 desiderata for argumentation protocols. These are criteria which dialogue game protocols need to adhere for basic fairness and efficiency. Each of the desiderata can be verified against our deliberation framework and protocol.

1. **Stated Dialogue Purpose.** The protocol is explicitly designed to decide on a course of action.
2. **Diversity of individual purpose.** Agents are allowed to have personal goals that possibly conflict with the stated mutual goal.
3. **Inclusiveness.** Many agents can join the deliberation dialogue and no roles are enforced upon them.
4. **Transparency.** The rules of our framework are fully explained, but it is up to an implementation to make sure every agents knows these rules and knows how to play the game.
5. **Fairness.** Every agent has equal rights in the dialogue and the framework allows for fair winner selection methods. Since an agent may always choose not to move (any more) at all, it is never forced to adopt or drop some belief or goal.
6. **Clarity of Argumentation Theory.** The reply structure and notion of relevance in our framework are not hidden implicitly in a protocol, but made explicit. Moreover, the structure of arguments is formalised in an explicitly defined argumentation logic and topic language.
7. **Separation of Syntax and Semantics.** The communication language is separately defined from the protocol. Also, dialogues in the framework are independent of the agent specification while their public behaviour can still be monitored.
8. **Rule Consistency.** We have not studied the rule consistency in detail, but the protocol will never lead to deadlocks; agents can always skip their turn or make a new proposal and within a proposal tree there is always a way to make a new contribution, as long as the top argue move was not conceded.
9. **Encouragement of Resolution.** Agents are encouraged to stay focussed on the dialogue topic through the notion of relevance. If agents still have something to say, there is always the opportunity to do so.
10. **Discouragement of Disruption.** Disruption is discouraged through the definition of legal speech acts, which are separated in attacks and replies. This restricts the available moves, for example agents cannot attack their own moves. However, it is still possible for aggressive agents to question everything that is claimed and no agent is compelled to accept any claim.
11. **Enablement of Self-Transformation.** Agents are allowed to adjust their beliefs or goals depending on the arguments that are moved and preferences that are expressed. Moreover, they are allowed to drop proposals and to retract or concede claims.
12. **System Simplicity.** Simplicity of the system is hard to prove or disprove. However, it is highly modular; communication and topic languages are separated and various alternative protocol rules may be adopted or dropped.

The winner function is left unspecified, but this may range from a dictator agent to a social welfare-based function.

13. **Computational Simplicity.** The computational implications of our framework have not yet been studied. However, the separation of agent and framework designs is at least one step towards simplifying the complexity.

Conforming to these guidelines does not yet mean that every dialogue will be fair and effective. A better understanding is needed of what fair and efficient deliberation dialogues are. Indeed, future work will need to assess how the deliberation process and outcome can be evaluated in relation to the initial situation. In contrast to beliefs, actions will never have an actual truth value but are rather more or less applicable in a specific situation. [5]

New research will also focus on more complete fairness and effectiveness results. For example it is interesting to see how agent attitudes [8] are influential in deliberation dialogues. Moreover, additional formal properties are interesting to study such as the correspondence between the dialogue outcome and the underlying logic of [10].

## 10 Related Work

The literature on argumentation theory for multi-agent systems includes several attempts at designing systems for deliberation dialogues. Earlier we already briefly discussed the most important work on argumentation in deliberation, i.e. that of McBurney et al. [5] They propose a very liberal protocol for agents to discuss proposals restricted by the advancement of a series of dialogue stages. The used speech acts are very similar to that of our framework, although no explicit logic is used to construct and evaluate arguments. Proposals can be forwarded or rejected, claims and arguments are made, questioned or retracted and preferences are expressed. The resulting commitments of agents are determined, but as in our model they are not used to restrict the legal moves.

Specific support is built into their system for discussion of different perspectives about the problem at hand. Perspectives are influential factors such as moral implications and costs. These perspectives can be integrated in our framework as well though the adopted topic language and logic. One model that could be adopted is proposed in [13].

Agents in the framework of McBurney et al. are constrained in their utterances only by preconditions of the different speech acts. For example, they may not state a preference on two actions before they are asserted. Our model accomplishes this through the explicit reply structure of moves rather than using preconditions. Moreover, our model can enforce dialogical coherence through the notion of move relevance.

To decide on a winning proposal agents need to unanimously accept some proposal or a voting system may be used. This way any knowledge of the arguments on proposals is discarded. In contrast, our model may utilise this knowledge on the multi-agent level to decide on a fair winner without the need for a consensus.

A dialogue protocol on proposals for action is introduced in the work of Atkinson et al. [2]. They list all the possible ways to attack a proposal for action, including the circumstances, the goal it achieves and the values it promotes. In our framework, both the goal and action itself are explicitly stated, while the circumstances appear within the arguments that are moved in our deliberation dialogues. As explained earlier, support for values, which are similar to the perspectives of McBurney et al. [5], will be added later.

Many locutions are available to attack proposals, like 'deny goal exists' or 'ask circumstances'. These are needed because no explicit reply structure is present. This also means that no direct relation between the attacks and the resolution of conflicting statements can be made. It is assumed that agents eventually agree on the subject at hand, agree to disagree or use a separate argumentation framework to establish the validity of the proposal. Moreover, the complete work only covers dialogues on a single proposal for action, which makes it persuasion rather than deliberation, albeit being about actions instead of beliefs.

A practical application of multi-agent deliberation dialogues was developed by Tolchinsky et al. [13]. A model for discussion on proposals is coupled to a dialogue game. In the model, agents are proponents or opponents of some proposal for action, while a mediator agent determines the legal moves and evaluates moved arguments to see if they are appropriate and how they support or criticize the proposal for action. Although the paper focusses on the translation and application of argument schemes, it is interesting to see how their work can be modelled inside our framework. The number of proposals is limited to a single action, namely to transplant some organ to some recipient, with a mutual goal to find the best organ donor. A dialogue has to start with propose, reject and why-reject moves after which agents can play argue moves. Whether the proposal is also the winner is determined by the authoritative mediator agent.

## 11 Conclusions

In this paper a framework for multi-agent deliberation dialogues has been proposed. The contribution is twofold.

The general framework for persuasion-type dialogues of Prakken [10] has been altered to provide support for multi-party deliberation dialogues. Consequently, non-trivial modifications have been made to the framework. First, support for moving, criticizing and preferring proposals for action was added. By reusing the explicit reply structure we represent deliberation dialogues as directed multiple trees. Second, the notions of dialogical status and relevance have been adapted for multiple agents. In particular, surrendering replies in a multi-agent context are studied and how strong and weak relevance can still be maintained.

Our framework also improves on the existing work on deliberation dialogues. In contrast with McBurney et al. [5], conflicts of interest are handled through a persuasion-style explicit move status. This allows for varying ways to impose coherence on the deliberating agents. Moreover, the status of proposals is used to define a classification so a preliminary ordering on them can be made.



This, together with the agents' explicit preferences, may be used to select a winning proposal.

The framework was checked against the desiderata for multi-agent argumentation protocols. Deliberation systems in our framework will adhere to those basic standards for efficiency and effectiveness. A more rigid study on formal properties of the framework will be valuable here as well as a study on how different agent strategies can affect fairness and effectiveness.

As an extension of our framework, we could allow agents to discuss not only beliefs but also goals, values and preferences. For example, attacking of preference moves could be allowed, by which a new argument tree is started. A preference-based argumentation framework [7] may be used to in turn evaluate the effect on the dialogical status of proposals. To support discussion on values the topic and communication languages can be extended. One option is to incorporate the work of Black and Atkinson [3], who explicitly allow discussion on promoted values.

**Acknowledgments.** This research was supported by the Netherlands Organisation for Scientific Research (NWO) under project number 612.066.823.

## References

1. Amgoud, L., Prade, H.: Using arguments for making and explaining decisions. *Artificial Intelligence* 173(3-4), 413–436 (2009)
2. Atkinson, K., Bench-Capon, T.J.M., McBurney, P.: A dialogue game protocol for multi-agent argument over proposals for action. *Autonomous Agents and Multi-Agent Systems* 11(2), 153–171 (2005)
3. Black, E., Atkinson, K.: Dialogues that account for different perspectives in collaborative argumentation. In: *Proceedings of the 8th International Conference on Autonomous Agents and Multiagent Systems*, Budapest, Hungary, pp. 867–874 (2009)
4. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artificial Intelligence* 77(2), 321–357 (1995)
5. McBurney, P., Hitchcock, D., Parsons, S.: The eightfold way of deliberation dialogue. *International Journal of Intelligent Systems* 22(1), 95–132 (2007)
6. McBurney, P., Parsons, S., Wooldridge, M.: Desiderata for agent argumentation protocols. In: *Proceedings of the First International Joint Conference on Autonomous Agents and Multiagent Systems*, Bologna, Italy, pp. 402–409 (2002)
7. Modgil, S.: Reasoning about preferences in argumentation frameworks. *Artificial Intelligence* 173(9-10), 901–934 (2009)
8. Parsons, S., Wooldridge, M., Amgoud, L.: Properties and complexity of some formal inter-agent dialogues. *Journal of Logic and Computation* 13(3), 347–376 (2003)
9. Pini, M.S., Rossi, F., Venable, K.B., Walsh, T.: Aggregating Partially Ordered Preferences. *Journal of Logic and Computation* 19(3), 475–502 (2008)
10. Prakken, H.: Coherence and Flexibility in Dialogue Games for Argumentation. *Journal of Logic and Computation* 15(6), 1009–1040 (2005)
11. Prakken, H.: An abstract framework for argumentation with structured arguments. *Argument and Computation* 1(2) ((to appear, 2010))

12. Tang, Y., Parsons, S.: Argumentation-based dialogues for deliberation. In: Proceedings of the 4th International Conference on Multi-Agent Systems, pp. 552–559. ACM Press, New York (2005)
13. Tolchinsky, P., Atkinson, K., McBurney, P., Modgil, S., Cortés, U.: Agents deliberating over action proposals using the *proCLAIM* model. In: Burkhard, H.-D., Lindemann, G., Verbrugge, R., Varga, L.Z. (eds.) CEEMAS 2007. LNCS (LNAI), vol. 4696, pp. 32–41. Springer, Heidelberg (2007)
14. Walton, D.N., Krabbe, E.C.W.: Commitment in dialogue: Basic concepts of interpersonal reasoning. State University of New York Press, New York (1995)
15. van der Weide, T.L., Dignum, F., Meyer, J.-J.C., Prakken, H., Vreeswijk, G.A.W.: Practical reasoning using values. In: McBurney, P., Rahwan, I., Parsons, S., Maudet, N. (eds.) ArgMAS 2009. LNCS, vol. 6057, pp. 79–93. Springer, Heidelberg (2010)

# Empirical Argumentation: Integrating Induction and Argumentation in MAS

Santiago Ontañón and Enric Plaza

IIIA, Artificial Intelligence Research Institute  
CSIC, Spanish Council for Scientific Research  
Campus UAB, 08193 Bellaterra, Catalonia (Spain)  
{santi,enric}@iiia.csic.es

**Abstract.** This paper presents an approach that integrates notions and techniques from two distinct fields of study —namely inductive learning and argumentation in multiagent systems (MAS). We will first discuss inductive learning and the role argumentation plays in multiagent inductive learning. Then we focus on how inductive learning can be used to realize argumentation in MAS based on empirical grounds. We present a MAS framework for empirical argumentation, A-MAIL, and then we show how this is applied to a particular task where two agents argue in order to reach agreement on a particular topic. Finally, an experimental evaluation of the approach is presented evaluating the quality of the agreements achieved by the empirical argumentation process.

## Categories and Subject Descriptors

I.2.11 [Artificial Intelligence]: Distributed Artificial Intelligence — Multiagent systems, Intelligent Agents.

I.2.6 [Artificial Intelligence]: Learning.

**General Terms:** Algorithms, Experimentation, Theory.

**Keywords:** Argumentation, Learning.

## 1 Introduction

This paper presents an approach that integrates notions and techniques from two distinct fields of study —namely inductive learning and argumentation in multiagent systems (MAS). We will first discuss inductive learning and the role argumentation may play in multiagent inductive learning, and later how inductive learning can be used to realize argumentation in MAS based on empirical grounds.

Multiagent inductive learning (MAIL) is the study of multiagent systems where individual agents have the ability to perform inductive learning, i.e. where agents are able to learn general descriptions from particular examples. Induction is a form of empirical-based inference, where what is true (or what is believed by

the agent) is derived from the experience of that agent in a particular domain (such experience is usually represented with “cases” or “examples”). Notice that inductive inference is not deductive, and specifically it is not truth-preserving<sup>1</sup>, and therefore it captures a form of empirical knowledge that can be called into question by new empirical data and thus needs to be revised.

The challenge of multiagent inductive learning is that several agents will inductively infer empirical knowledge that in principle may not be the same, since that knowledge is dependent on each individual in two ways: the concrete empirical data an agent has encountered and the specific inductive method an agent employs.

Communication among agents is necessary in order to reach shared and agreed-upon empirical knowledge that is based on, and consistent with, all the empirical data available to a collection of agents. Agents could simply communicate all the data to the other agents, and then each agent could just use induction individually. However, data redistribution might have a high cost, or might not even be feasible in some domains due to organizational or privacy issues. In this paper we propose an argumentation-based communication process where agents can propose, compare and challenge the empirical knowledge of other agents, with the goal of achieving a more accurate, shared, and agreed-upon body of empirical knowledge without having to share all of their empirical data.

From the point of view of argumentation in MAS, inductive learning provides a basis for automating, in empirical domains, a collection of activities necessary for implementing artificial agents that support argumentation: how to generate arguments, how to attack and defend arguments, and how to change an agent’s beliefs as a result of the arguments exchanged. Logic-based approaches to argumentation like DeLP [2] amend classical deductive logic to support defeasible reasoning. Our approach takes a different path, assuming agents that *learn their knowledge* (by using induction over empirical data) instead of assuming agents have been *programmed* (by giving them a rule-based knowledge base). Therefore, we need to specify empirical methods that are able to perform the required activities of argumentation (generating arguments and attacks, comparing arguments and revising an agent’s beliefs).

This paper presents a MAS framework for empirical argumentation called A-MAIL, which implements those activities on the basis of the inductive inference techniques developed in the field of machine learning. The main idea behind A-MAIL is the following: given two agents with inductive learning capabilities, they can use induction to generate hypotheses from examples. These hypotheses can be used as arguments in a computational argumentation framework. Argumentation helps the agents reach an agreement over the induced knowledge, thus reaching hypotheses that are consistent with the data known to both agents. Effectively, A-MAIL integrates inductive learning and computational argumentation to let groups of agents perform multiagent induction. This means that agents can reach hypotheses consistent with the data known to a set of agents without having to share all this data.

---

<sup>1</sup> Inductive inference is not truth-preserving, since new and unseen examples may contradict past generalizations, albeit it is falsity-preserving.

The structure of the paper starts by introducing the needed notions of inductive learning (Section 2). Then, Section 3 presents our empirical argumentation framework, A-MAIL, while Section 4 shows the utility of the framework in the task of concept convergence (in which two agents argue with the goal of achieving an agreement on a particular topic); an experimental evaluation of the approach is presented evaluating the quality of the agreements achieved by argumentation. The paper closes with sections on related work and conclusions.

## 2 Concept Induction

Inductive learning, and in particular concept learning, is the process by which given an *extensional definition* of a concept  $C$  (a collection of examples of  $C$  and a collection of examples that are not  $C$ ) an *intensional definition* (or generalization) of a concept  $C$  can be found. Formally, an *induction domain* is characterized as pair  $\langle \mathcal{E}, \mathcal{G} \rangle$  where  $\mathcal{E}$  is the language describing examples or instances and  $\mathcal{G}$  is the language for describing generalizations; usually  $\mathcal{E} \subset \mathcal{G}$  is assumed, but this is not necessary. A language is understood as the set of well formed formulas built from a domain vocabulary or ontology  $\mathcal{O}$ . The relation between languages  $\mathcal{E}$  and  $\mathcal{G}$  is established by the subsumption relation ( $\sqsubseteq$ ); we say a generalization  $g \in \mathcal{G}$  subsumes (or covers) an example  $e \in \mathcal{E}$ ,  $g \sqsubseteq e$ , whenever  $e$  satisfies the properties described by  $g$  [8]. Different approaches to induction work with different languages, from propositional languages (attribute value vectors) to subsets of predicate logic (like Inductive Logic Programming that uses a sublanguage of Horn logic).

Given a collection of examples  $E = \{e_1, \dots, e_M\}$  described in a language  $\mathcal{E}$ , an extensional definition of a concept  $C$  is a function  $C : E \rightarrow \{+, -\}$ , that determines the subset  $E^+$  of (positive) examples of  $C$ , and the subset  $E^-$  of counterexamples (or negative examples) of  $C$ . An inductive concept learning method is a function  $I : \mathcal{P}(E) \times C \rightarrow \mathcal{G}$  such that, given a collection of examples and a target concept  $C$ , yields an intensional definition  $h \in \mathcal{G}$ ; generally one single formula in  $\mathcal{G}$  is not sufficient to describe an intensional definition so it is usually described as a disjunction of generalizations  $C = h_1 \vee \dots \vee h_n$ .

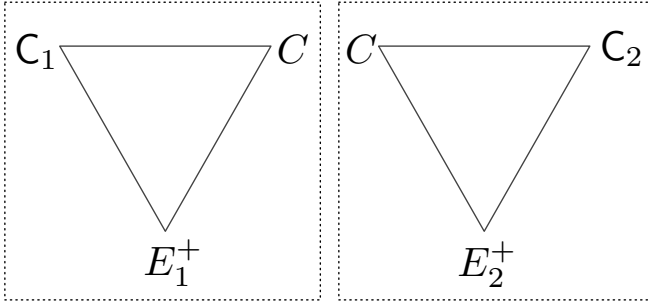
**Definition 1.** An intensional definition  $C$  of a concept  $C$  is a disjunct  $C = h_1 \vee \dots \vee h_n$ , such that  $\forall e_j \in E^+ \exists h_i : h_i \sqsubseteq e_j$  and  $\forall e_j \in E^- \forall h_i : h_i \not\sqsubseteq e_j$

That is to say, that each positive example of  $C$  is subsumed by at least one generalization  $h_i$ , and no counterexample of  $C$  is subsumed by any  $h_i$ .

For simplicity, we will shorten the previous expression as follows:  $C \sqsubseteq E^+ \wedge C \not\sqsubseteq E^-$ . Moreover, in the remainder of this paper we will refer to each  $h_i$  as a generalization or as a hypothesis.

### 2.1 Inductive Agents with Empirical Beliefs

In this paper we will focus on argumentation between two agents (say  $A_1$  and  $A_2$ ) that are interested in learning an intensional definition for a particular concept based on the experience of both agents. Each agent will have certain beliefs



**Fig. 1.** Schema for two agents where a concept name ( $C$ ) is shared while intensional descriptions are, in general, not equivalent ( $C_1 \not\equiv C_2$ )

according to what they have learnt. Thus, we will now explore how differences between these two agents relate to induction and argumentation. First, we will assume each agent has its own set of examples from which they may learn by induction (say  $E_1$  and  $E_2$ ) and they are both in principle unrelated although expressed in the same language  $\mathcal{E}$ . Furthermore, each agent may use, in principle, different induction techniques but they obtain generalizations in the same language  $\mathcal{G}$ . Thus, for any particular concept  $C$  two agents will have intensional descriptions  $C_1$  and  $C_2$  that are, in general, not equal or equivalent. Figure 1 depicts these relationships between two agents beliefs ( $C_1$  and  $C_2$ ) about what  $C$  is based on their empirical data  $E_1$  and  $E_2$ .

Finally, since Definition 1 is too restrictive for practical purposes, machine learning approaches allow the intensional definitions to subsume less than 100% of positive examples by defining a confidence measure. The goal of induction is then, given as a target the function  $C : E \rightarrow \{+, -\}$ , to find a new function  $C$ , which is a good approximation of  $C$ , in the sense of yielding a small error in determining when an example is a positive or negative example of  $C$ .

In the remainder of this paper we will use a confidence measure that assesses the confidence of each individual hypothesis  $h$  in an intensional definition.

**Definition 2.** *The individual confidence of a hypothesis  $h$  for an agent  $A_i$ :*

$$B_i(h) = \frac{|\{e \in E_i^+ | h \sqsubseteq e\}| + 1}{|\{e \in E_i | h \sqsubseteq e\}| + 2}$$

$B_i(h)$  is the ratio of positive examples correctly covered by  $h$  over the total number examples covered by  $h$ ; moreover, we add 1 to the numerator and 2 to the denominator following the Laplace probability estimation procedure (which prevents estimations too close to 0 or 1 when very few examples are covered). Other confidence measures could be used, our framework only requires that the confidence measure reflects how much the set of examples known to an agent endorses a hypothesis  $h$ .

Finally, a threshold  $\tau$  is established, and only hypotheses with confidence  $B_i(h) > \tau$  are accepted as valid outcomes of the inductive process.

**Definition 3.** A hypothesis  $h$  is  $\tau$ -acceptable for an agent  $A_i$  if  $B_i(h) \geq \tau$ , where  $0 \leq \tau \leq 1$ .

Thus, intensional definitions ( $C_1$  and  $C_2$ ) consist of a disjunction of hypotheses, each of them being  $\tau$ -acceptable. In the rest of this paper we will say that a hypothesis is *consistent* with a set of examples, if the hypothesis is  $\tau$ -acceptable with respect to that set of examples.

### 3 An Empirical Approach to MAS Argumentation

This section will focus on how to integrate argumentation with inductive agents in scenarios where the goal is to achieve an agreement between two agents on the basis of their empirical knowledge. Here the *empirical* adjective refers to the observations of the real world that each agent has had access to and that is embodied in the set of examples  $E_1$  and  $E_2$  represented using a language  $\mathcal{E}$ .

Argumentation in Multiagent Inductive Learning (A-MAIL) is a framework where argumentation is used as a communication mechanism for agents that want to perform collaborative inductive tasks such as concept convergence (see Section 4). We do not claim, however, that A-MAIL is a new “argumentation framework” in the sense of Dung [6], it is intended as a framework to that integrates argumentation processes and inductive processes in MAS.

According to Dung, an argumentation framework  $AF = \langle A, R \rangle$  is composed by a set of arguments  $A$  and an attack relation  $R$  among the arguments. A-MAIL is not a general logic framework and, although certainly we will define what we mean as arguments and attack relations, we take an empirical approach to argumentation. Thus, the main difference from Dung’s framework is that, since arguments are generated from examples, our approach necessarily defines a specific relation between arguments and examples, which is not part of the usual interpretations of Dung’s framework<sup>2</sup>.

#### 3.1 The A-MAIL Framework

A-MAIL is a framework that allows groups of agents to perform collaborative induction tasks. A typical collaborative induction task is multiagent induction, where a group of agents wants to find an intensional definition of a concept and where each agents has a different set of positive and negative examples of that concept. A simple way to solve this problem is by sharing all the examples and then just using induction in a centralized way. However, that solution might not be feasible in some scenarios. Imagine, for instance, that a group of physicians needed to share the data concerning all of their patients to a centralized location in order to draw inductive inferences from that data. Another approach could be

<sup>2</sup> Some approaches may consider “counter-examples” as a kind of arguments. This is certainly true, but in our approach there is a constitutive relation between examples and arguments (the “empirical” approach) that is different from merely accepting counter-examples as arguments.

use *ensemble learning* [5] techniques, where each agent would learn a local intensional definition, and then those definitions can be combined at problem solving time using some sort of voting mechanism. A-MAIL is an alternative approach where agents first use induction individually, and then use computational argumentation to argue about the individually induced hypothesis. Nevertheless, in this paper we focus on scenarios with only two agents; extending A-MAIL for more than two agents is part of our future work.

The main idea behind A-MAIL is that the arguments to be used in an argumentation process can be generated from examples by inductive learning methods. Agents using A-MAIL use induction to generate an initial set of hypotheses explaining the data known to them, and then communicate those hypotheses to other agents, starting an argumentation process where arguments and counterarguments (also generated by induction) are exchanged until an agreement is reached. While sharing arguments and counterarguments, the agents learn new information from the data known to the other agents, and may need to revise their beliefs accordingly; once the argumentation process is over, the agents will have agreed on a set of hypotheses that are consistent with the data known to each other (including the exchanged in the process).

Summarily, there are three main processes in the A-MAIL framework: 1) generation of arguments from examples using inductive learning, 2) computational argumentation using the previously generated arguments, and 3) belief revision, for revising the hypotheses generated by induction in front of new arguments received from other agents. Let us address each one of them in turn.

### 3.2 Arguments and Counterarguments

We first define the kinds of arguments employed in A-MAIL and their attack relation. There are two kinds of arguments in A-MAIL:

**Example Argument:**  $\alpha = \langle e, \overline{C} \rangle$  is a pair where an example  $e \in \mathcal{E}$  is related to a concept  $\overline{C} \in \{C, \neg C\}$ , where  $\overline{C} = C$  if  $e$  is a positive example of  $C$ , and  $\overline{C} = \neg C$  otherwise.

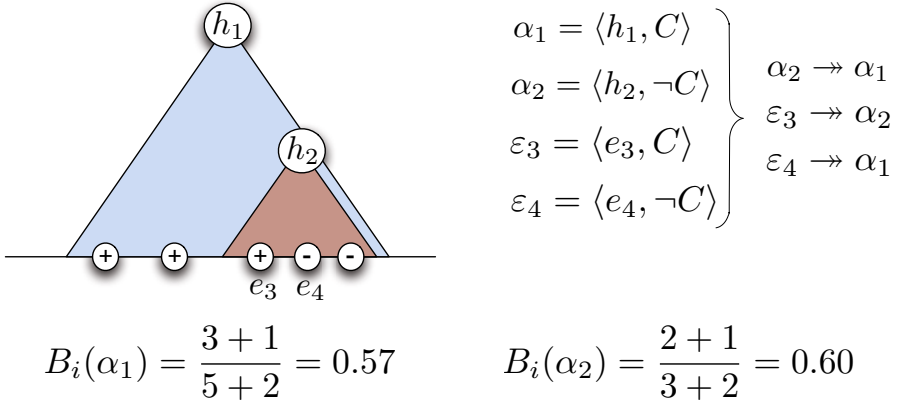
**Hypothesis Argument:**  $\alpha = \langle h, \overline{C} \rangle$  is a pair where  $h$  is a  $\tau$ -acceptable hypothesis and  $\overline{C} \in \{C, \neg C\}$ . An argument  $\langle h, C \rangle$  states that  $h$  is a hypothesis of  $C$ , while  $\langle h, \neg C \rangle$  states that  $h$  is a hypothesis of  $\neg C$ , i.e. that examples covered by  $h$  do not belong to  $C$ .

Since hypotheses in arguments are generated by induction, they have an associated degree of confidence for an individual agent:

**Definition 4.** *The confidence of a hypothesis argument  $\alpha = \langle h, \overline{C} \rangle$  for an agent  $A_i$  is:*

$$B_i(\alpha) = \begin{cases} \frac{|\{e \in E_i^+ | h \sqsubseteq e\}| + 1}{|\{e \in E_i | h \sqsubseteq e\}| + 2} & \text{if } \overline{C} = C \\ \frac{|\{e \in E_i^- | h \sqsubseteq e\}| + 1}{|\{e \in E_i | h \sqsubseteq e\}| + 2} & \text{if } \overline{C} = \neg C \end{cases}$$





**Fig. 2.** An illustration of the different argument types, their confidences and attacks

Consequently, we can use the threshold  $\tau$  to impose that only arguments with a strong confidence are acceptable in the argumentation process.

**Definition 5.** An argument  $\alpha$  generated by an agent  $A_i$  is  $\tau$ -acceptable iff  $\alpha$  is a hypothesis argument and  $B_i(\alpha) > \tau$ , or if  $\alpha$  is an example argument.

From now on, only  $\tau$ -acceptable arguments will be considered within the A-MAIL framework. Moreover, notice that we require arguments to be  $\tau$ -acceptable for the agent who generates them. An argument generated by one agents might not be  $\tau$ -acceptable for another agent.

Next we define the attack relation between arguments:

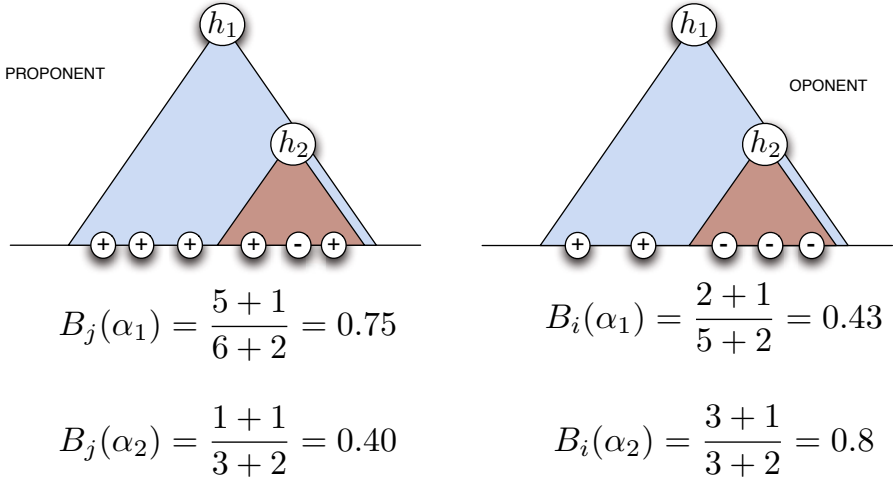
**Definition 6.** An attack relation ( $\alpha \rightarrow \beta$ ) between two  $\tau$ -acceptable arguments  $\alpha, \beta$  holds when:

1.  $\langle h_1, \widehat{C} \rangle \rightarrow \langle h_2, \overline{C} \rangle \iff \widehat{C} = \overline{C} \wedge h_2 \sqsubset h_1$ , or
2.  $\langle e, \overline{C} \rangle \rightarrow \langle h, \widehat{C} \rangle \iff \overline{C} = \neg \widehat{C} \wedge h \sqsubseteq e$

where  $\overline{C}, \widehat{C} \in \{C, \neg C\}$ .

Notice that a hypothesis argument  $\alpha = \langle h_1, \widehat{C} \rangle$  only attacks another argument  $\beta = \langle h_2, \overline{C} \rangle$  if  $h_2 \sqsubset h_1$ , i.e. when  $\alpha$  is (strictly) more specific than  $\beta$ . This is required since it implies that all the examples covered by  $\alpha$  are also covered by  $\beta$ , and thus if one supports  $C$  and the other  $\neg C$ , they must be in conflict.

Figure 2 shows some examples of arguments and attacks. Positive examples of the concept  $C$  are marked with a positive sign, whereas negative examples are marked with a negative sign. Hypothesis arguments are represented as triangles covering examples; when an argument  $\alpha_1$  subsumes another argument  $\alpha_2$ , we draw  $\alpha_2$  inside of the triangle representing  $\alpha_1$ . Argument  $\alpha_1$  has a hypothesis  $h_1$  supporting  $C$ , which covers 3 positive examples and 2 negative examples, and



**Fig. 3.** An comparison of two individual viewpoints on arguments, attacks, and acceptability

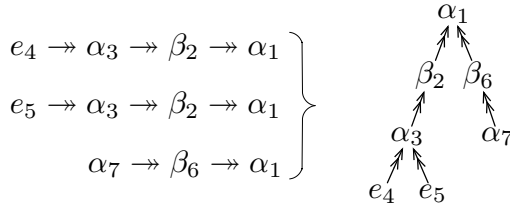
thus has confidence 0.57, while argument  $\alpha_2$  has a hypothesis  $h_2$  supporting  $\neg C$  with confidence 0.60, since  $h_2$  covers 2 negative examples and only one positive example. Now, the attack  $\alpha_2 \rightarrow \alpha_1$  holds because  $\alpha_2$  supports  $\neg C$ ,  $\alpha_1$  supports  $C$  and  $h_1 \sqsubseteq h_2$ . Moreover,  $\varepsilon_3 \rightarrow \alpha_2$ , since  $e_3$  is a positive example of  $C$  while  $\alpha_2$  supports  $\neg C$  and covers this example ( $h_2 \sqsubseteq e_3$ ).

Notice that the viewpoint on the (empirical) acceptability of an argument or of an attack depends on each individual agent, as shown in Fig 3, where two agents  $A_i$  and  $A_j$  compare arguments  $\alpha_1$  and  $\alpha_2$  for hypotheses  $h_1$  and  $h_2$ , assuming  $\tau = 0.6$ . From the point of view of agent  $A_i$  (the Opponent), proposing argument  $\alpha_2$  as an attack against argument  $\alpha_1$  of agent  $A_j$  (the Proponent) is a sound decision, since for  $A_i$ ,  $\alpha_1$  is not  $\tau$ -acceptable, while  $\alpha_2$  is. However, from the point of view of the Proponent of  $\alpha_1$ ,  $\alpha_2$  is not  $\tau$ -acceptable. Thus,  $A_j$  does not accept  $\alpha_2$  and will proceed by attacking it.

Next we will define when arguments *defeat* other arguments, based on the notion of argumentation lines [2].

**Definition 7.** An Argumentation Line  $\alpha_n \rightarrow \alpha_{n-1} \rightarrow \dots \rightarrow \alpha_1$  is a sequence of  $\tau$ -acceptable arguments where  $\alpha_i$  attacks  $\alpha_{i-1}$ , and  $\alpha_1$  is called the root.

Notice that odd-numbered arguments are generated by the agent whose hypothesis is under attack (the Proponent of the root argument  $\alpha_1$ ) and the even-numbered arguments are generated by the Opponent agent attacking  $\alpha_1$ . Moreover, since hypothesis arguments can only attack other hypothesis arguments, and example arguments can only attack hypothesis arguments, example arguments can only appear as the left-most argument (e.g.  $\alpha_n$ ) in an argumentation line.



**Fig. 4.** Multiple argumentation lines rooted in the same argument  $\alpha_1$  can be composed into an argumentation tree

**Definition 8.** An  $\alpha$ -rooted argumentation tree  $T$  is a tree where each path from the root node  $\alpha$  to one of the leaves constitutes an argumentation line rooted on  $\alpha$ . The example-free argumentation tree  $T^f$  corresponding to  $T$  is a tree rooted in  $\alpha$  that contains the same hypothesis arguments of  $T$  and no example argument.

Therefore, a set of argumentation lines rooted in the same argument  $\alpha_1$  can be represented as an argumentation tree, and vice versa. Notice that example arguments can only appear as leaves in any argumentation tree.

Figure 4 illustrates this idea, where three different argumentation lines rooted in the same  $\alpha_1$  are shown with their corresponding argumentation tree. The  $\alpha_i$  arguments are provided by the Proponent agent (the one proposing the root argument) while  $\beta_i$  arguments are provided by the Opponent trying to attack the Proponent's arguments.

**Definition 9.** Let  $T$  be an  $\alpha$ -rooted argumentation tree, where argument  $\alpha$  belongs to an agent  $A_i$ , and let  $T^f$  be the example-free argumentation tree corresponding to  $T$ . Then the root argument  $\alpha$  is warranted (or undefeated) iff all the leaves of  $T^f$  are arguments belonging to  $A_i$ ; otherwise  $\alpha$  is defeated.

In A-MAIL agents will exchange arguments and counterarguments following some interaction protocol. The protocol might be different depending on the task the agents are trying to achieve (be it concept convergence, multiagent induction, or any other). Nevertheless, independently of the protocol being used, we can define the state of the argumentation two agents  $A_i$  and  $A_j$  at an instant  $t$  as the tuple  $\langle R_i^t, R_j^t, G^t \rangle$ , consisting of:

- $R_i^t = \{ \langle h, C \rangle \mid h \in \{h_1, \dots, h_n\} \}$ , the set of arguments defending the current intensional definition  $C_i^t = h_1 \vee \dots \vee h_n$  of agent  $A_i$ ;
- $R_j^t$  is the same for  $A_j$ .
- $G^t$  contains the collection of arguments generated before  $t$  by either agent, and belonging to one argumentation tree rooted in an argument in  $R_i^t \cup R_j^t$ .

### 3.3 Argument Generation through Induction

Agents need two kinds of argument generation capabilities: generating an intensional definition from the individual examples known to an agent, and generating

arguments that attack arguments provided by other agents; notice that a defense argument is simply  $\alpha' \rightarrow \beta \rightarrow \alpha$ , i.e. an attack on the argument attacking a previous argument. Thus, defense need not be considered separately.

An agent  $A_i$  can generate an intensional definition of  $C$  by using any inductive learning algorithm capable of learning concepts as a disjunction of hypothesis, e.g. learning algorithms such as CN2 [3] or FOIL [13].

Attack arguments, however, require a more sophisticated form of induction. When an agent  $A_i$  wants to generate an argument  $\beta = \langle h_2, \overline{C} \rangle$  to attack another argument  $\alpha = \langle h_1, \widehat{C} \rangle$ , i.e.  $\beta \rightarrow \alpha$ ,  $A_i$  has to find an inductive hypothesis  $h_2$  for  $\beta$  that satisfies four conditions:

1.  $h_2$  should support the opposite concept than  $\alpha$ : namely  $\overline{C} = \neg \widehat{C}$ ,
2.  $\beta$  should have a high confidence  $B_i(\beta)$  (at least being  $\tau$ -acceptable),
3.  $h_2$  should satisfy  $h_1 \sqsubset h_2$ , and
4.  $\beta$  should not be attacked by any undefeated argument in  $G^t$ .

Currently existing inductive learning techniques cannot be applied out of the box, mainly because they do not satisfy the last two conditions.

In previous work, we developed the Argumentation-based Bottom-up Induction (ABUI) algorithm, capable of performing such task [12]; this is the inductive algorithm used in our experiments in Section 4.2. However, any algorithm which can search the space of hypotheses looking for a hypothesis which satisfies the four previous conditions would work in our framework.

Let  $L$  be the inductive algorithm used by an agent  $A_i$ ; when the goal is to attack an argument  $\alpha = \langle h_1, \widehat{C} \rangle$  then  $L$  has to generate an argument  $\beta = \langle h_2, \overline{C} \rangle$  such that  $\beta \rightarrow \alpha$ . The uses  $L$  trying to find such a hypothesis  $h_2$ :

- If  $L$  returns an individually  $\tau$ -acceptable  $h_2$ , then  $\beta$  is the attacking argument to be used.
- If  $L$  fails to find a suitable  $h_2$ , then  $A_i$  looks for examples in  $E_i$  that attack  $\alpha$ . If any exist, then one such example  $e$  is randomly chosen to be used as an attacking argument  $\beta = \langle e, \overline{C} \rangle$ .

Otherwise,  $A_i$  is unable to generate any argument attacking  $\alpha$ .

If a hypothesis or example argument is not enough to defeat another argument, additional arguments or examples could be sent in subsequent rounds of the interaction protocol (as long as the protocol allows it).

### 3.4 Empirical Belief Revision

During argumentation, agents exchange arguments which contain new hypotheses and examples. These exchanges contain empirical knowledge that agents will integrate with their previous empirical beliefs. Consequently, their beliefs will change in such a way that their hypotheses are consistent with the accrued empirical evidence: we call this process empirical belief revision.

The belief revision process of an agent  $A_i$  at an instant  $t$ , with an argumentation state  $\langle R_i^t, R_j^t, G^t \rangle$  starts whenever  $A_i$  receives an argument from another agent:

1. If it is an example argument  $\varepsilon = \langle e, \widehat{C} \rangle$  then  $e$  is added as a new example into  $E_i$ , i.e.  $A_i$  expands its extensional definition of  $C$ .
2. Whether the received argument is an example or an hypothesis, the agent re-evaluates the confidence of the arguments in  $R_i^t$  and  $G^t$ : if any of these arguments becomes no longer  $\tau$ -acceptable for  $A_i$  they are removed from  $R_i^{t+1}$  or  $G^{t+1}$ .
3. If any argument  $\alpha = \langle h, \widehat{C} \rangle$  in  $R_i^t$  became defeated, and  $A_i$  is not able to expand the argumentation tree rooted in  $\alpha$  to defend it, then the hypothesis  $h$  will be removed from  $C_i$ . This means that some positive examples in  $E_i$  will not be covered by  $C_i$  any longer. The inductive learning algorithm is called again to generate new hypotheses  $h'$  for the now uncovered examples.

We would like to remark that, as shown in Figure 5, all aspects of the argumentation process (generating arguments and attacks, accepting arguments, determining defeat, and revising beliefs) are supported on an empirical basis and, from the point of view of MAS, implemented by autonomous decision making of artificial agents. The activities in Figure 5 permit the MAS to be self-sufficient in a domain of empirical enquiry, since individual agents are autonomous and every decision is based on the empirical knowledge available to them.

The next section presents an application of this MAS framework to reach agreements in MAS.

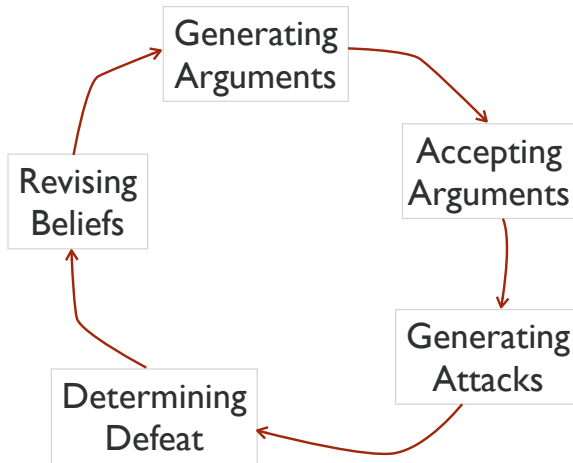


Fig. 5. The closed loop of empirically based activities used in argumentation

## 4 Concept Convergence

We have developed A-MAIL as part of our research line on deliberative agreement<sup>3</sup>, in which 2 or more artificial agents use argumentation to reach different

<sup>3</sup> This is part the project Agreement Technologies: <http://www.agreement-technologies.org/>

forms of agreement. In this section we will present a particular task of deliberative agreement called concept convergence. The task of *Concept Convergence* is defined as follows: Given two or more individuals which have individually learned non-equivalent meanings of a concept  $C$  from their individual experience, find a shared, equivalent, agreed-upon meaning of  $C$ .

**Definition 10.** *Concept Convergence (between 2 agents) is the task defined as follows:*

**Given** two agents ( $A_i$  and  $A_j$ ) with individually different intensional ( $C_i \not\cong C_j$ ) and extensional definitions ( $E_i^+ \neq E_j^+$ ) of a concept  $C$ ,

**Find** a convergent, shared and agreed-upon intensional description ( $C_i \cong C_j$ ) for  $C$  that is consistent with the extensional descriptions ( $E_i^+$  and  $E_j^+$ ) of each individual.

For example, in the experiments reported in this paper, we used the domain of marine sponge identification. The two agents need to agree on the definition of the target concept  $C = \textit{Hadromerida}$ , among others. While in ontology alignment the focus is on establishing a mapping between the ontologies of the two agents, here we assume that the ontology is shared, i.e. both agents share the concept name *Hadromerida*. Each agent may have experience in a different area (say, one in the Atlantic, and the other in the Mediterranean), so they have collected different samples of *Hadromerida* sponges, those samples constitute their extensional definitions (which are different, since each agent has collected sponges on their own). Now, they would like to agree on an intensional definition  $C$ , which describes such sponges and is consistent with their individual experience. In our experiments, one such intensional definition reached by one of the agents is:  $C =$  “all those sponges which do not have gemmules in their external features, whose megascleres had a tylostyle smooth form and that do not have a uniform length in their spikulate skeleton”.

Concept convergence is assessed individually by an agent  $A_i$  by computing the *individual degree of convergence* among two definitions  $C_i$  and  $C_j$ , as follows:

**Definition 11.** *The individual degree of convergence among two intensional definitions  $C_i$  and  $C_j$  for an agent  $A_i$  is:*

$$K_i(C_i, C_j) = \frac{|\{e \in E_i \mid C_i \sqsubseteq e \wedge C_j \sqsubseteq e\}|}{|\{e \in E_i \mid C_i \sqsubseteq e \vee C_j \sqsubseteq e\}|}$$

where  $K_i$  is 0 if the two definitions are totally divergent, and 1 when the two definitions are totally convergent. The degree of convergence corresponds to the ratio between the number examples covered by both definitions (intersection) and the number of examples covered by at least one definition (union). The closer the intersection is to the union, the more similar the definitions are.

**Definition 12.** *The joint degree of convergence of two intensional definitions  $C_i$  and  $C_j$  is:*

$$K(C_i, C_j) = \min(K_i(C_i, C_j), K_j(C_j, C_i))$$

Concept convergence is defined as follows:

**Definition 13.** *Two intensional definitions are convergent ( $C_i \cong_\epsilon C_j$ ) if  $K(C_i, C_j) \geq \epsilon$ , where  $0 \leq \epsilon \leq 1$  is a the degree of convergence required.*

The next section describes the protocol to achieve concept convergence.

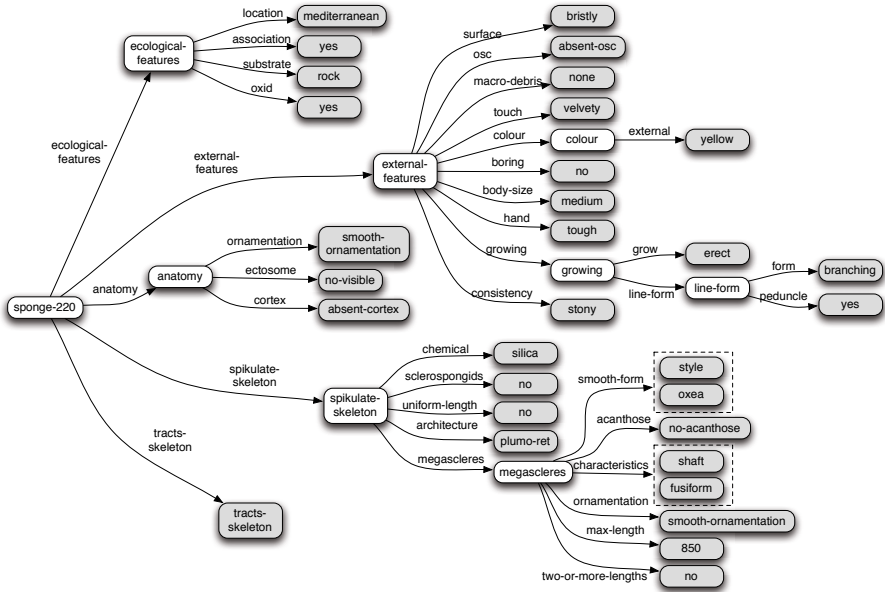
#### 4.1 Argumentation Protocol for Concept Convergence

The concept convergence (CC) argumentation process follows an iteration protocol composed of a series of rounds, during which two agents will argue about the individual hypotheses that compose their intensional definitions of a concept  $C$ . At each round  $t$  of the protocol, each agent  $A_i$  holds a particular intensional definition  $C_i^t$ , and only one agent will hold a *token*. The holder of the token can assert new arguments and then the token will be passed on to the other agent. This cycle will continue until  $C_i \cong C_j$ .

The protocol starts at round  $t = 0$  with a value set for  $\epsilon$  and works as follows:

1. Each agent  $A_i$  communicates to the other their current intensional definition by sharing  $R_i^0$ . The token is given to one agent at random, and the protocol moves to 2.
2. The agents share  $K_i(C_i, C_j)$  and  $K_j(C_j, C_i)$ , their individual convergence degrees. If  $C_i \cong_\epsilon C_j$  the protocol ends with success; if no agent has produced a new attack in the last two rounds then the protocol ends with failure; otherwise it moves to 3.
3. The agent with the token,  $A_i$ , checks if belief revision has modified  $C_i^t$ , and if so sends a message communicating its current intensional definition  $R_i^t$ . Then, the protocol moves to 4.
4. If any argument  $\alpha \in R_i^t$  is defeated, and  $A_i$  can generate an argument  $\alpha'$  to defend  $\alpha$ , the argument  $\alpha'$  will be sent to the other agent. Also, if any of the undefeated arguments  $\beta \in R_j^t$  is not individually  $\tau$ -acceptable for  $A_i$ , and  $A_i$  can find an argument  $\beta'$  to extend any argumentation line rooted in  $\beta$ , in order to attack it, then  $\beta'$  is sent to the other agent. If at least one of these arguments was sent, a new round  $t + 1$  starts; the token is given to the other agent, and the protocol moves back to 2. Otherwise, if none of these arguments could be found, the protocol moves to 5.
5. If there is any example  $e \in E_i^+$  such that  $C_j^t \not\sqsubseteq e$  (i.e. a positive example not covered by the definition of  $A_j$ ),  $A_i$  will send  $e$  to the other agent, stating that the intentional definition of  $A_j$  does not cover  $e$ . A new round  $t + 1$  starts, the token is given to the other agent, and the protocol moves to 2.

Moreover, in order to ensure termination, no argument is allowed to be sent twice by the same agent. A-MAIL ensures that the joint degree of convergence of the resulting concepts is at least  $\tau$  if (1) the number of examples is finite, (2) the number of hypotheses that can be generated is finite. Joint convergence degrees higher of than  $\tau$  cannot be ensured, since  $100 \times (1 - \tau)\%$  of the examples covered by a  $\tau$ -acceptable hypothesis might be negative, causing divergence. Therefore,



**Fig. 6.** A description of one of the sponges of the Axinellida order used in our experiments

when  $\epsilon > \tau$ , we cannot theoretically ensure convergence. However, as we will show in our experiments, in practical scenarios, convergence is almost always reached. Notice that increasing  $\tau$  too much in order to ensure convergence could be detrimental, since that would impose a too strong restriction on the inductive learning algorithms. And, although convergence would be reached, the concept definitions might cover only a small subset of the positive examples.

Termination is assured even when both agents use different inductive algorithms because of the following reason. By assumption, agents use the same finite generalization space, and thus there is no hypothesis  $\tau$ -acceptable by one agent that could not be  $\tau$ -acceptable by the other agent when both use the same acceptability condition over the same collection of examples. Thus, in the extreme, if the agents reach the point when they have exchanged all their examples, their  $\tau$ -acceptability criteria will be identical, and thus all rules acceptable to one are also acceptable to the other.

## 4.2 Experimental Evaluation

In order to empirically evaluate A-MAIL with the purpose of concept convergence we used the marine sponge identification problem. Sponge classification is interesting because the difficulties arise from the morphological plasticity of the species, and from the incomplete knowledge of many of their biological and cytological features. Moreover, benthology specialists are distributed around the world and they have experience in different benthos that spawn species with



**Table 1.** Precision (P), Recall (R) and degree of convergence (K) for the intensional definitions obtained using A-MAIL versus those obtained using

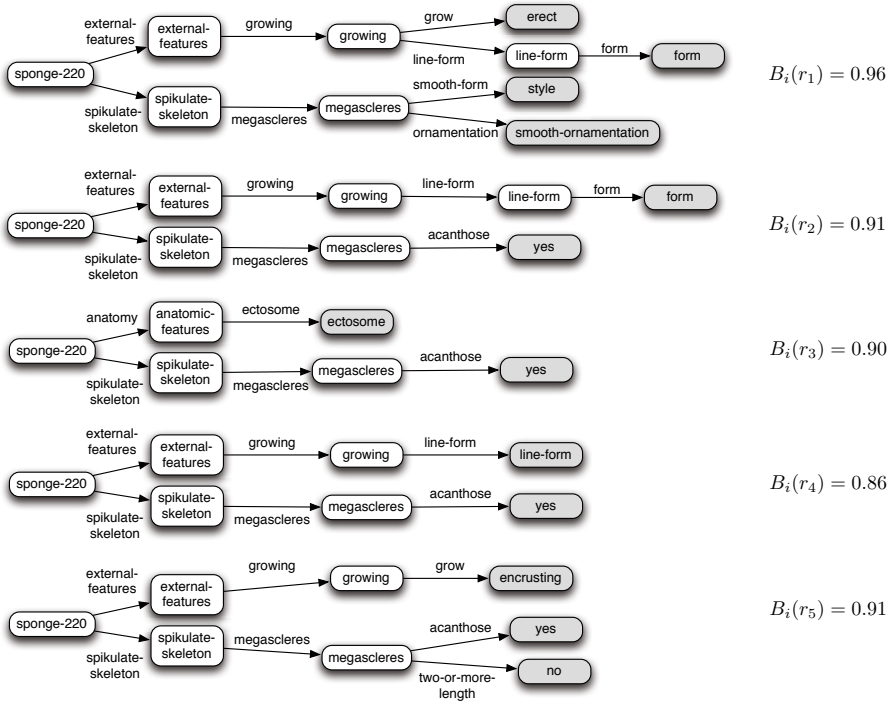
<i>C</i>	<i>Centralized</i>		<i>Individual</i>			<i>A-MAIL</i>		
	P	R	P	R	K	P	R	K
Axinellida	0.98	1.00	0.97	0.95	0.80	0.97	0.95	0.89
Hadromerida	0.85	0.98	0.89	0.91	0.78	0.92	0.96	0.97
Astrophorida	0.98	1.00	0.97	0.97	0.93	0.98	0.99	0.97

different characteristics due to the local habitat conditions. The specific problem we target in these experiments is that of agreeing upon a shared description of the features that distinguish one order of sponges from the others.

To have an idea of the complexity of this problem, Figure 6 shows a description of one of the sponges collected from the Mediterranean sea used in our experiments. As Figure 6 shows, a sponge is defined by five groups of attributes: ecological features, external features, anatomy, features of its spikulate skeleton, and features of its tracts skeleton. Specifically, we used a collection of 280 sponges belonging to three different orders of the demospongiae family: axinellida, Hadromerida and astrophorida. Such sponges were collected from both the Mediterranean sea and Atlantic ocean. In order to evaluate A-MAIL, we used each of the three orders as target concepts for concept convergence —namely Axinellida, Hadromerida and Astrophorida. In an experimental run, we split the 280 sponges randomly among the two agents and, given as target concept one of the orders, the goal of the agents is to reach a convergent definition of such concept. The experiments model the process that two human experts undertake when they to discuss over which features determine whether a sponge belongs to a particular order.

We compared the results of A-MAIL with respect to agents which do not perform argumentation (*Individual*), and to the result of centralizing all the examples and performing centralized induction (*Centralized*). Thus, the difference between the results of *individual* agents and agents using A-MAIL should provide a measure of the benefits of A-MAIL for concept convergence, where as comparing with *Centralized* gives a measure of the quality of the outcome. All the results are the average of 10 executions, with  $\epsilon = 0.95$  and  $\tau = 0.75$ .

Table 1 shows one row for each of the 3 concepts we used in our evaluation: Axinellida, Hadromerida and Astrophorida, and setting we show for them three values: precision (P), recall (R), and convergence degree (K). Precision measures how many of the examples covered by the definition are actually positive examples; recall measures how many of the total number of positive examples in the data set are covered by the definition; and convergence degree is as in Definition 12. The first thing we see in Table 1 is that A-MAIL is able to increase convergence from the initial value appearing in the Individual setting. For two concepts (the exception is Axinellida) the convergence was higher than  $\epsilon = 0.95$ . Total convergence was not reached for because in our experiments  $\tau = 0.75$ , allowing hypotheses to cover some negative examples and preventing overfitting.



**Fig. 7.** Set of rules forming the definition of Axinellida and obtained by one of the agents using A-MAIL in our experiments

This means that acceptable hypotheses can cover some negative examples, and thus generate some divergence. Increasing  $\tau$  could improve convergence but it would make finding hypotheses by induction more difficult, and thus recall might suffer. Moreover, both precision and recall are maintained or improve thanks to argumentation, reaching values close to the ones in a Centralized setting.

Moreover, during argumentation, agents exchanged an average of 10.7 examples to argue about Axinellida, 18.5 for Hadromerida and only 4.1 for Astrophorida. Thus, compared to a centralized approach where all the examples would have to be exchanged, i.e. 280, only a very small fraction of examples are exchanged.

Figure 7 shows the set of rules that one of the agents using A-MAIL obtained in our experiments as the definition of the concept Axinellida. For instance, the first rule states that “all the sponges with an erect and line-form growing, and with megascleres in the spikulate skeleton which had style smooth form and smooth ornamentation belong to the Axinellida order”. By looking at those rules, we can clearly see that both the growing external features and the characteristics of the megascleres are the distinctive features of the Axinellida order.

In summary, we can conclude that A-MAIL successfully achieves concept convergence by integrating argumentation and inductive learning, while maintainig

or improving the quality of the intensional definition (precision and recall). This is achieved by exchanging only a small percentage of the examples the agents know (as opposed to the Centralized setting where all the examples are given to a single agent, which might not be feasible in some applications).

## 5 Related Work

Concerning argumentation in MAS, previous work focuses on several issues like a) logics, protocols and languages that support argumentation, b) argument selection and c) argument interpretation, a recent overview can be found at [14].

The idea that argumentation might be useful for machine learning was discussed in [7], but no concrete proposal has followed, since the authors goal was propose that a defeasible logic approach to argumentation could provide a sound formalization for both expressing and reasoning with uncertain and incomplete information as appears in machine learning. Since the possible hypotheses can be induced from data could be considered an argument, and then by defining a proper attack and defeat relation, a sound hypotheses can be found. However, they did not develop the idea, or attempted the actual integration of an argumentation framework with any particular machine learning technique. Amgoud and Serrurier [1] elaborated on the same idea, proposing an argumentation framework for classification. Their focus is on classifying examples based on all the possible classification rules (in the form of arguments) rather than on a single one learned by a machine learning method.

A related idea is that of *argument-based machine learning* [9], where some examples are augmented with a justification or “supporting argument”. The idea is that those supporting arguments are then used to constrain the search in the hypotheses space: only those hypotheses which classify examples following the provided justification are considered. Notice that in this approach, arguments are used to augment the information contained in an example. A-MAIL uses arguments in a different way. A-MAIL does not require examples to be augmented with such supporting arguments; in A-MAIL the inductive process itself generates arguments. Notice, however, that both approaches could be merged, and that A-MAIL could also be designed to exploit extra information in the form of examples augmented with justifications. Moreover, A-MAIL is a model for multiagent induction, whereas argument-based machine learning is a framework for centralized induction which exploits additional annotations in the examples in the form of arguments.

The idea of using argumentation with case-based reasoning in multiagent systems has been explored by [11] in the AMAL framework. Compared to A-MAIL, AMAL focuses on lazy learning techniques where the goal is to argue about the classification of particular examples, whereas A-MAIL, although uses case bases, allows agents to argue about rules generated through inductive learning techniques. Moreover, the AMAL framework explored a related idea to A-MAIL, namely learning from communication [10]. An approach similar to AMAL is PADUA [15], an argumentation framework that allows agents to use examples

to argue about the classification of particular problems, but they generate association rules and do not perform concept learning.

## 6 Conclusions

The two main contributions of this paper are the definition of an argumentation framework for agents with inductive learning capabilities, and the introduction of the concept convergence task. Since our argumentation framework is based on reasoning from examples, we introduced the idea of *argument acceptability*, which measures how much empirical support an argument has, which is used to define an *attack* relation among arguments. A main contribution of the paper has been to show the feasibility of a completely automatic and autonomous approach to argumentation in empirical tasks. All necessary processes are autonomously performed by artificial agents: generating arguments from their experience, generating attacks to defeat or defend, changing their beliefs as a result of the argumentation process — they are all empirically based and autonomously undertaken by individual agents.

The A-MAIL framework has been applied in this paper to the concept convergence task. However, it can also be seen as a multi-agent induction technique to share inductive inferences [4]. As part of our future work, we want to extend our framework to deal with more complex inductive tasks, such achieving convergence on a collection of interrelated concepts, as well as scenarios with more than 2 agents. Additionally, we would like to explore the use of argumentation frameworks which support weights or strengths in the arguments, in order to take into account the confidence of each agent during the argumentation process.

Our long term goal is to study the relation and integration of inductive inference and communication processes among groups of intelligent agents into a coherent unified MAS framework.

## Acknowledgments

This research was partially supported by projects Next-CBR TIN2009-13692-C03-01 and Agreement Technologies CONSOLIDER CSD2007-0022.

## References

1. Amgoud, L., Serrurier, M.: Arguing and explaining classifications. In: Rahwan, I., Parsons, S., Reed, C. (eds.) ArgMAS 2007. LNCS (LNAI), vol. 4946, pp. 164–177. Springer, Heidelberg (2008)
2. Chesñevar, C.I., Simari, G.R., Godo, L.: Computing dialectical trees efficiently in possibilistic defeasible logic programming. In: Baral, C., Greco, G., Leone, N., Terracina, G. (eds.) LPNMR 2005. LNCS (LNAI), vol. 3662, pp. 158–171. Springer, Heidelberg (2005)
3. Clark, P., Niblett, T.: The CN2 induction algorithm. *Machine Learning*, 261–283 (1989)

4. Davies, W., Edwards, P.: The communication of inductive inferences. In: Weiss, G. (ed.) ECAI 1996 Workshops. LNCS, vol. 1221, pp. 223–241. Springer, Heidelberg (1997)
5. Dietterich, T.G.: Ensemble methods in machine learning. In: Kittler, J., Roli, F. (eds.) MCS 2000. LNCS, vol. 1857, pp. 1–15. Springer, Heidelberg (2000)
6. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artificial Intelligence* 77(2), 321–357 (1995)
7. Gómez, S.A., Chesñevar, C.I.: Integrating defeasible argumentation and machine learning techniques. In: CoRR, cs.AI/0402057 (2004)
8. Mitchell, T.: *Machine Learning*. McGraw-Hill, New York (1997)
9. Mozina, M., Zabkar, J., Bratko, I.: Argument based machine learning. *Artificial Intelligence* 171(10-15), 922–937 (2007)
10. Ontañón, S., Plaza, E.: Case-based learning from proactive communication. In: IJCAI, pp. 999–1004 (2007)
11. Ontañón, S., Plaza, E.: Learning and joint deliberation through argumentation in multiagent systems. In: AAMAS 2007, pp. 971–978 (2007)
12. Ontañón, S., Plaza, E.: Multiagent inductive learning: an argumentation-based approach. In: ICML, pp. 839–846 (2010)
13. Quinlan, J.R.: Learning logical definitions from relations. *Machine Learning* 5, 239–266 (1990)
14. Rahwan, I., Simari, G.R.: *Argumentation in Artificial Intelligence*. Springer, Heidelberg (2009)
15. Wardeh, M., Bench-Capon, T.J.M., Coenen, F.: PADUA: a protocol for argumentation dialogue using association rules. *Artificial Intelligence in Law* 17(3), 183–215 (2009)

# Arguing about Preferences and Decisions<sup>\*</sup>

T.L. van der Weide, F. Dignum, J.-J. Ch. Meyer,  
H. Prakken, and G.A.W. Vreeswijk

Universiteit Utrecht  
{tweide,dignum,jj,henry,gv}@cs.uu.nl

**Abstract.** Complex decisions involve many aspects that need to be considered, which complicates determining what decision has the most preferred outcome. Artificial agents may be required to justify and discuss their decisions to others. Designers must communicate their wishes to artificial agents. Research in argumentation theory has examined how agents can argue about what decision is best using goals and values. Decisions can be justified with the goals they achieve, and goals can be justified by the values they promote. Agents may agree on having a value, but disagree about what constitutes that value. In existing work, however, it is not possible to discuss what constitutes a specific value, whether a goal promotes a value, why an agent has a value and why an agent has specific priorities over goals. This paper introduces several argument schemes, formalised in an argumentation system, to overcome these problems. The techniques presented in this paper are inspired by multi attribute decision theory.

## Categories and Subject Descriptors

I.2.4 [**Artificial Intelligence**]: Knowledge Representation Formalisms and Methods.

**General Terms:** Design.

**Keywords:** Argumentation, Decision-Making, Practical Reasoning, Preferences.

## 1 Introduction

In complex situations, decisions involve many aspects that need to be considered. These aspects are typically different in nature and therefore difficult to compare. This complicates determining what outcome is the most preferable. Throughout the paper, we will use buying a house as an example, which involves many different aspects. For example, an agent may care about the costs

---

<sup>\*</sup> The research reported here is part of the Interactive Collaborative Information Systems (ICIS) project, supported by the Dutch Ministry of Economic Affairs, grant nr: BSIK03024.

of a house, but also about how fun, comfortable, close to shops, and how beautiful a house is. Artificial agents are expected to act in the designer's or user's best interest. This requires designers or users to communicate their wishes to the agent and the agent to explain and discuss why a certain decision was made. A significant amount of research has been concerned with using argumentation theory for decision-making and practical reasoning to determine which decisions are defensible from a given motivation, see for example [7, 2, 1].

A possible argumentation framework for decision-making for this purpose is the one proposed in [1]. Several decision principles are formalised to select the best decision using arguments in favour and against the available decisions. Agents are assumed to have a set of prioritised goals, which are used to construct argument in favour and against decisions. For example, agent  $\alpha$  has the goal to live in a house that is downtown and the less important goal to live in a house bigger than  $60m^2$ . In complex situations it is useful to argue about *what* goals should be pursued. Why have the goal to live downtown and why not in a village? Why is living downtown more important than living in a bigger house? However, justifying and attacking goals is not possible using the framework of [1].

In order to solve this problem, we could use the framework described in [2], where goals can be justified and attacked using the values they promote and demote. People use their values as standards or criteria to guide selection and evaluation of actions [10, 12]. The values of agents reflect their preferences. For example, agent  $\alpha$  has the value of fun and of comfort. The goal to live downtown promotes the value of fun and the goal to live in a bigger house promotes the value of comfort. In [2] an argument scheme is proposed for practical reasoning in which goals and actions can be justified and attacked by the values they promote and demote. What constitutes a specific value like fun, comfort, justice, or health often is disputable and therefore it is also disputable whether goals and actions promote values. Namely, another agent may find that living downtown demotes the value of fun because of the noise and lack of parking space. However, in [2] it is not possible to explain or discuss what constitutes a value and consequently it is also not possible to justify or attack that a goal or action promotes or demotes a value.

This paper presents an argumentation approach to discuss what constitutes a specific value and its effects on agent's goals and preferences over outcomes. To argue about decisions, an argumentation system is described in Section 2. Since the subject of argumentation is making decisions, some basic notions of decision theory are also described in Section 2. Next, we propose a model to specify the meaning of values and their relation to preferences in Section 3. This model is based on previous work [15] and inspired by techniques from decision theory to find an appropriate multi-attribute utility function [8, 9]. A value is seen as an aspect over which an agent has preferences and can be decomposed into the aspects it contains. Given the meaning of a value, several argument schemes are proposed in Section 4 to justify that goals promote or demote values. The introduced formalism is demonstrated with an example of buying houses in Section 5. The paper is concluded with some discussion and conclusions.

## 2 Background

In Section 2.1, an argumentation system is described that will be used to argue about what decision is best. Outcomes describe the effects of decisions and attributes describe properties of outcomes. Attributes of outcomes can be used to describe what constitutes a value and to justify goals. To argue about what decision is best, the notions of outcomes and attributes from decision theory are introduced in our argumentation system in Section 2.2.

### 2.1 Argumentation

Argument schemes are stereotypical patterns of defeasible reasoning [14]. An argument scheme consists of a set of premises, a conclusion, and is associated to a set of critical questions that can be used to critically evaluate the inference. In later sections, argument schemes are proposed to reason about what decision is best.

We introduce an argumentation system to reason defeasibly and in which argument schemes can be expressed. For the largest part, this argumentation system is based on [5]. We will use both defeasible and strict inference rules. The informal reading of a strict inference rule is that if its antecedent holds, then its conclusion holds without exception. The informal reading of a defeasible inference rule is that if its antecedent holds, then its conclusion tends to hold.

**Definition 1 (Argumentation System).** *An argumentation system is a tuple  $\mathcal{AS} = (\mathcal{L}, \mathcal{R})$  with  $\mathcal{L}$  the language of first-order logic and  $\mathcal{R}$  a set of strict and defeasible inference rules.*

We will use  $\phi$  and  $\psi$  as typical elements of  $\mathcal{L}$  and say that  $\phi$  and  $\neg\phi$  are each other's complements. In the meta-language,  $\sim\phi$  denotes the complement of any formula  $\phi$ , positive or negative. Furthermore,  $\rightarrow$  denotes the material implication.

**Definition 2 (Strict and defeasible rules).** *A strict rule is an expression of the form  $s(x_1, \dots, x_n) : \phi_1, \dots, \phi_m \Rightarrow \phi$  and a defeasible rule is an expression of the form  $d(x_1, \dots, x_n) : \phi_1, \dots, \phi_m \rightsquigarrow \phi$ , with  $m \geq 0$  and  $x_1, \dots, x_n$  all variables in  $\phi_1, \dots, \phi_m, \phi$ .*

We call  $\phi_1, \dots, \phi_m$  the antecedent,  $\phi$  the conclusion, and both  $s(x_1, \dots, x_n)$  and  $d(x_1, \dots, x_n)$  the identifier of a rule.

Arguments are inference trees constructed from a knowledge-base  $\mathcal{K} \subset \mathcal{L}$ . If an argument  $A$  was constructed using no defeasible inference rules, then  $A$  is called a *strict argument*, otherwise  $A$  is called a *defeasible argument*.

*Example 1.* Let  $\mathcal{AS} = (\mathcal{L}, \mathcal{R})$  be an argumentation system such that  $\mathcal{L} = \{\phi_1, \phi_2, \phi_3\}$  and  $\mathcal{R} = \{s() : \phi_1 \Rightarrow \phi_2; d() : \phi_1, \phi_2 \rightsquigarrow \phi_3\}$ . From the knowledge-base  $\mathcal{K} = \{\phi_1\}$ , we can construct 3 arguments. Argument  $A_1$  has conclusion  $\phi_1$ , no premises, and no last applied inference rule. Argument  $A_2$  is constructed



by applying  $s()$ . Consequently,  $A_2$  has premise  $A_1$ , conclusion  $\phi_2$ , and last applied inference rule  $s()$ . Argument  $A_3$  can then be constructed using  $d()$  and has premises  $A_1$  and  $A_2$ , conclusion  $\phi_3$ , and last applied rule  $d()$ . Arguments  $A_1$  and  $A_2$  are strict arguments and argument  $A_3$  is a defeasible argument.  $A_3$  can be visualised as follows:

$$\frac{\phi_1 \quad \frac{\phi_1 \quad \phi_2}{\phi_3} s()}{\phi_3} d()$$

All arguments can be attacked by rebutting one of their premises. Defeasible arguments can also be attacked by attacking the application of a defeasible rule. For example, let  $d(c_1, \dots, c_n) : \phi_1, \dots, \phi_m \rightsquigarrow \phi$  be a defeasible inference rule that was applied in argument  $A$ . We can attack  $A$  in three ways: by rebutting a premise of  $A$ , by rebutting  $A$ 's conclusion, and by undercutting a defeasible inference rule that was applied in  $A$ . The application of a defeasible inference rule can be undercut when there is an exception to the rule. An argument concluding  $\sim d(c_1, \dots, c_n)$  undercuts  $A$ .

Following [4], argument schemes are formalised as defeasible inference rules. Critical questions point to counterarguments that either rebut the scheme's premises or undercut the scheme.

Argumentation Frameworks (AF) were introduced by [6] and provide a formal means to determine what arguments are justified given a set of arguments and a set of attack relations between them.

**Definition 3.** An Argumentation Framework (AF) is a tuple  $(\text{Args}, R)$  with  $\text{Args}$  a set of arguments and  $R \subseteq \text{Args} \times \text{Args}$  an attack relation.

**Definition 4 (Conflict-free, Defense).** Let  $(\text{Args}, R)$  be an AF and  $S \subseteq \text{Args}$ .

- $S$  is called conflict-free iff there are no  $A, B \in S$  s.t.  $(A, B) \in R$ .
- $S$  defends argument  $A$  iff for all  $B \in \text{Args}$  s.t.  $(B, A) \in R$ , there is a  $C \in S$  s.t.  $(C, B) \in R$ .

**Definition 5 (Acceptability).** Let  $(\text{Args}, R)$  be an AF,  $S \subseteq \text{Args}$  conflict-free and  $F : 2^{\text{Args}} \rightarrow 2^{\text{Args}}$  a function s.t.  $F(S) = \{A \mid S \text{ defends } A\}$ . Then:

- $S$  is admissible if  $S \subseteq F(S)$ .
- $S$  is a complete extension iff  $S = F(S)$ .
- $S$  is a grounded extension iff  $S$  is the smallest (w.r.t. set inclusion) complete extension,
- $S$  is a preferred extension iff  $S$  is a maximal (w.r.t. set inclusion) complete extension.
- $S$  is a stable extension iff  $S$  is a preferred extension that attacks all arguments in  $\text{Args} \setminus S$ .

## 2.2 Outcomes and Attributes

The notion of *outcomes* is one of the main notions in decision theory [8, 11] and is used to represent the possible consequences of an agent's decisions. The

set  $\Omega$  of possible outcomes should distinguish all consequences that matter to the agent and are possibly affected by its actions. Agents have preferences over outcomes and decision theory postulates that a rational agent should make the decision that leads to the most preferred expected outcome.

The notion of *attribute* is used to denote a feature, characteristic or property of an outcome. For example, when buying a house, relevant attributes could be price, neighbourhood in which it is located, size, or type of house. An attribute has a domain of ‘attribute-values’ outcomes can have. Every outcome has exactly one attribute-value of each attribute. It cannot be that an outcome has two attribute-values of the same attribute.

*Example 2 (Buying a house).* There are 2 houses on the market and buying one of them results in one of the two outcomes  $\Omega = \{\omega_1, \omega_2\}$ . Consider the attributes ‘price’, ‘size’, ‘neighbourhood’, and whether there is a garden. Price is expressed in dollar and size in  $m^2$ . The neighbourhood can either be ‘downtown’ or ‘suburb’ and ‘yes’ represents there is a garden and ‘no’ that there is not.

Outcome  $\omega_1$  has the following attribute-values: price is 150.000, size is 50, neighbourhood is ‘suburb’ and garden is ‘yes’. On the other, outcome  $\omega_2$ ’s price is 200.000, size is also 50, neighbourhood is ‘downtown’ and garden is ‘no’.

Each attribute is a term in  $\mathcal{L}$  and we use  $\mathcal{A}$  to denote the set containing all attributes. If  $x$  is an attribute, we will also say  $x$ -values instead of the attribute values of attribute  $x$ . We define several functions concerning attributes and outcomes.

- The function  $\text{domain}(x)$  returns a set of attribute-values that the attribute  $x$  can have. For example, let attribute `nbhd` denote the neighbourhood of a house, then  $\text{domain}(\text{nbhd}) = \{\text{downtown}, \text{suburb}\}$  or for the attribute price,  $\text{domain}(\text{price}) = \mathbb{R}^+$ .
- For each attribute  $x$ , the function  $\bar{x} : \Omega \rightarrow \text{domain}(x)$  gives the attribute-value of the given outcome for the attribute  $x$ . For example,  $\overline{\text{price}}(\omega_1) = 150.000$ .

*Example 3.* Suppose that  $\Omega = \{\omega_1, \omega_2\}$  is true,  $x$  is an attribute and  $\text{domain}(x) = \{1, 2, 3\}$ . In that case, the function  $\bar{x}$  returns the following:  $\bar{x}(\omega_1) = 3$  and  $\bar{x}(\omega_2) = 1$ .

### 3 Justification of Preferences over Outcomes

Preferences can be expressed in terms of outcomes, e.g. outcome  $A$  is preferred to outcome  $B$ . The more aspects are involved, the more difficult it becomes to directly express preferences over outcomes. Luckily, it is also natural to express preferences in terms of attributes of outcomes. For example, maximising the attribute profit is preferred. From such statements, preferences over outcomes can be justified, e.g. outcome  $A$  is preferred to outcome  $B$  because the  $A$ ’s profit is higher. Typically, outcomes have many attributes, yet agents care only about

a subset. What set of attributes an agent cares about determines the preferences over outcomes. Using argumentation, agents can discuss why certain attributes should and others should not be used.

Justification for a preference statement like “agent  $\alpha$  prefers living downtown to living in a suburb”, is useful to better understand  $\alpha$ ’s preferences. Namely,  $\alpha$  could argue that the centrality of a house positively influences the amount of fun of a house and that  $\alpha$  wants to maximise fun. If it is better understood why  $\alpha$  prefers something, then one could disagree (centrality is not fun because it is very noisy) and give alternatives perspectives (living near nature is also fun and do you not also care about quietness).

In complex situations, the preferences of agents depend on multiple attributes. By decomposing an agent’s preferences into the different aspects it involves, the number of attributes an aspect depends on becomes smaller. By recursively decomposing preferences, we will arrive at aspects that depend on a single attribute. For example, an agent  $\alpha$  decomposes its preferences concerning houses into the aspects costs and comfort. The perspective of costs is determined by the attribute acquisition price. Comfort however, depends on multiple aspects. Therefore, comfort is decomposed into location and size. Location is then connected to the attribute neighbourhood and size to the surface area of the house. One may argue that  $\alpha$  forgets that other aspects also influence costs, e.g. maintenance, taxes, heating costs, and so on. On the other hand, another agent may decompose comfort differently. For example, for agent  $\beta$  comfort is influenced by the closeness to highway and whether there is central heating.

In Section 3.1 we will introduce *perspectives* to represent preferences and aspects of preferences, after which we introduce perspectives on attribute values in Section 3.2. In Section 3.3 we introduce influence between perspective to denote that one perspective is an aspect of another. Finally in Section 3.4, we slightly adapt Value-based Argumentation Frameworks, see [3], to determine what conclusions are justified to make.

### 3.1 Perspectives

An ordering over outcomes can represent an agent’s preferences. In that case, if an outcome is higher in the order, the agent prefers that outcome. Similarly, outcomes can be ordered according to some criterion. For example, outcomes can be ordered by how fun they are, or how fair they are. To talk about these different orderings, we introduce the notion of *perspective*. With buying houses, an outcome may be better than another from the perspective of costs, worse from the perspective of its centrality, indifferent from the perspective of comfort, and perhaps incomparable from the perspective of fun.

**Definition 6 (Perspective).** *A perspective  $p$  is associated with a preorder,  $\leq_p$  over outcomes  $\Omega$ . The set  $\mathcal{P}$  denotes the set of all perspectives.*

In other words, a perspective  $p$  is associated to a binary relation  $\leq_p \subseteq \Omega \times \Omega$  that is transitive and reflexive. If  $\omega_1 \leq_p \omega_2$  is true, we say that  $\omega_2$  is weakly preferred

to  $\omega_1$  from perspective  $p$ . Strong preference from perspective  $p$  is denoted as  $\omega_1 <_p \omega_2$  and stands for  $\omega_1 \leq_p \omega_2$  and  $\omega_2 \not\leq_p \omega_1$ . Equivalence from perspective  $p$  is denoted as  $\omega_1 \approx_p \omega_2$  and stands for  $\omega_1 \leq_p \omega_2$  and  $\omega_2 \leq_p \omega_1$ .

Each agent  $\alpha$  is associated with a perspective  $\hat{\alpha}$  representing  $\alpha$ 's preferences over outcomes. If  $\omega_1 <_{\hat{\alpha}} \omega_2$  is true, then we either say that  $\omega_2$  is preferred to  $\omega_1$  from agent  $\alpha$ 's perspective, or we say that  $\alpha$  prefers  $\omega_2$  to  $\omega_1$ . Since perspectives are the main notion in this paper,  $\hat{\alpha}$  is abbreviated to  $\alpha$ , so that  $\alpha$  denotes a perspective.

Not only the preferences of agents can be represented with perspectives, we will also use perspectives to represent aspects of outcomes and the values of agents. For example, the value of 'safety' is represented with a perspective that orders outcomes according to how safe they are or the aspect of comfort is represented with a perspective that orders outcomes by the amount of comfort.

*Example 4.* Agent  $\alpha$  wants to buy a new house and wants to minimise costs, maximise fun and maximise comfort. In Figure 1a, we sketch how the preferences of agent  $\alpha$  can be decomposed and how attributes of outcomes can be assigned. Costs are determined by the attribute acquisition price. Fun is influenced by the centrality of the house, i.e. the more central the neighbourhood, the more fun it is. Comfort is influenced by how quiet it is and the size of the house. Again, the attribute neighbourhood is ordered but now by how quiet the neighbourhood is. The size is determined by the surface area of the house.

Agent  $\beta$  just won the lottery and does not care about costs. To  $\beta$  fun is being close to nature, which is completely different from  $\alpha$ 's idea about fun. Also, because  $\beta$  has a car and  $\alpha$  does not,  $\beta$  cares about whether there is enough parking space in the area. Figure 1b sketches the decomposition of  $\beta$ 's preferences.

### 3.2 Attributes Determine Perspectives

Attributes can be used to determine how outcomes should be ordered from a perspective. For example, if you want to order houses from the perspective of size, then the attribute 'surface area of the house' is an appropriate attribute.

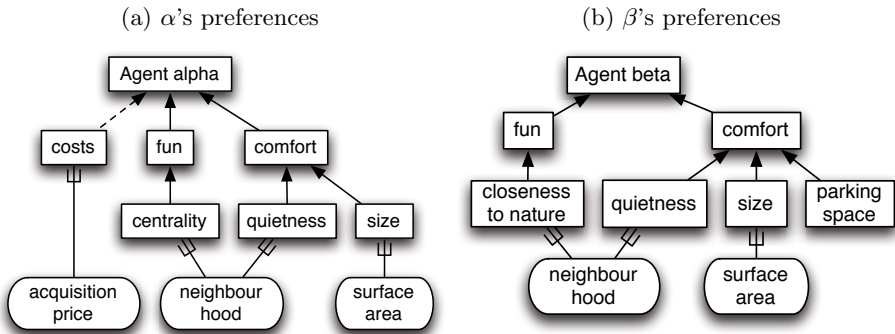


Fig. 1. Different Perspectives Using Different Attributes

In that case, if house A has a higher surface area than house B, then A is preferred to B from the perspective of size. To use an attribute  $x$  to determine a perspective,  $x$ 's attribute values need to be ordered.

**Definition 7 (Attribute Perspective).** *An attribute perspective  $p_x$  is a perspective that is associated with a partial preorder  $\preceq_x^p$  over the domain of attribute  $x$ .*

Note that there can be different attribute perspectives on the same attribute.

*Example 5.* Let the attribute `nbhd` denote the neighbourhood of the house with  $\text{domain}(\text{nbhd}) = \{\text{dwntwn}, \text{vllg}, \text{sbrb}\}$ . Furthermore, let  $\text{social}_{\text{nbhd}}$  and  $\text{quiet}_{\text{nbhd}}$  be attribute perspectives denoting the sociableness and the quietness of the neighbourhood respectively. The different neighbourhoods are preferred from each attribute perspective as follows:

$$\text{sbrb} \prec_{\text{nbhd}}^{\text{social}} \text{vllg} \prec_{\text{nbhd}}^{\text{social}} \text{dwntwn} \qquad \text{dwntwn} \prec_{\text{nbhd}}^{\text{quiet}} \text{sbrb} \prec_{\text{nbhd}}^{\text{quiet}} \text{vllg}$$

If an attribute value is preferred to another attribute value from  $p_x$ , then outcomes with the preferred attribute value are preferred from  $p_x$ . This order between outcomes from an attribute perspective  $p_x$  can be inferred with the following strict inference rule.

$$s_{ap}(p_x, \omega_1, \omega_2) : \bar{x}(\omega_1) \prec_x^p \bar{x}(\omega_2) \Rightarrow \omega_1 <_{p_x} \omega_2$$

*Example 6.* Consider the attributes and attributes perspectives from the previous example. Let there be two outcomes  $\omega_1$  and  $\omega_2$  such that  $\overline{\text{nbhd}}(\omega_1) = \text{dwntwn}$  and  $\overline{\text{nbhd}}(\omega_2) = \text{vllg}$ . To determine the order between  $\omega_1$  from perspective  $\text{social}_{\text{nbhd}}$  and from perspective  $\text{quiet}_{\text{nbhd}}$ , the following two arguments can be constructed.

$$\frac{\overline{\text{nbhd}}(\omega_2) \prec_{\text{nbhd}}^{\text{social}} \overline{\text{nbhd}}(\omega_1)}{\omega_2 <_{\text{social}_{\text{nbhd}}} \omega_1} \quad s_{ap} \qquad \frac{\overline{\text{nbhd}}(\omega_1) \prec_{\text{nbhd}}^{\text{quiet}} \overline{\text{nbhd}}(\omega_2)}{\omega_1 <_{\text{quiet}_{\text{nbhd}}} \omega_2} \quad s_{ap}$$

### 3.3 Influence between Perspectives

Some perspectives involve different aspects such that not one attribute can be assigned. For example, the perspective `comfort` of a house may involve size, location, the type of heating, and so on. In general, the more abstract a perspective is, the more aspects it has. Furthermore, the more abstract a perspective, the more disputable it may be. Thus it becomes important to specify all the different aspects so that one can communicate clearly.

By decomposing an abstract perspective into several more concrete perspectives, one makes explicit what an abstract perspective means and makes it easier to assign attributes to. For example, although  $\alpha$  may not be able to express its preferences over houses,  $\alpha$  does want to minimise costs, maximise comfort, and maximise fun. These perspectives may be decomposed further, e.g. fun means maximising the centrality of the house, until an attribute can be assigned.

To decompose a perspective into other perspectives, we introduce two relations between perspectives in  $\mathcal{L}$  to denote ‘influence’ between perspectives:

- the binary relation  $\uparrow \subseteq \mathcal{P} \times \mathcal{P}$  is written as  $p \uparrow q$  and denotes that perspective  $p$  *positively influences* perspective  $q$ .
- the binary relation  $\downarrow \subseteq \mathcal{P} \times \mathcal{P}$  is written as  $p \downarrow q$  and denotes that perspective  $p$  *negatively influences* perspective  $q$ .

If perspective  $p$  positively influences perspective  $q$ , then outcomes that are better from perspective  $p$  tend to be better from perspective  $q$ . In other words, if an outcome is better from  $p$  and  $p$  positively influences  $q$ , then this is a reason to believe that the outcome is better from  $q$ . For example, the size of a house positively influences the comfort of the house, i.e. the more size, the more comfort.

The following argument scheme reasons with influence: *outcome  $\omega_2$  is preferred to  $\omega_1$  from  $p$  and  $p$  positively influences  $q$ , therefore  $\omega_2$  is preferred to  $\omega_1$  from  $q$* . In other words, if a perspective  $p$  positively influences perspective  $q$ , then an outcome being preferred from  $p$  is a reason to believe that that outcome is also preferred from  $q$ . We formalise this argument scheme with the following defeasible inference rule:

$$d_{<, \uparrow}(p, q, \omega_1, \omega_2) : \omega_1 <_p \omega_2, p \uparrow q \rightsquigarrow \omega_1 <_q \omega_2$$

The argument scheme to propagate negative influence is similar, except that if a perspective  $p$  negatively influences perspective  $q$ , then outcomes that are better from perspective  $p$  tend to be worse from perspective  $q$ . For example, costs negatively influences agent  $\alpha$ 's preferences, i.e. the more costs, the less preferable for  $\alpha$ . This argument scheme is formalised with the following defeasible inference rule:

$$d_{<, \downarrow}(p, q, \omega_1, \omega_2) : \omega_1 <_p \omega_2, p \downarrow q \rightsquigarrow \omega_2 <_q \omega_1$$

Using these inference rules, arguments can be constructed to justify preferences over outcomes using both positive influence and negative influence between perspectives.

*Example 7.* Agent  $\alpha$  wants to buy a new house and to minimise costs, i.e.  $\text{costs} \downarrow \alpha$  is true. There are two outcomes,  $\omega_1$  and  $\omega_2$ , such that the acquisition price, denoted attribute  $\text{acq}$ , of  $\omega_1$  is \$200k and of  $\omega_2$  \$150k. The acquisition price of a house positively influences the costs, so  $\text{costs}_{\text{acq}} \uparrow \text{costs}$  is true.

$$A = \frac{\text{costs} \downarrow \alpha}{\omega_1 <_{\alpha} \omega_2} \frac{\frac{\text{costs}_{\text{acq}} \uparrow \text{costs}}{\omega_2 <_{\text{costs}_{\text{acq}}} \omega_1} \frac{\overline{\text{acq}}(\omega_2) \prec_{\text{acq}}^{\text{costs}} \overline{\text{acq}}(\omega_1)}{s_{ap}}}{d_{<, \downarrow}} d_{<, \uparrow}$$

The relation  $\uparrow$  is irreflexive and transitive, meaning that  $p \uparrow p$  is never true and that if  $p \uparrow q$  and  $q \uparrow r$  are true, then  $p \uparrow r$  is true. If  $p \uparrow p$  would be true and for any two outcomes  $\omega_1 <_p \omega_2$  is true, then we can defeasibly infer that  $\omega_1 <_p \omega_2$  is true using inference rule  $d_{<, \uparrow}$ . Such an argument concludes one of its premises, which is useless. Furthermore, if  $p \uparrow p$  can be true, then this may cause infinite loops in implementations. For this reason, the relation  $\uparrow$  is irreflexive.

The relation  $\uparrow$  should be transitive. Firstly, this is intuitive. For example, let the location of a house positively influence the fun of that house and let the fun of a house positively influence agent  $\alpha$ 's preferences. Then we can also say that the location of a house positively influences  $\alpha$ 's preferences. Secondly, this leads to inferences we could already infer. Namely, if  $p \uparrow q$ ,  $q \uparrow r$  and  $\omega_1 <_p \omega_2$  are true, then we can defeasibly infer that  $\omega_1 <_q \omega_2$  is true. Similarly,  $\omega_1 <_r \omega_2$  can be defeasibly inferred from  $q \uparrow r$  and  $\omega_1 <_q \omega_2$ . If  $\uparrow$  is transitive, then  $p \uparrow r$  is also true. From  $p \uparrow r$  and  $\omega_1 <_p \omega_2$  we can defeasibly infer that  $\omega_1 <_r \omega_2$  is true.

The relation  $\downarrow$  is irreflexive and antitransitive, meaning that  $p \downarrow p$  is not allowed and that if  $p \downarrow q$  and  $q \downarrow r$  are true, then  $p \downarrow r$  is not true. If for some perspective  $p$  it would be true that  $p \downarrow p$  and for any two outcomes  $\omega_1 <_p \omega_2$  is true, then the following argument  $A$  can be constructed.

$$A = \frac{p \downarrow p \quad \omega_1 <_p \omega_2}{\omega_2 <_p \omega_1} d_{<, \downarrow}$$

The conclusion of  $A$  conflicts with  $A$ 's premise  $\omega_1 <_p \omega_2$ . Consequently,  $A$  attacks itself. Allowing  $p \downarrow p$  to be true adds nothing useful and can only result in contradictions. Therefore, the relation  $\downarrow$  is irreflexive.

The relation  $\downarrow$  should be antitransitive. Firstly, this is intuitive. For example, the amount of discount on a house negatively influences the costs of that house (the more discount the less costs) and the costs of a house negatively influences agent  $\alpha$ 's preferences. Then we can also say that the amount of discount does not negatively influence  $\alpha$ 's preferences. Secondly, if  $\downarrow$  could be transitive, then this could lead to false inferences. Let  $p \downarrow q$ ,  $q \downarrow r$  and  $\omega_1 <_p \omega_2$  be true. From  $p \downarrow q$  and  $\omega_1 <_p \omega_2$  we can infer that  $\omega_2 <_q \omega_1$  is true and from  $\omega_2 <_q \omega_1$  and  $q \downarrow r$  we can infer that  $\omega_1 <_r \omega_2$  is true. If  $p \downarrow r$  would be true, then we could infer that  $\omega_2 <_r \omega_1$  is true, which conflicts with  $\omega_1 <_r \omega_2$ .

### 3.4 Acceptability of Arguments

In Section 2.1 different acceptability semantics are given for argumentation frameworks. Like agents have preferences over values in [3], it is natural to assume that agents have preferences over perspectives. The preferences over perspectives should determine whether attacks between arguments based on those perspectives should be successful. For this reason, this section introduces a defeat relation based on agent's preferences between sets of perspectives.

Let the function  $\eta : \text{Args} \rightarrow 2^{\mathcal{P}}$  map an argument to the set of perspectives that are used in that argument. For example,  $\eta(A) = \{\text{costs}_{\text{acq}}, \text{costs}, \alpha\}$  in Example 7. Using what perspectives an argument uses and an agent's preferences over sets of perspectives, we will now define a defeat relation between arguments with respect to an agent.

**Definition 8 (Defeat).** *Let  $\text{Args}$  be a set of arguments and  $\preceq_{\alpha} \subseteq 2^{\mathcal{P}} \times 2^{\mathcal{P}}$  agent  $\alpha$ 's preferences over sets of perspectives. Argument  $A \in \text{Args}$   $\alpha$ -defeats argument  $B \in \text{Args}$  iff  $A$  undercuts  $B$  or  $A$  rebuts  $B$  on  $B'$  and not  $\eta(A) \preceq_{\alpha} \eta(B')$ .*

Let  $AF = (\text{Args}, R)$  be an argumentation framework. If the arguments in  $\text{Args}$  are constructed using the set of perspectives  $\mathcal{P}$  and the argumentation system as introduced in this section, then  $AF$  is said to be based on perspectives  $\mathcal{P}$ . Furthermore, if  $AF$  is based on  $\mathcal{P}$ ,  $\alpha$  an agent and  $R$  the  $\alpha$ -defeat relation as defined in Definition 8, then we say that  $AF$  corresponds to agent  $\alpha$ . In that case,  $AF$  can be used to determine what arguments  $\alpha$  finds acceptable.

Note that because different agents care differently about perspectives, they have different preference orderings over sets of perspectives. Since the defeat relation is built using an agent's preferences over sets of perspectives, the defeat relation  $R$  therefore is subjective. What argument is acceptable is thus also subjective.

## 4 Justification of Goals

In this section, we propose how an agent  $\alpha$  can justify having a goal given  $\alpha$ 's preferences. In [13], Simon views goals as threshold aspiration levels that signal satisfactory of utility. A goal thus does not have to be optimal. Following [16], we see goals as expressions of the desirability of attribute values of a single attribute signaling that these attribute values are 'satisfactory'. For example, an agent may have the goal to live in a house that is located downtown. This expresses that the attribute value 'downtown' of the attribute 'location' is satisfactory to that agent. Another attribute value, e.g. 'suburb', does not achieve that goal and is thus not satisfactory.

The predicate  $\text{goal}(\alpha, x, G)$  is introduced in  $\mathcal{L}$  and denotes that agent  $\alpha$  should have the goal to achieve an outcome that has an  $x$ -value in  $G \subset \text{domain}(x)$ . If agent  $\alpha$  has the goal to achieve an  $x$ -value in  $G$  (i.e.  $\text{goal}(\alpha, x, G)$  is true) and outcome  $\omega_1$  has an  $x$ -value in  $G$ , i.e.  $\bar{x}(\omega_1) \in G$  is true, then we say that goal  $\text{goal}(\alpha, x, G)$  is achieved in outcome  $\omega_1$ . Consequently, a subset of  $\Omega$  achieves a goal and the other outcomes in  $\Omega$  do not achieve that goal.

**Justification Is Subjective.** What justification for a goal an agent accepts, depends on the type of agent. For example, a very ambitious but realistic agent only accepts goals that aim for the best achievable  $x$ -value, whereas a less ambitious agent may accept goals that just improves the current situation or does better than doing nothing. Another agent may set its standard on a value that is realistic and challenging, i.e. not too easy and not too difficult.

We introduce two argument schemes to distinguish between satisficing goals and optimising goals. The following argument scheme justifies the goal to achieve an  $x$ -value that is the best possible. The basis for this justification is that agents should aim to achieve their maximal potential.

*Agent  $\alpha$  wants to maximise attribute  $x$ -values from perspective  $p_x$ ,  
 $v$  is most preferred  $x$ -value from  $p_x$  that is achievable,  
 therefore,  $\alpha$  pursues the goal to achieve  $x$ -values of  $v$  or better from  $p_x$*

If the predicates  $\text{max}(\alpha, p_x, v)$  and  $\text{min}(\alpha, p_x, v)$  denote that  $v$  is the maximal / minimal  $x$ -value from  $p_x$  that  $\alpha$  can achieve, then the optimistic goal argument



scheme can be modelled with the following defeasible inference rules:

$$d_{\text{optim},\uparrow}(\alpha, p_x, v) : p_x \uparrow \alpha, \max(\alpha, v, p_x) \rightsquigarrow \text{goal}(\alpha, x, \{g \in \text{domain}(x) \mid v \preceq_x^p g\})$$

$$d_{\text{optim},\downarrow}(\alpha, p_x, v) : p_x \downarrow \alpha, \min(\alpha, v, p_x) \rightsquigarrow \text{goal}(\alpha, x, \{g \in \text{domain}(x) \mid g \preceq_x^p v\})$$

A possible undercutter of the optimistic argument scheme is that achieving the goal is too unlikely. Therefore, the agent should adopt the goal to achieve an easier  $x$ -value. Another undercutter would be that achieving  $v$  is too costly and that  $\alpha$  does not care that much about  $p_x$ .

The following argument scheme justifies a goal in a satisficing manner. This scheme's underlying motivation is that agents should adopt goals that achieve outcomes that are satisfactory rather than the best outcome.

*Agent  $\alpha$  wants to maximise attribute  $x$ -values from perspective  $p_x$ ,  
 $v$  is a satisfactory and achievable  $x$ -value for  $\alpha$ ,  
therefore,  $\alpha$  pursues the goal to achieve  $x$ -values of  $v$  or better from  $p_x$*

A possible undercutter for the satisficing argument scheme is that it is too easy and that the agent should adopt a more challenging goal. Another undercutter could be that the perspective  $p_x$  is important to  $\alpha$  and therefore  $\alpha$  should set a higher goal.

Let the predicate  $\text{satisf}(\alpha, x, v)$  denote that  $x$ -value  $v$  is satisfactory for agent  $\alpha$ . Then this argument scheme can be modelled with the following defeasible inference rule.

$$d_{\text{satisf},\uparrow}(\alpha, p_x, v) : p_x \uparrow \alpha, \text{satisf}(\alpha, x, v) \rightsquigarrow \text{goal}(\alpha, x, \{g \in \text{domain}(x) \mid v \preceq_x^p g\})$$

$$d_{\text{satisf},\downarrow}(\alpha, p_x, v) : p_x \downarrow \alpha, \text{satisf}(\alpha, x, v) \rightsquigarrow \text{goal}(\alpha, x, \{g \in \text{domain}(x) \mid g \preceq_x^p v\})$$

This only solves part of the problem because how can an agent justify that an attribute value is satisfactory? We can think of several justifications of a satisfaction level: anything better than the current situation is satisfactory, it is better than some standard action such as 'do nothing', it is better than what other agents achieve, or the agent is obliged to achieve at least  $v$ . This is however still an open issue that is left for future work.

**Priorities Of Goals.** As explained in Section 3.4, agents have preferences over sets of perspectives. This information can be used to give goals priorities. Namely suppose agent  $\alpha$  has goal  $G$  because there is an 'influence-path' between  $p_x, p_1, \dots, p_n, \alpha$  (i.e.  $p_x$  influences  $p_1$ ,  $p_1$  influences  $p_2$ , and so on) and agent  $\alpha$  has goal  $H$  because of an influence path between  $q_x, q_1, \dots, q_m, \alpha$ . If agent  $\alpha$  prefers the set of perspectives  $\{p_x, p_1, \dots, p_n, \alpha\}$  to the set of perspectives  $\{q_x, q_1, \dots, q_m, \alpha\}$ , then  $\alpha$  should give a higher priority to goal  $G$  than to goal  $H$ .

Goals are created using an attribute perspective that influences an agent. For the same attribute perspective, optimistic goals are stricter than satisficer goals since they do not include satisfactory attribute values upon which the agent can improve. For this reason, achieving an optimistic goal should have a higher priority than achieving a satisficer goal for the same attribute perspective.

## 5 Buying a House

Agent  $\alpha$ , who lives in a suburb, recently got a raise in income and wants to buy a new house to live in. The broker shows two houses that are for sale, one in a village and one downtown, represented with outcomes  $\omega_v$  and  $\omega_d$  respectively. Of course,  $\alpha$  has the possibility not to buy a new house and stay in its current house. This is represented with outcome  $\omega_0$ . Consequently,  $\Omega = \{\omega_0, \omega_d, \omega_v\}$ . Except for its own house,  $\alpha$  is unfamiliar with these houses and can therefore not express whether it prefers one of the new houses to its own house.

The broker includes the following attributes of each house: the neighbourhood, the size, and the acquisition price. The attribute `nbhd` denotes the neighbourhood of the house and  $\text{domain}(\text{nbhd}) = \{\text{dwntwn}, \text{sbrb}, \text{vllg}\}$ . The attribute `area` denotes how big the house is in square meters. Consequently,  $\text{domain}(\text{area}) = \mathbb{R}^+$ . The attribute `acq` denotes the price of the acquisition of the house and  $\text{domain}(\text{acq}) = \mathbb{R}^+$ . The set of all attributes is the following:  $\mathcal{A} = \{\text{nbhd}, \text{area}, \text{acq}\}$ . The attribute values for each outcome can be found in Table [1](#).

**Table 1.** Attribute Values Of Outcomes

Attribute	Domain	$\omega_0$	$\omega_d$	$\omega_v$
<code>nbhd</code>	$\{\text{dwntwn}, \text{sbrb}, \text{vllg}\}$	<code>sbrb</code>	<code>dwntwn</code>	<code>vllg</code>
<code>area</code>	$\mathbb{R}^+$ in $m^2$	60	50	100
<code>acq</code>	$\mathbb{R}^+$ in \$1000	0	220	190

### 5.1 Decomposing Perspectives

Agent  $\alpha$  starts reasoning about its preferences over  $\Omega$  by expressing what aspects it finds important. Namely,  $\alpha$  wants to minimise costs, maximise fun and maximise comfort. By doing so,  $\alpha$ 's perspective is decomposed into other perspectives that are more concrete. Because  $\alpha$  wants to minimise costs, the perspective `costs` negatively influences the perspective of  $\alpha$ , i.e. `costs`  $\downarrow$   $\alpha$  is true. Also,  $\alpha$  wants to maximise fun and comfort, so `fun`  $\uparrow$   $\alpha$  and `comfort`  $\uparrow$   $\alpha$  are true.

Agent  $\alpha$  figures that the acquisition price attribute is appropriate to determine the perspective of `costs` such that the higher the acquisition price, the higher the costs. The attribute perspective `costsacq` prefers an `acq`-value if it is higher. Therefore, `costsacq`  $\uparrow$  `costs` is true.

For  $\alpha$  fun means having people around him. The centrality of a house positively influences fun since  $\alpha$  is more likely to out for dinner or drinks with his friends. Therefore,  $\alpha$  decomposes the perspective `fun` into the perspective of the centrality of the neighbourhood, denoted with the attribute perspective `cntrlnbhd` on the attribute `nbhd`. Consequently, `cntrlnbhd`  $\uparrow$  `fun` is true.

There is however no attribute that  $\alpha$  finds adequate to determine the perspective of comfort. Therefore, `comfort` is decomposed into the quietness around the house and its size. Size is measured by the attribute perspective `sizearea`

that orders the attribute *area* (denoting the surface area in  $m^2$ ) according to size. The attribute perspective  $\text{quiet}_{\text{nbhd}}$  orders neighbourhoods by their quietness. Both attributes positively influence comfort, i.e.  $\text{size}_{\text{area}} \uparrow \text{comfort}$  and  $\text{quiet}_{\text{nbhd}} \uparrow \text{comfort}$  are true.

The attribute perspectives  $\text{cntrl}_{\text{nbhd}}$  and  $\text{quiet}_{\text{nbhd}}$  both order the attribute values of the attribute ‘neighbourhood’ and are as follows:

$$\text{sbrb} \prec_{\text{nbhd}}^{\text{cntrl}} \text{vllg} \prec_{\text{nbhd}}^{\text{cntrl}} \text{downtwn} \quad \text{downtwn} \prec_{\text{nbhd}}^{\text{quiet}} \text{sbrb} \prec_{\text{nbhd}}^{\text{quiet}} \text{vllg}$$

## 5.2 Arguments about Preference

Now,  $\alpha$  starts constructing arguments concerning its preferences over houses. The following argument concludes that  $\alpha$  should prefer staying in its house, outcome  $\omega_0$ , to buying the house downtown,  $\omega_d$ , because the costs of not buying are lower.

$$A_{\text{costs}} = \frac{\text{costs} \downarrow \alpha \quad \frac{\overline{\text{acq}}(\omega_0) \prec_{\text{acq}}^{\text{costs}} \overline{\text{acq}}(\omega_d) \quad s_{ap}}{\omega_0 \prec_{\text{costs}_{\text{acq}}} \omega_d}}{\omega_d \prec_{\alpha} \omega_0} d_{\downarrow, <}$$

However, the following argument concludes that  $\alpha$  should prefer  $\omega_d$ , which conflicts with  $A_{\text{costs}}$ ’s conclusion, because  $\omega_d$  is more fun since it is located in a more central neighbourhood.

$$A_{\text{fun}} = \frac{\text{fun} \uparrow \alpha \quad \frac{\text{cntrl}_{\text{nbhd}} \uparrow \text{fun} \quad \frac{\overline{\text{nbhd}}(\omega_0) \prec_{\text{nbhd}}^{\text{cntrl}} \overline{\text{nbhd}}(\omega_d) \quad s_{ap}}{\omega_0 \prec_{\text{cntrl}_{\text{nbhd}}} \omega_d}}{\omega_0 \prec_{\text{fun}} \omega_d}}{\omega_0 \prec_{\alpha} \omega_d} d_{\uparrow, <}$$

Agent  $\alpha$  keeps thinking and comes up with the following argument that concludes that its current house is actually more comfortable since it is in a neighbourhood that is more quiet.

$$A_{\text{comfort}} = \frac{\text{comfort} \uparrow \alpha \quad \frac{\text{quiet}_{\text{nbhd}} \uparrow \text{comfort} \quad \frac{\overline{\text{nbhd}}(\omega_d) \prec_{\text{nbhd}}^{\text{quiet}} \overline{\text{nbhd}}(\omega_0) \quad s_{ap}}{\omega_d \prec_{\text{quiet}_{\text{nbhd}}} \omega_0}}{\omega_d \prec_{\text{comfort}} \omega_0}}{\omega_d \prec_{\alpha} \omega_0} d_{\uparrow, <}$$

Given these three arguments, we want to determine what conclusions are justified. For this we construct the PerspAF  $\langle H(\text{Args}, R), \mathcal{P}, \eta \rangle$ , with:

$$\begin{aligned} \text{Args} &= \{A_{\text{costs}}, A_{\text{fun}}, A_{\text{comfort}}\} \\ R &= \{(A_{\text{costs}}, A_{\text{fun}}), (A_{\text{fun}}, A_{\text{costs}}), (A_{\text{comfort}}, A_{\text{fun}}), (A_{\text{fun}}, A_{\text{comfort}})\} \\ \mathcal{P} &= \{\alpha, \text{fun}, \text{comfort}, \text{costs}, \text{quiet}_{\text{nbhd}}, \text{cntrl}_{\text{nbhd}}, \text{costs}_{\text{acq}}, \text{size}_{\text{area}}\} \end{aligned}$$

and function  $\eta$ , that maps an argument to the perspectives it contains, is as follows.

$$\begin{aligned} \eta(A_{\text{costs}}) &= \{\alpha, \text{costs}, \text{costs}_{\text{acq}}\} \\ \eta(A_{\text{fun}}) &= \{\alpha, \text{fun}, \text{cntrl}_{\text{nbhd}}\} \\ \eta(A_{\text{comfort}}) &= \{\alpha, \text{comfort}, \text{quiet}_{\text{nbhd}}\} \end{aligned}$$

Let  $\alpha$  find fun more important than comfort and costs more important than fun, i.e.  $\text{comfort} \triangleleft_{\alpha} \text{fun}$  and  $\text{fun} \triangleleft_{\alpha} \text{costs}$  are true. Then  $A_{\text{fun}}$   $\alpha$ -defeats  $A_{\text{comfort}}$  and  $A_{\text{costs}}$   $\alpha$ -defeats  $A_{\text{fun}}$ . Then the set  $\{A_{\text{comfort}}, A_{\text{costs}}\}$  is the preferred extension, so the conclusion that  $\alpha$  prefers staying in its current house to buying a house downtown is justified.

### 5.3 Goals

If  $\alpha$  will also visit other brokers and thus considers more houses, it can be computationally efficient for  $\alpha$  to generate a number of goals that can easily be checked when evaluating a new house. Given a number of goals, evaluating an outcome involves checking whether its attribute values are in the goals. If no goals are used, then evaluating an outcome involves constructing arguments for all relevant perspectives to check whether it is better than some other outcome(s).

The current house of  $\alpha$  is  $60m^2$ , i.e.  $\overline{\text{area}}(\omega_0) = 60$ , and  $\alpha$  finds this size satisfactory. Since  $\alpha$  does not feel very strongly about the size of its house,  $\alpha$  uses the satisfying argument scheme to justify the following goal.

$$\frac{\text{size}_{\text{area}} \uparrow \alpha \quad \text{satisf}(\alpha, \text{area}, 60)}{\text{goal}(\alpha, \text{area}, \{g \in \text{domain}(\text{area}) \mid g \geq 60\})} d_{\text{satisf}}$$

With its new job,  $\alpha$  can maximally lend 200 thousand dollar for the acquisition of a house and therefore  $\alpha$  sets its aspiration level for the acquisition on 200. Given this information,  $\alpha$  justifies having the following goal:

$$\frac{\text{costs}_{\text{acq}} \downarrow \alpha \quad \text{satisf}(\alpha, \text{acq}, 200)}{\text{goal}(\alpha, \text{acq}, \{g \in \text{domain}(\text{acq}) \mid g \leq 200\})} d_{\text{satisf}}$$

The current house of  $\alpha$  is in a suburb and  $\alpha$  wants to maximise neighbourhood with respect to both centrality and quietness, i.e.  $\text{cntrl}_{\text{nbhd}} \uparrow \alpha$  and  $\text{quiet}_{\text{nbhd}} \uparrow \alpha$  are true. Agent  $\alpha$  cares a lot about the centrality of its house and less about its quietness. Therefore,  $\alpha$  uses the optimising argument scheme to justify its goal to live downtown:

$$\frac{\text{cntrl}_{\text{nbhd}} \uparrow \alpha \quad \max(\alpha, \text{dwntwn}, \text{cntrl}_{\text{nbhd}})}{\text{goal}(\alpha, \text{nbhd}, \{\text{dwntwn}\})} d_{\text{optim}}$$

About the quietness  $\alpha$  cares less and therefore uses the satisfying argument scheme:

$$\frac{\text{quiet}_{\text{nbhd}} \uparrow \alpha \quad \text{satisf}(\alpha, \text{nbhd}, \text{sbrb})}{\text{goal}(\alpha, \text{nbhd}, \{\text{sbrb}, \text{vllg}\})}$$

It is impossible for  $\alpha$  to achieve both goals. However,  $\alpha$  finds costs more important than fun and fun more important than comfort. Consequently,  $\alpha$  finds the goal to live downtown more important than the goal to live in a quiet suburb.

## 6 Discussion

**Outcomes Compared To States.** In [3, 2], states are used to reason about decisions over actions, rather than outcomes. A state is a truth assignment to a set of propositions. In a state  $r$ , an agent can perform an action  $a$ , which results in another state  $s$ . If the agent performs  $a$ , there is a state transition from state  $r$  to state  $s$ .

In decision theory, making a decision results in an outcome. Outcomes represent all possible consequences of a decision. Outcomes can represent the state resulting from the action performed, effects in the far future, how pleasant the action was, and possibly the history of all preceding states. An outcome is thus a more general notion than a state, because outcomes can contain all information in states and even more.

**Values Versus Perspectives.** In [2], there is a valuation function  $\delta$  that takes a state transition and a value and returns whether that state transition either promotes, demotes, or is neutral towards that value. More specifically,  $\delta : S \times S \times V \rightarrow \{+, -, 0\}$  with  $S$  the set of states and  $V$  the set of values. Note that a state transition either promotes, demotes, or is neutral towards a value resembles Simon's simple valuation function, which values an outcome either as 'satisfactory', 'indifferent' or 'unsatisfactory'.

The valuation function must be specified for all state transitions and all values, which can become time consuming when the number of states or values increases. Namely, if there are  $n$  states and  $m$  values, then the valuation function must be specified for  $m \cdot n^2$  different inputs. Furthermore, if two agents disagree about whether a state transition promotes a value, e.g. whether performing an action promotes the value of fun, then they can only explain that that is the outcome of their valuation function. Since values typically are abstract, it is important to explain and discuss what a value means. This is not possible in the approach of [2].

In our approach, a value is represented with a perspective, which is associated with an ordering over outcomes. A perspective can be decomposed into other perspectives and a perspective can be associated with an attribute of outcomes. This allows agents to explain and argue why a transition or goal promotes one of their values. For example, an agent can explain that its value of 'fun' means maximising spending time with friends and minimising time at work. whereas another agent can then explain that to him fun means spending time in nature and accomplishing things at work.

Furthermore, decomposing an abstract perspective into more specific perspectives for which an ordering is more easy to specify, makes it less demanding to specify whether a transition promotes, demotes or is neutral towards a value.

If a perspective  $p$  represents a value, then its associated ordering  $\leq_p$  can be used to define the valuation function  $\delta$  for  $p$  in the following way

$$\delta(q_1, q_2, v) = \begin{cases} + & \text{if } q_1 <_v q_2 \\ - & \text{if } q_2 <_v q_1 \\ 0 & \text{if } q_1 \equiv_v q_2 \end{cases}$$

If  $\leq_v$  is a total order, i.e. no elements are incomparable, then  $\delta$  is a normal function, otherwise  $\delta$  is a partial function.

## 7 Conclusion

In this paper we have proposed several argument schemes to argue about what decision is best for an agent based on its preferences over outcomes. An agent's preferences are expressed in terms of values and goals and we propose a model to represent what a value means and how it affects an agent's preferences. If the meaning of a value is clear, goals can be justified or attacked by arguing that they promote or demote a value.

We represent values as perspectives over outcomes. By recursively decomposing the different aspects of a perspective into other perspectives until they are decomposed in attribute perspectives, the meaning of a perspective and thus a value is made explicit. In this way, an agent can explain what a value exactly means to him, which allows other agents to argue that some aspect is wrong or forgotten or that the wrong attribute is used. Agents can justify pursuing a goal using the perspectives that are important to an agent and the attributes that are associated to those perspectives. We have discussed a satisficing and a optimistic argument scheme to justify a goal. Furthermore, priorities between goals can be justified using the priorities agents have over perspectives.

In future work, the relation between values and goals may be explored further. Different agent types and different situations may lead to pursuing different goals. An optimistic goal may be undercut by stating that it is too hard to achieve, but when is a goal too hard to achieve? Moreover, how can an agent justify that an attribute value is satisfactory and how is that influenced by circumstances?

When an agent finds costs more important than the centrality of the neighbourhood, and a house in a suburb is \$1 cheaper than a house downtown, then the costs argument is stronger than the centrality argument. By extending the formalism in this paper with 'distances' between attribute values, such weird results might be solved.

## References

- [1] Amgoud, L., Prade, H.: Using arguments for making and explaining decisions. *Artificial Intelligence* 173(3-4), 413–436 (2009)
- [2] Atkinson, K., Bench-Capon, T., McBurney, P.: Computational representation of practical argument. *Synthese* 152(2), 157–206 (2006)
- [3] Bench-Capon, T.J.M.: Persuasion in practical argument using value-based argumentation frameworks. *Journal of Logic and Computation* 13(3), 429–448 (2003)
- [4] Bex, F., Prakken, H., Reed, C., Walton, D.: Towards a formal account of reasoning about evidence: Argumentation schemes and generalisations. *Artificial Intelligence and Law* 11(2), 125–165 (2003)
- [5] Caminada, M., Amgoud, L.: On the evaluation of argumentation formalisms. *Artificial Intelligence* 171(5-6), 286–310 (2007)

- [6] Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artificial Intelligence* 77(2), 321–358 (1995)
- [7] Kakas, A., Moraitis, P.: Argumentation based decision making for autonomous agents. In: *Proc. of 2nd Int. Conf. on Autonomous Agents and Multiagent Systems (AAMAS 2003)*, pp. 883–890 (2003)
- [8] Keeney, R., Raiffa, H.: *Decisions with Multiple Objectives*. Wiley, Chichester (1976)
- [9] Keeney, R.L.: *Value-Focused Thinking: A Path to Creative Decisionmaking*. Harvard University Press, Cambridge (1992)
- [10] Rokeach, M.: *The nature of human values*. Free Press, New York (1973)
- [11] Savage, L.J.: *The foundations of statistics*. Wiley, New York (1954)
- [12] Schwartz, S.H.: Universals in the content and structure of values: theoretical advances and empirical tests in 20 countries. *Advances in Experimental Social Psychology* 25, 1–65 (1992)
- [13] Simon, H.A.: A behavioral model of rational choice. *The Quarterly Journal of Economics*, 99–118 (1955)
- [14] Walton, D.N.: *Argumentation Schemes for Presumptive Reasoning*. Lawrence Erlbaum, Mahwah (1996)
- [15] van der Weide, T.L., Dignum, F., Meyer, J.-J.C., Prakken, G.A.W., Vreeswijk, H.: Practical reasoning using values. In: *McBurney, P., Rahwan, I., Parsons, S., Maudet, N. (eds.) ArgMAS 2009. LNCS, vol. 6057*, pp. 225–240. Springer, Heidelberg (2010)
- [16] Wellman, M.P., Doyle, J.: Preferential semantics for goals. In: *Proceedings of the National Conference on Artificial Intelligence*, pp. 698–703 (1991)

# On the Benefits of Argumentation-Derived Evidence in Learning Policies

Chukwuemeka David Emele<sup>1</sup>, Timothy J. Norman<sup>1</sup>,  
Frank Guerin<sup>1</sup>, and Simon Parsons<sup>2</sup>

<sup>1</sup> University of Aberdeen, Aberdeen, AB24 3UE, UK  
{c.emele,t.j.norman,f.guerin}@abdn.ac.uk

<sup>2</sup> Brooklyn College, City University of New York, 11210 NY, USA  
parsons@sci.brooklyn.cuny.edu

**Abstract.** An important and non-trivial factor for effectively developing and resourcing plans in a collaborative context is an understanding of the policy and resource availability constraints under which others operate. We present an efficient approach for identifying, learning and modeling the policies of others during collaborative problem solving activities. The mechanisms presented in this paper will enable agents to build more effective argumentation strategies by keeping track of who might have, and be willing to provide the resources required for the enactment of a plan. We argue that agents can improve their argumentation strategies by building accurate models of others' policies regarding resource use, information provision, etc. In a set of experiments, we demonstrate the utility of this novel combination of techniques through empirical evaluation, in which we demonstrate that more accurate models of others' policies (or norms) can be developed more rapidly using various forms of evidence from argumentation-based dialogue.

**Keywords:** Argumentation, Machine learning, Policies, Norms, Evidence.

## 1 Introduction

Distributed problem solving activities often require the formation of a team of collaborating agents. In such scenarios agents often operate under constraints placed on them by the organisations or interests that they represent. Such constraints, typically, determine the behaviour of representatives of organisations. When these constraints are part of the standard operating procedures of the agents or the organisations in question, we refer to them as *policies* (also known as norms). Members of the team agree to collaborate and perform joint activities in a mutually acceptable fashion. Often, agents in the team represent different organisations, and so there are different organisational constraints imposed on them. Even within a single organisation, team members often represent sub-organisations with different procedures and constraints. For example, the sales department of an organisation may possess certain operating procedures, which differ from those of the purchasing department. Furthermore, team members



Example 1:	Example 2:
<i>i</i> : Can I have a <i>screw-driver</i> ?	<i>i</i> : Can I have a <i>screw-driver</i> ?
<i>j</i> : What do you want to use it for?	<i>j</i> : What do you want to use it for?
<i>i</i> : To hang a picture.	<i>i</i> : To hang a picture.
<i>j</i> : No.	<i>j</i> : I can provide you with a <i>hammer</i> instead.
	<i>i</i> : I accept a <i>hammer</i> .

**Fig. 1.** Dialogue between Collaborating Agents (Examples 1 and 2)

may possess individual interests and goals that they seek to satisfy, which are not necessarily shared with other members of the team. These individual motivations largely determine the way in which members carry-out tasks assigned to them during joint activities.

In this paper, we focus on policy and resource availability constraints, and define policies as explicit prohibitions that members of the team are required to adhere to. Policy constraints may be team-wide or individual. We focus on individual policies. These policies are often private to that individual member or subset of the team, and are not necessarily shared with other members of the team. In order to develop effective plans, an understanding of the policy and resource availability constraints of other members in the team is beneficial. However, tracking and reasoning about such information is non-trivial.

Our conjecture is that machine learning techniques may be employed to aid decision making in this regard. Although this is not a new claim [7], it is novel to combine it with evidence derived from argumentation-based dialogue, which we call argumentation-derived evidence (ADE). We present a system where agents learn from dialogue by automatically extracting useful information (*evidence*) from the dialogue and using these to model the policies of others in order to adapt their behaviour in the future. We describe an experimental framework and present results of our evaluation in a resource provisioning scenario [4], which show empirically: (1) that evidence derived from argumentation-based dialogue can indeed be effectively exploited to learn better (more complete and correct) models of the policy constraints that other agents operate within; and (2) that through the use of appropriate machine learning techniques more accurate and stable models of others' policies can be derived more rapidly than with simple memorisation of past experiences.

To illustrate the sorts of interaction between agents in this framework, consider the following examples. Let *i* and *j* be two agents collaborating to hang a picture [11].

Following from the interaction in example 1 (see Figure 1), there is very little that we can learn from that encounter. It is unclear why agent *j* said no to agent *i*'s request. It could be that there exists some policy that forbids agent *j* from providing the *screw-driver* to agent *i*, or it could be that the *screw-driver* is not available at the moment. On the other hand, suppose we have an argumentation framework that allows agents to suggest alternatives as in

Example 3:	Example 4:
<i>i</i> : Can I have a <i>screw-driver</i> ?	<i>i</i> : Can I have a <i>screw-driver</i> ?
<i>j</i> : What do you want to use it for?	<i>j</i> : What do you want to use it for?
<i>i</i> : To hang a picture.	<i>i</i> : To hang a picture.
<i>j</i> : No.	<i>j</i> : No.
<i>i</i> : Why?	<i>i</i> : Why?
<i>j</i> : I'm not permitted to release <i>screw-driver</i> .	<i>j</i> : <i>Screw-driver</i> is not available.

**Fig. 2.** Dialogue between Collaborating Agents (Examples 3 and 4)

example 2 (see Figure 1) or ask for and receive explanations as in examples 3 and 4 (see Figure 2), then agent *i* can, potentially, gather more evidence regarding the provision of the resources involved.

Considering examples 3 and 4 (see Figure 2), it is worth noting that without the additional evidence, obtained by the information-seeking dialogue, the two cases are indistinguishable. This means that the agent will effectively be guessing which class these cases fall into. The additional evidence allows the agent to learn the right classification for each of the cases. It should be noted here that although in example 3, we now have an explanation that the resource is not to be provided for policy reasons, the question remains: what are the important characteristics of the prevailing circumstances that characterise this policy? Additional evidence is beneficial, and could be used to identify the important features that characterise an agent's prevailing circumstance. Each piece of evidence can be used to improve the model, hence, the quality of decisions made in future episodes. Section 3 discusses how we do this.

In a domain where there are underlying constraints that could yield similar results, standard machine learning techniques will have limited efficacy. Using argumentation to gather additional evidence could improve the accuracy of the information learned about the policies of others. We claim that significant improvements can be achieved because argumentation can help clarify reasons behind decisions made by others (e.g. the *provider*).

In the research presented in this paper, we intend to validate the following hypotheses: (1) Allowing agents to exchange arguments during practical dialogue (like negotiation) will mean that the proportion of correct policies learned during interaction will increase faster than when there is no exchange of arguments; and (2) Through the use of appropriate machine learning techniques more accurate and stable models of others' policies can be derived more rapidly than with simple memorisation of past experiences.

The remainder of this paper is organised as follows: In section 2 we briefly describe argumentation-based dialogue and introduce the protocol employed. Learning policies is discussed in section 3 and section 4 describes our simulation environment. Experimental results are reported in section 5. Section 6 discusses related work and future direction, and the paper concludes in section 7.

## 2 Argumentation-Based Dialogue

In this section we present the argumentation-based negotiation protocol which will be used in guiding the negotiation process, and for obtaining additional evidence from the interaction. This protocol uses information-seeking dialogue [13,18] to probe for additional evidence.

### 2.1 The Negotiation Protocol

The negotiation for resources takes place in a turn-taking fashion, where the *seeker* agent sends a request for resource to a *provider* agent. Figure 3 captures the negotiation protocol in a AUML-like interaction diagram ([www.fipa.org](http://www.fipa.org)). If the *provider* agent has the requested resource in its resource pool and it is in a usable state then it checks whether there is any policy constraint that forbids it from providing the resource to the *seeker* or not. If the *provider* agent needs more information from the *seeker* in order to make a decision, the *provider* agent would ask for more information to be provided. This is the information gathering stage. The information gathering cycle will continue until the *provider* has acquired enough information (necessary to make the decision), or the *seeker* refuses to provide more information and the negotiation ends.

The *provider* agent releases the resource to the *seeker* agent if there is no policy that prohibits the *provider* agent from doing so. Otherwise, the *provider* agent offers an alternative resource (if there are no policies that forbid that line of action and the alternative resource is available). When an alternative resource is suggested by the *provider* agent, the *seeker* agent evaluates it. If it

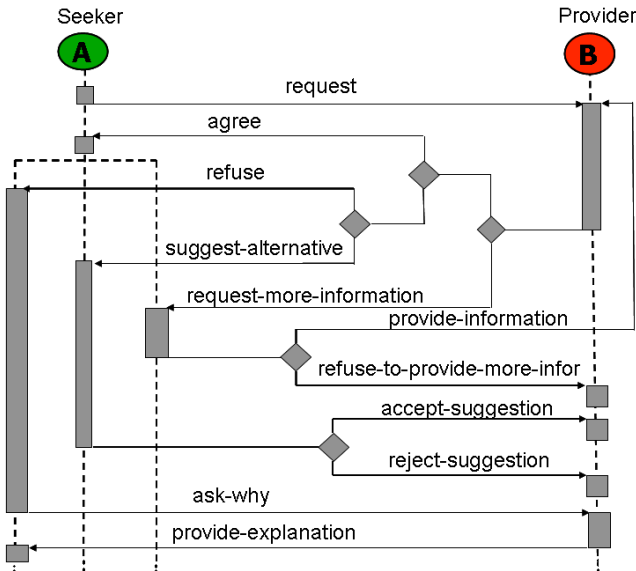


Fig. 3. The negotiation protocol

is acceptable, the *seeker* agent accepts it and the negotiation ends. Otherwise, the *seeker* agent refuses the alternative (in principle, this cycle may be repeated until an alternative is accepted or the negotiation ends). However, for simplicity and brevity, only one suggest-refuse cycle is permitted per request.

From a learning point of view, the suggestion of alternative resources could be a positive evidence that the *provider* agent does not have any policy that forbids the provision of the alternative resource to the *seeker*. In addition, it provides an evidence that the alternative resource is also available. This extra evidence, we anticipate, may help to improve the performance of the learner in predicting the policy constraints of *provider* agents in future encounters.

If there is a policy constraint that forbids the provision of the resource, or the resource is not available then the *provider* agent will refuse to provide the resource to the *seeker* agent. From the *seeker*'s perspective, the refusal could be as a result of policy constraint or because the resource is not available. In order to disambiguate which of these constraints are responsible for the refusal, the *seeker* agent switches to argumentation-based dialogue. The *seeker* agent asks for explanations for the refusal so as to gather further evidence and thereby identify the underlying constraints. The *provider* agent, therefore, responds with some explanations and the negotiation ends. Three categories of explanations are investigated in this framework: (1) Policy constraints (whereby the agent attributes the refusal to policy constraints); (2) Resource not available (here, the agent attributes the refusal to resource availability constraints); (3) Won't tell you (in this case, the agent refuses to give any explanation for the refusal). These pieces of evidence will be explored further in Section 2.3.

## 2.2 Policies

In our model, policies are conditional entities (or rules) and so are relevant to an agent under specific circumstances only. These circumstances are characterised by a set of features. Some examples of features may include: (1) the height of a tower; (2) the temperature of a room; (3) the manufacturer of a car, etc. In other words, policies map feature vectors,  $\bar{F}$ , of agents to appropriate policy decisions.

In order to illustrate the way policies may be captured in this model, we present the following examples (see Figure 4). Let  $F$  be the set of all features that characterise an agent's prevailing circumstance such that  $f_1, f_2, \dots \in F$ . Assuming,  $f_1$  is resource,  $f_2$  is purpose,  $f_3$  is location,  $f_4$  is the affiliation of the

$\mathbb{P}_1$ : You are <b>permitted</b> to release a <i>ladder</i> to an agent if the <i>ladder</i> is required for the purpose of <i>hanging a picture</i> .
$\mathbb{P}_2$ : You are <b>prohibited</b> from releasing a <i>screw-driver</i> to an agent if the <i>screw-driver</i> is to be used for <i>hanging a picture</i> .
$\mathbb{P}_3$ : You are <b>permitted</b> to release a <i>hammer</i> to an agent.
$\mathbb{P}_4$ : You are <b>permitted</b> to release a <i>nail</i> to an agent.
$\mathbb{P}_5$ : You are <b>permitted</b> to release a <i>table</i> to an agent.

Fig. 4. Sample Agent Policies

Example A
<i>i</i> : <b>request</b> ( <i>i</i> , <i>j</i> , <i>screw-driver</i> )
<i>j</i> : <b>ask-infor</b> ( <i>j</i> , <i>i</i> , need( <i>screw-driver</i> , <i>f</i> <sub>2</sub> , <i>f</i> <sub>3</sub> , <i>f</i> <sub>5</sub> ))
<i>i</i> : <b>provide-infor</b> ( <i>i</i> , <i>j</i> , need( <i>screw-driver</i> , <i>f</i> <sub>2</sub> = <i>x</i> , <i>f</i> <sub>3</sub> = <i>y</i> , <i>f</i> <sub>5</sub> = <i>z</i> ))
<i>j</i> : <b>refuse</b> ( <i>j</i> , <i>i</i> , <i>screw-driver</i> )
<i>i</i> : <b>why</b> ( <i>i</i> , <i>j</i> , refuse( <i>screw-driver</i> ))
<i>j</i> : <b>inform</b> ( <i>j</i> , <i>i</i> , <i>screw-driver</i> , reason(policy-constraints))
Example B
<i>i</i> : <b>request</b> ( <i>i</i> , <i>j</i> , <i>nail</i> )
<i>j</i> : <b>refuse</b> ( <i>j</i> , <i>i</i> , <i>nail</i> )
<i>i</i> : <b>why</b> ( <i>i</i> , <i>j</i> , refuse( <i>nail</i> ))
<i>j</i> : <b>inform</b> ( <i>j</i> , <i>i</i> , <i>nail</i> , reason(wont-tell-you))
<i>i</i> : <b>request</b> ( <i>i</i> , <i>j</i> , <i>table</i> )
<i>j</i> : <b>agree</b> ( <i>j</i> , <i>i</i> , <i>table</i> )

**Fig. 5.** Dialogue between agents *i* and *j*

agent, and *f*<sub>5</sub> is the day the resource is required then an agent’s policies may be captured as shown in Figure 4.

### 2.3 Argumentation-Derived Evidence

Following the argumentation-based negotiation protocol described earlier, the agents could ask for more information (with respect to a request or the response to a request), which indicates what constraints others may be operating within. For instance, let us assume that a *provider* agent has a policy that forbids it from providing a *screw-driver* to any *seeker* agent that intends to use it for *hanging a picture*. Then, whenever a *screw-driver* is requested the *provider* agent will probe for more information to ascertain that the purpose the *seeker* intends to use the *screw-driver* for is not *hanging a picture*. This extra evidence could be useful. Similarly, whenever a *seeker* agent’s request is refused then the *seeker* agent will ask for explanations/justifications for the refusal. These additional evidence are beneficial, and we expect them to improve the quality of the models of other agents that can be inferred in future encounters.

Figure 5 shows two simple examples of the kind of dialogue that may occur between two collaborating agents, *i* and *j*. For the purpose of the example, we use **need**(*R*, *f*<sub>2</sub>, *f*<sub>3</sub>, *f*<sub>5</sub>) to denote that the *seeker* agent intends to use the resource *R* for a purpose (captured by feature *f*<sub>2</sub>) at a location (represented by *f*<sub>3</sub>) on a given day (captured by *f*<sub>5</sub>). Note that although this is presented as a dialogue between two agents, in reality the initiator (agent *i*, the agent that wishes to resource its plan) may engage in multiple instances of this dialogue with other agents.

## 3 Learning Policies

In this section we discuss the machine learning techniques that we have explored for learning policies through argumentation-derived evidence. These techniques

include decision tree learning (C4.5), instance-based learning ( $k$ -Nearest Neighbours, abbreviated as  $k$ -NN) and rule-based learning (Sequential Covering, abbreviated as SC).

Our approach seeks ways to combine argumentation analysis with already existing machine learning techniques with a view to improving the performance of agents at predicting the policy constraints of others. We anticipate that this could enable them to build more effective argumentation strategies. In other words, we argue that evidence derived from argumentation-based dialogue can indeed be effectively exploited to learn better (more complete and correct) models of the policy constraints that other agents operate within. Also, we claim that through the use of appropriate machine learning techniques more accurate and stable models of others' policies can be derived more rapidly than with simple memorisation of past experiences. In future encounters, the *seeker* agent attempts to predict the policies of *provider* agents based on the model it has built.

When an agent has a collection of experiences with other agents described by feature vectors, we can make use of existing machine learning techniques for learning associations between sets of discrete attributes (e.g.  $f_1, f_2, f_3, f_4, f_5$ ) and policy decisions. The following sections discuss specifically, the three classes of machine learning algorithms<sup>1</sup> [8], namely: decision tree learning (using C4.5), instance-based learning (using  $k$ -nearest neighbours), and rule-based learning (using sequential covering).

We define a *learning interval*,  $\phi$ , which determines the number of interactions<sup>2</sup> an agent must engage in before building (or re-building) its policy model<sup>3</sup>. Once an agent has had  $\phi$  interactions, the policy learning process proceeds as follows. For each interaction, which involves collaborating with provider  $j$  for the provision of a resource, we add the example  $(\bar{F}_j, grant)$  or  $(\bar{F}_j, deny)$  to the training set, depending on the evidence obtained from the interaction. The model is then constructed. In this way, an agent may build a model of the relationship between observable features of agents and the policies they are operating under. Subsequently, when faced with resourcing a new task, the policy model can be used to obtain a prediction of whether a particular provider has a policy that permits him to provide the resource or not. This satisfies our requirement for a policy learning mechanism.

### 3.1 Decision Tree Learning (C4.5)

C4.5 [14] builds decision trees from a set of training data, using the concept of information entropy [8] (beyond the scope of this paper). Generally, the training data is a set  $S = s_1, s_2, \dots, s_n$  of already classified samples. Each sample  $s_i = x_1, x_2, \dots, x_m$  is a vector where  $x_1, x_2, \dots, x_m$  represent attributes of the

<sup>1</sup> We use the Weka [19] implementation of these algorithms. Weka is a popular open-source machine learning toolkit written in Java.

<sup>2</sup> By interaction, we mean, an entire plan resourcing episode.

<sup>3</sup> This requirement relates, specifically, to C4.5 and sequential covering because they are non-incremental learning algorithms. A detailed discussion of non-incremental learning is covered in [19].

---

**Algorithm 1.** The C4.5 algorithm

---

- 1: **Check** for base cases
  - 2: **For** each attribute  $D$ ,  
    Find the normalised information gain from splitting on  $D$
  - 3: **Let**  $D_{best}$  be the attribute with the highest normalised information gain
  - 4: **Create** a decision node that splits on  $D_{best}$
  - 5: **Recurse** on the sublists obtained by splitting on  $D_{best}$ , and add those nodes as children of the node
- 

**Fig. 6.** The C4.5 algorithm

sample. The training data is augmented with a vector  $C = c_1, c_2, \dots, c_n$  where  $c_1, c_2, \dots, c_n$  represent the class to which each sample belongs.

Integrating this algorithm into our system with the intention of learning policies is appropriate since the algorithm supports concept learning and policies can be conceived as concepts/features of an agent. Agent policies are represented as a vector of attributes (e.g. resource, purpose, location, etc.) and these attributes are communicated back and forth during negotiation. The C4.5 algorithm is then used to classify each set of attributes (policy instance) into a class. There are two classes: grant and deny. Grant means that the *provider* agent will possibly provide the resource that is requested while deny implies that the *provider* agent will potentially refuse. The leaf nodes of a decision tree hold the class labels of the instances while the non-leaf nodes hold the test attributes. In order to classify a test instance, the C4.5 algorithm searches from the root node by examining the value of test attributes until a leaf node is reached and the label of that node becomes the class of the test instance. Figure 6 outlines the C4.5 algorithm in pseudo-code.

The C4.5 algorithm has three base cases.

- All the samples in the list belong to the same class. When this happens, it simply creates a leaf node for the decision tree saying to choose that class.
- None of the features provide any information gain. In this case, C4.5 creates a decision node higher up the tree using the expected value of the class.
- Instance of previously-unseen class encountered. Again, C4.5 creates a decision node higher up the tree using the expected value.

The problem with this algorithm is that it is not incremental, which means all the training examples should exist before learning. To overcome this problem, the system keeps track of the *provider* agent's responses. After a number of interactions, predefined by  $\phi$ , the decision tree is rebuilt. Without doubt, there is a computational drawback involved in periodically reconstructing the decision tree. However, in practice, we have evaluated C4.5 to be fast and the reconstruction cost to be small. Our approach is similar to the incremental induction of decision trees proposed in [17].

### 3.2 Instance-Based Learning ( $k$ -NN)

The  $k$ -nearest neighbours algorithm ( $k$ -NN) [3] is a type of instance-based learning, or lazy learning, where the function is only approximated locally and all

computation is deferred until classification. The universal set of all the policies an agent may be operating within could be conceived as a feature space (or a grid) and the various policy instances represent points on the grid. Using  $k$ -NN, a policy instance is classified by a majority vote of its neighbours, with the policy instance being assigned to the class most common amongst its  $k$  nearest neighbours, where  $k$  is a positive integer, typically small. The  $k$ -NN algorithm is incremental, which means all the training examples need not exist at the beginning of the learning process. This is a good feature because the policy model could be updated as new knowledge is learned.

The  $k$ -nearest neighbour algorithm is sensitive to the local structure of the data and this, interestingly, makes  $k$ -NN a good candidate for learning policies because slight changes in the variables/attributes of a policy could trigger different action. For example:

**Policy1:** You are *permitted* to release resource  $R$  to team member  $X$  if his affiliation is  $O$  and  $R$  is to be deployed at location  $L$  for purpose  $P$  on day 1.

**Policy2:** You are *prohibited* from releasing resource  $R$  to team member  $X$  if his affiliation is  $O$  and  $R$  is to be deployed at location  $L$  for purpose  $P$  on day 2.

In order to identify neighbours, the policy instances are represented by position vectors in a multidimensional feature space. In this approach, new policy instances are classified based on the closest training examples in the feature space. A policy instance is assigned to the class  $c$  if it is the most frequent class label among the  $k$  nearest training samples. It is usual to use the Euclidean distance, though other distance measures, such as the Manhattan distance, Hamming distance could in principle be used instead. The training phase of the algorithm consists only of storing the feature vectors and class labels of the training samples. In the actual classification phase, the test sample is represented as a vector in the feature space. Distances from the new vector to all stored vectors are computed and  $k$  closest samples are selected.

A major drawback to using this technique to classify a new vector to a class is that the classes with the more frequent examples tend to dominate the prediction of the new vector, as they tend to come up in the  $k$  nearest neighbours when the neighbours are computed due to their large number. The distance-weighted  $k$ -NN algorithm, which weights the contribution of each of the  $k$  neighbours according to their distance to the new vector, uses distance weights to minimise the bias caused by the imbalance in the training examples by giving greater weight to closer neighbours. In our work, the weight of a neighbour is computed as the inverse of its distance from the new vector.

### 3.3 Rule-Based Learning (Sequential Covering)

Since policies guide the way entities within a community (or domain) act by providing rules for their behaviour it makes sense to learn policies as rules. Sequential covering algorithm [28] is a rule-based learning technique, which constructs rules by sequentially covering the examples. The sequential covering algorithm,



---

**Algorithm 2.** The Sequential Covering Algorithm

---

- 1: **Input** the training data ( $D$ ) and the classes ( $C$ )
- 2: **For** each class  $c \in C$
- 3:   **Initialise**  $E$  to the instance set
- 4:   **Repeat**
- 5:     **Create** a rule  $R$  with an empty left-hand side (LHS) that predicts class  $c$ :
- 6:     **Repeat**
- 7:       **For** each ( $Attribute, Value$ ) pair found in  $E$
- 8:         **Consider** adding the condition  $Attribute = Value$  to the LHS of  $R$
- 9:         **Find**  $Attribute = Value$  that maximises  $p/t$
- 10:         (break ties by choosing the condition with the largest  $p$ )
- 11:         **Add**  $Attribute = Value$  to  $R$
- 12:         **Until**  $R$  is perfect (or no more attributes to use)
- 13:         **Remove** the instances covered by  $R$  from  $E$
- 14:     **Until**  $E$  contains no more instances that belong to  $c$

---

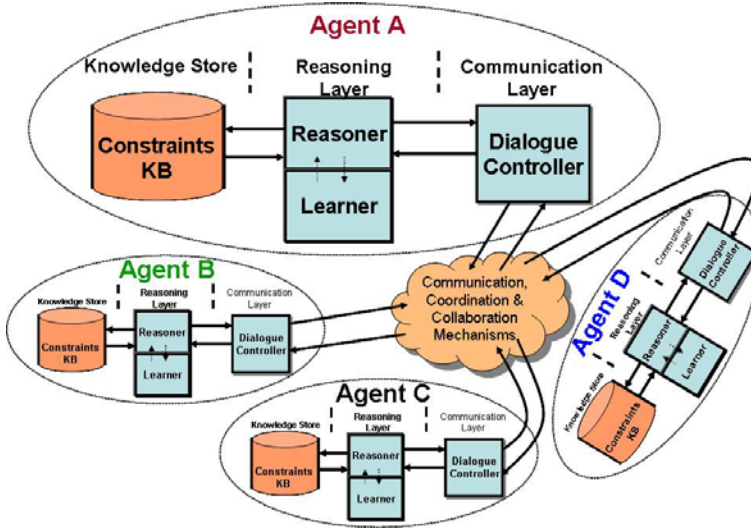
**Fig. 7.** The Sequential Covering Algorithm

SC for short, is a method that induces one rule at a time (by selecting attribute-value pairs that satisfy the rule), removes the data covered by the rule and then iterates the process. SC generates rules for each class by looking at the training data and adding rules that completely describe all tuples in that class. For each class value, rule antecedents are initially empty sets, augmented gradually for covering as many examples as possible. Figure 7 outlines the sequential covering algorithm in pseudo-code.

In this study we used three different machine learning mechanisms: *Decision tree learning*, *Instance-based learning* and *Rule-based learning*. These three mechanisms represent very different classes of machine learning algorithms. The rationale for exploring a range of learning techniques is to demonstrate the utility of argumentation-derived evidence regardless of the machine learning technique employed. Thus, we hypothesize that the use of evidence acquired through argumentation significantly improves the performance of machine learning in the development and refinement of models of other agents. Also, we claim that through the use of appropriate machine learning techniques more accurate and stable models of others' policies can be derived more rapidly than with simple memorisation of past experiences.

## 4 Simulation Environment

To test our hypotheses, we developed a simulation environment that combines mechanisms for agents to engage in argumentative dialogue and to learn from dialogical encounters with other agents. For the purpose of resourcing plans, agents may act as resource seekers, which collaborate and communicate with



**Fig. 8.** Architecture of the framework for learning policies in team-based activities using dialogue

potential providers to perform joint actions. The enactment of both *seeker* and *provider* roles are governed by individual policies that regulate their actions. A *seeker* agent requires resources in order to carry out some assigned tasks. The *seeker* agent generates requests in accordance with its policies and negotiates with the *provider* agents based on these constraints. On the other hand, *provider* agents have access to certain resources and may have policies that govern the provision of such resources to other members of the team.

Although agents may have prior assumptions about the policies that constrain the activities of others, these models are often incomplete and may be inaccurate. *Provider* agents do not have an unlimited pool of resources and so some resources may be temporarily unavailable. By a resource being available we mean that it is not committed to another task (or agent) at the time requested and the resource is in a usable state. Both *seeker* and *provider* agents have access to the team-wide policies but not the individual policies of others. Agents in this domain play the role of a *seeker* or a *provider* in different interactions.

#### 4.1 Architecture

Figure 8 depicts our architecture. Each agent has two main layers, the communication layer and the reasoning layer. The communication layer embodies the dialogue controller, which handles all communication with other agents in the domain. The dialogue controller sends/receives messages to/from other agents, and the reasoning layer reasons over the dialogue. If an agent is playing the role of a *seeker* agent then the dialogue controller sends out the request for resources. On the other hand, if the agent is a *provider* agent then the dialogue controller receives a request and passes it on to the reasoning layer.

The reasoning layer consists of two modules: the reasoner and the learner. Upon receiving a message (e.g. a request), the reasoner evaluates the message and determines the response of the agent. In most cases, the reasoner looks up policy constraints from the knowledge-base and generates the appropriate response for the agent. Policy and non-policy constraints are stored in the constraints knowledge-base. Whenever the agent observes a new pattern of behaviour the agent uses this experience as evidence for learning, and updates the model of the other agent accordingly. The learner uses standard machine learning techniques to learn policies based on the perceived actions of other agents. The learning techniques are discussed in Section 3.

The knowledge store in Figure 8 acts as a repository where an agent stores the constraints it has learned by interacting with other agents in the domain. The information includes the features that an agent requires in order to make a decision about providing a resource or not. For example, following from [11], a *provider* agent  $j$  may need to know what the purpose for requesting a *screw-driver* is before deciding whether to release the *screw-driver* or not. The *seeker* agent stores such information about agent  $j$  in the knowledge store. Also, the decision of  $j$  after the purpose has been revealed will also be learned for future interactions.

To achieve this, we have developed a simple dialogue game<sup>4</sup> involving *seeker* agents and *provider* agents operating under different constraints. The players take turns and the game starts with an agent,  $i$ , sending a request to another agent,  $j$ , for the use of some resources needed to fulfill a plan. The other agent ( $j$ ) responds with an agree or refuse based on the prevailing context, e.g. policy constraints. The requesting agent could ask for explanations and reasons for an action, and so on until the game ends.

## 4.2 Implementation

We implemented a simulation environment for agent support in team-based problem solving and integrated our learning and argumentation mechanisms into the framework. The policies are encoded as rules in a rule engine [6]. The application programming interface in Weka [19] was used to integrate standard machine learning algorithms into the framework. We note that, although these three learning algorithms were used, the framework is configured such that other machine learning algorithms can be plugged in. As discussed in the previous section, we evaluated the performance of a decision tree learner (C4.5), an Instance based learner ( $k$ -Nearest Neighbour algorithm) and a rule based learner (Sequential Covering) in learning policies through argumentation-derived evidence.

The simulation environment allows us to generate multiple *providers* with randomised policies, *seeker* agents with randomised initial models of the policies of *providers* in the simulation and randomised problems for the *seeker* to solve (that is, random resource requirements). The *seeker* predicts (based on the model of the *provider*) whether the *provider* has a policy that forbids/permits the

---

<sup>4</sup> Dialogue games have proven extremely useful for modeling various forms of reasoning in many domains [1].

---

**Assume** *seeker* *i* requests resource *R* from *provider* *j*

---

```

IF    ( is_available(R) ∧ NOT (forbid(release(R, i)) )
THEN
      agree( release(R, i))
ELSE
      refuse( release(R, i))

```

---

**Fig. 9.** *Provider* agents' pseudo decision function

provision of such resource in that context. The *seeker* requests the required resource from the *provider* agent and the *provider* uses a simple decision function (See Figure 9) to decide whether to grant or deny the request.

If the decision of the *provider* agent deviates from the predictions of the *seeker* agent then the *seeker* agent seeks additional evidence (through dialogue) to disambiguate whether the deviation was as a result of policy or resource availability constraints. The dialogue follows the protocol specified in Figure 3, and at the end of the interaction the outcome is learned by the *seeker* and the model of the *provider* is updated accordingly. This adaptive learning process serves to improve the quality of the models of the other agents that can be inferred from their observable actions in future interactions.

## 5 Experiments and Results

In a series of experiments, we show how learning techniques and argumentation can support agents engaging in collaborative activities, increase their predictive accuracy, avoid unnecessary policy conflicts, hence improve their performance. The experiments show that agents can effectively and rapidly increase their predictive accuracy of the learned model through the use of dialogue.

The scenario adopted in this research involves a team of five software agents (one *seeker* and four *provider* agents) collaborating to complete a joint activity in a region over a period of three days. The region is divided into five locations. There are five resource types, and five purposes that a resource could be used to fulfill. A task involves the *seeker* agent identifying resource needs for a plan and collaborating with the *provider* agents to see how that plan can be resourced.

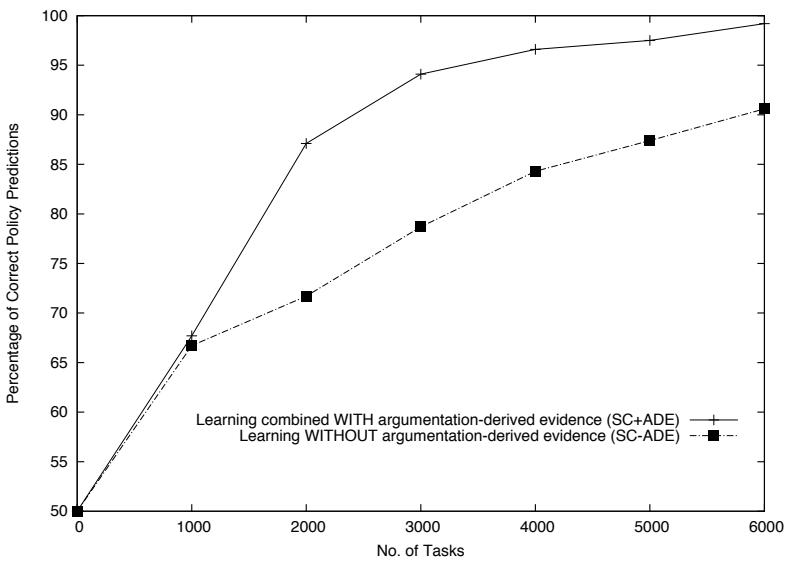
Argumentation-derived evidence (ADE) was incorporated into the learning process of the three machine learning techniques (C4.5, *k*-NN, and SC) described earlier, and their performances in learning the policy constraints of others were evaluated. A simple lookup table (hereafter called, LT) was used as a control condition and it serves as a structure for simple memorisation of outcomes from past encounters.

### 5.1 Results

This section presents the results of the experiments carried out to validate this work. Experiments were conducted with *seeker* agents initialised with random models of the policies of *provider* agents. 100 runs were conducted for each

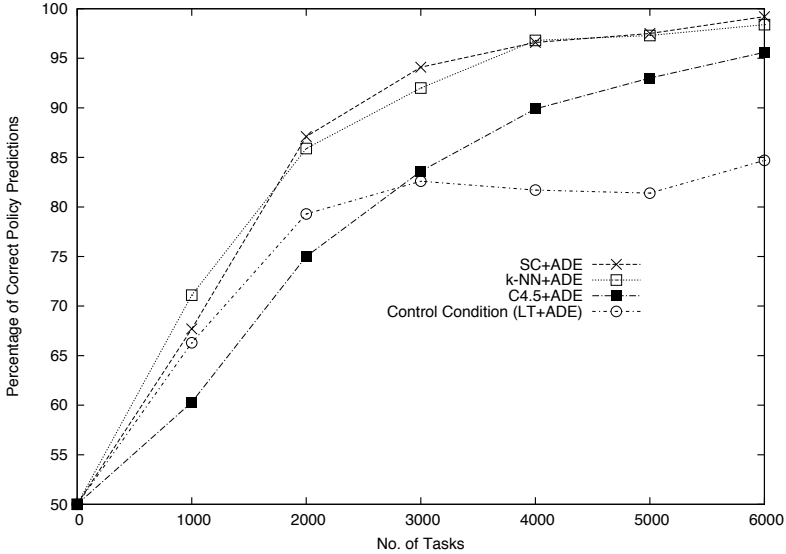
**Table 1.** Average percentage of policies classified correctly and standard deviation

Approach \ Tasks	1000	2000	3000	4000	5000	6000
LT-ADE	65.1±6.5	70.3±10.3	75.6±6.7	78.1±10.2	79.3±8.3	81.3±10.1
LT+ADE	66.3±6.0	79.3±9.3	83.6±8.2	81.7±11.2	81.4±7.8	84.7±9.1
C4.5-ADE	58.3±15.1	69.2±16.6	75.1±12.0	82.1±12.3	85.3±8.9	88.2±8.2
C4.5+ADE	60.3±14.4	75.0±12.6	83.6±6.5	89.9±5.2	93.0±3.4	95.6±5.1
$k$ -NN-ADE	65.2±9.8	71.0±7.8	75.3±5.3	80.7±3.8	81.0±4.1	82.0±3.8
$k$ -NN+ADE	71.1±9.0	85.9±7.3	92.0±4.6	96.8±3.1	97.3±3.6	98.4±1.7
SC-ADE	66.7±8.2	71.7±6.0	78.7±8.4	84.3±6.5	87.4±6.0	90.6±5.3
SC+ADE	67.7±7.7	87.1±6.4	94.1±4.2	96.6±4.1	97.5±2.6	99.2±1.0

**Fig. 10.** Graph showing the effectiveness of allowing the exchange of arguments in learning policies

case, and tasks were randomly created during each run from 375 possible configurations.

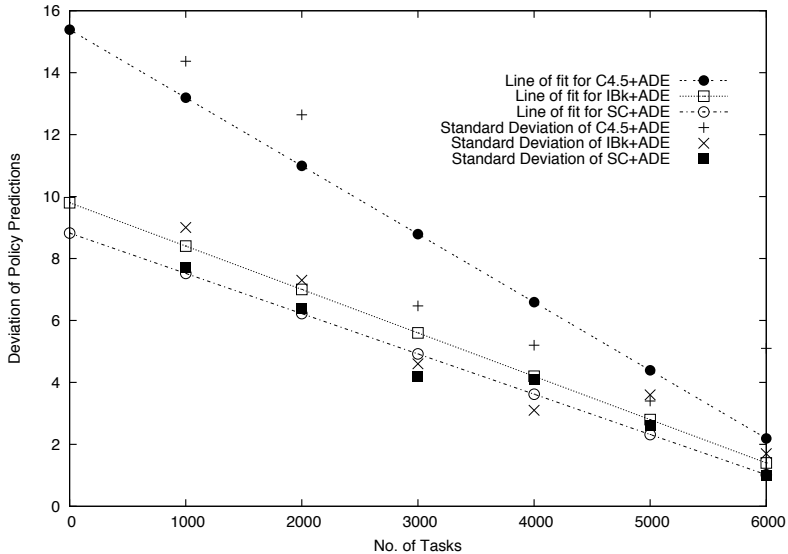
Table 1 illustrates the effectiveness of identifying and learning policies through argumentation-derived evidence using the three machine learning techniques described earlier, and the control condition (lookup table). It shows the average percentage of policies classified correctly and the standard deviations for each of the approaches, namely: Lookup Table without the aid of argumentation-derived evidence (LT-ADE), Lookup Table enhanced with argumentation-derived evidence (LT+ADE), C4.5-ADE, C4.5+ADE,  $k$ -NN-ADE,  $k$ -NN+ADE, SC-ADE, and SC+ADE. In each case, the model of others' policies is recomputed after each set of 1000 tasks. For all three machine learning techniques considered, the



**Fig. 11.** Graph showing the effectiveness of learning policies with the aid of argumentation-derived evidence using various techniques (LT+ADE, C4.5+ADE,  $k$ -NN+ADE & SC+ADE)

percentage of policies predicted correctly as a result of exploiting evidence derived from argumentation was consistently and significantly higher than those predicted without such evidence. Figure 10 gives a graphical illustration of the effectiveness of learning policies with the aid of argumentation-derived evidence using rule-based learning technique, for instance. After 3000 tasks, the accuracy of the approach with additional evidence had risen above 94% while the configuration without additional evidence was approaching 79%. It is easy to see that the experiments where additional evidence was combined with machine learning significantly and consistently outperformed those without additional evidence. These results show that the exchange of arguments during practical dialogue enabled agents to learn and build more accurate models of others' policies much faster than scenarios where there was no exchange of arguments.

Figure 11 captures the effectiveness of the three machine learning techniques described earlier, and a simple memorisation technique (a lookup table) in learning policies. The result shows that both instance-based learning ( $k$ -NN+ADE) and rule-based learning (SC+ADE) constantly and consistently outperform the control condition (LT+ADE) throughout the experiment. It is interesting to see that, with relatively small training set, the control condition performed better than the decision tree learner (C4.5+ADE). This is, we believe, because the model built by the decision tree learner overfit the data. The tree was pruned after each set of 1000 tasks and after 3000 tasks the accuracy of the C4.5+ADE model rose to about 83% to tie with the control condition and from then the decision tree learner performed better than the control condition.



**Fig. 12.** Graph showing the rate of convergence of the three techniques enhanced with ADE in learning policies (C4.5+ADE,  $k$ -NN+ADE, & SC+ADE)

The performance of the control condition dropped to about 81% after 4000 tasks. After 6000 tasks the accuracy of the decision tree learner had risen above 95% while that of the control condition was just over 84%.

Tests of statistical significance were applied to the results. The standard deviations of the results were analysed and the trend line plotted. (See Figure 12). Using linear regression, the analysis of variance (ANOVA) shows that as the number of tasks increases, each of the three machine learning techniques (with or without argumentation-derived evidence) consistently converges with a 95% confidence interval. Furthermore, for all the pairwise comparisons, the scenarios where argumentation-derived evidence was combined with machine learning techniques consistently yielded higher rates of convergence ( $p < 0.02$ ) than those without additional evidence. Specifically, the decision tree learner enhanced with argumentation-derived evidence (C4.5+ADE) converges ( $y = 15.3944 - 0.0022x$ ) with a  $F$  value of 15.66 and significance  $p = 0.0167$ . The  $k$ -NN+ADE converges ( $y = 9.7983 - 0.0014x$ ) with a  $F$  value of 38.58 and significance  $p = 0.0034$ , and the SC+ADE ( $y = 8.819 - 0.0013x$ ) converges with a  $F$  value of 136.45 and significance  $p = 0.0003$ . On the other hand, with a significance  $p = 0.3957$ , there is no statistical significance as to whether LT+ADE converges or not. These results confirm our hypotheses.

## 6 Discussion and Related Work

The research presented in this paper represents the first model for using evidence derived from argumentation to learn underlying social characteristics

(e.g. policies/norms) of others. There is, however, some prior research in combining machine learning and argumentation, and in using argument structures for machine learning. In that research, Možina et al. [9] propose a novel induction-based machine learning mechanism using argumentation. The work implemented an argument-based extension of CN2 rule learning (ABCN2) and showed that ABCN2 out-performed CN2 in most tasks. However, the framework developed in that research will struggle to disambiguate between constraints that may produce similar outcome/effect, which is the main issue we are addressing in our work. Also, the authors assume that the agent knows and has access to the arguments required to improve the prediction accuracy, but we argue that it is not always the case. As a result, we employ information-seeking dialogue to tease out evidence that could be used to improve performance.

In related research, Rovatsos et al. [15] use hierarchical reinforcement learning in modifying symbolic constructs (*interaction frames*) that regulate agent *conversation patterns*, and argue that their approach could improve an agent’s conversation strategy. In our work, we used information-seeking dialogue to obtain evidence from the interaction and learned the entire sequence as against a segment (frame) of the interaction [15]. We have demonstrated the effectiveness of using argumentation-derived evidence to learn underlying social characteristics (e.g. policies) without assuming that those underlying features are public knowledge.

In recent research, Sycara et al. [16] investigate agent support for human teams in which software agents aid the decision making of team members during collaborative planning. One area of support that was identified as important in this context is guidance in making policy-compliant decisions. This prior research focuses on giving guidance to humans regarding their own policies. An important and open question, however, is how can agents support human decision makers in developing models of others’ policies and using these in guiding the decision maker? Our work is aimed at bridging this gap (a preliminary version was presented in [5]). We employ a novel combination of techniques in identifying, learning and building accurate models of others’ policies, with a view to exploiting these in supporting human decision making.

In our future work, we plan to develop strategies for advising human decision makers on how a plan may be resourced and who to talk to on the basis of policy and resource availability constraints learned [10]. Parsons et al. [12] investigated the properties of argumentation-based dialogues and examined how different classes of protocols can have different outcomes. Furthermore, we plan to explore ideas from this work to see which class of protocol will yield the “best” result in this kind of task. We are hoping that some of these ideas will drive the work on developing strategies for choosing who to talk to.

## 7 Conclusions

In this paper, we have presented a technique that combines machine learning and argumentation for learning policies in a team of collaborating agents engaging in joint activities. We believe, to the best of our knowledge, that this is



the first study into learning models of other agents using argumentation-derived evidence. The results of our empirical investigations show that evidence derived from argumentation can have a statistically significant positive impact on identifying, learning and modeling others' policies during collaborative activities. The results also demonstrate that through the use of appropriate machine learning techniques more accurate and stable models of others' policies can be derived more rapidly than with simple memorisation of past experiences. Accurate policy models can inform strategies for advising human decision makers on how a plan may be resourced and who to talk to [16], and may aid in the development of more effective strategies for agents [10]. Our results demonstrate that significant improvements can be achieved by combining machine learning techniques with argumentation-derived evidence. Having shown that accurate models of others' policies could be learned through argumentation-derived evidence, we conjecture that one could, in principle, learn accurate models of other agents' properties (e.g. priorities, preferences, and so on).

**Acknowledgements.** This research was sponsored by the U.S. Army Research Laboratory and the U.K. Ministry of Defence and was accomplished under Agreement Number W911NF-06-3-0001. The views and conclusions contained in this document are those of the author(s) and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Army Research Laboratory, the U.S. Government, the U.K. Ministry of Defence or the U.K. Government. The U.S. and U.K. Governments are authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation hereon.

## References

1. Bench-Capon, T.J.M., Freeman, J.B., Hohmann, H., Prakken, H.: Computational models, argumentation theories and legal practice. In: Reed, C., Norman, T.J. (eds.) *Argumentation Machines. New Frontiers in Argument and Computation*, pp. 85–120. Kluwer Academic Publishers, Dordrecht (2003)
2. Cendrowska, J.: Prism: An algorithm for inducing modular rules. *International Journal of Man-Machine Studies* 27(4), 349–370 (1987)
3. Cover, T., Hart, P.: Nearest neighbor pattern classification. *IEEE Transaction on Information Theory* 13(1), 21–27 (1967)
4. Emele, C.D., Norman, T.J., Guerin, F., Parsons, S.: Learning policy constraints through dialogue. In: *Proc. of the AAAI Fall Symposium on The Uses of Computational Argumentation, USA*, pp. 20–26 (2009)
5. Emele, C.D., Norman, T.J., Guerin, F., Parsons, S.: Learning policies through argumentation-derived evidence (extended abstract). In: van der Hoek, Lesprance, Kaminka, Luck, Sen (eds.) *Proc. of 9th Intl. Conf. on Autonomous Agents and Multiagent Systems (AAMAS 2010)*, Toronto, Canada (2010)
6. Friedman-Hill, E.: *Jess in Action*. Manning (2003)
7. Kelemen, A., Liang, Y., Franklin, S.: A comparative study of different machine learning approaches for decision making. In: Mastorakis, N.E. (ed.) *Recent Advances in Simulation. Computational Methods and Soft Computing*, pp. 181–186. WSEAS Press, Piraeus (2002)

8. Mitchell, T.M.: *Machine Learning*. McGraw-Hill, New York (1997)
9. Možina, M., Žabkar, J., Bratko, I.: Argument based machine learning. *Artificial Intelligence* 171(10-15), 922–937 (2007)
10. Oren, N., Norman, T.J., Preece, A.: Loose lips sink ships: A heuristic for argumentation. In: Maudet, N., Parsons, S., Rahwan, I. (eds.) *ArgMAS 2006*. LNCS (LNAI), vol. 4766, pp. 121–134. Springer, Heidelberg (2007)
11. Parsons, S., Jennings, N.R.: Negotiation through argumentation-A preliminary report. In: *Proc. of the 2nd Int'l Conference Multi-Agent Systems (ICMAS 1996)*, Kyoto, Japan, pp. 267–274 (1996)
12. Parsons, S., Wooldridge, M., Amgoud, L.: Properties and complexities of some formal inter-agent dialogues. *Journal of Logic and Comp.* 13(3), 347–376 (2003)
13. Perrussel, L., Doutre, S., Thévenin, J.-M., McBurney, P.: A persuasion dialog for gaining access to information. In: Rahwan, I., Parsons, S., Reed, C. (eds.) *Argumentation in Multi-Agent Systems*. LNCS (LNAI), vol. 4946, pp. 63–79. Springer, Heidelberg (2008)
14. Quinlan, J.R.: *C4.5: programs for machine learning*. Morgan Kaufmann Publishers Inc., San Francisco (1993)
15. Rovatsos, M., Rahwan, I., Fischer, F., Weiss, G.: Practical strategic reasoning and adaptation in rational argument-based negotiation. In: Parsons, S., Maudet, N., Moraitis, P., Rahwan, I. (eds.) *ArgMAS 2005*. LNCS (LNAI), vol. 4049, pp. 122–137. Springer, Heidelberg (2006)
16. Sycara, K., Norman, T.J., Giampapa, J.A., Kollingbaum, M.J., Burnett, C., Masato, D., McCallum, M., Strub, M.H.: Agent support for policy-driven collaborative mission planning. *The Computer Journal* 53(1), 528–540 (2009)
17. Utgoff, P.E.: Incremental induction of decision trees. *Machine Learning* 4(2), 161–186 (1989)
18. Walton, D.N., Krabbe, E.C.W.: *Commitment in Dialogue: Basic Concepts of Interpersonal Reasoning*. SUNY Press, USA (1995)
19. Witten, I.H., Frank, E.: *Data Mining: Practical machine learning tools and techniques*, 2nd edn. Morgan Kaufmann, San Francisco (2005)

# Argumentative Alternating Offers

Nabila Hadidi<sup>1</sup>, Yannis Dimopoulos<sup>2</sup>, and Pavlos Moraitis<sup>1</sup>

<sup>1</sup> Paris Descartes University, 45 rue des Saints-Pères, 75270 Paris 06, France  
{nabila.hadidi, pavlos}@mi.parisdescartes.fr

<sup>2</sup> University of Cyprus, 75 Kallipoleos Str. POBox 20537, Nicosia, Cyprus  
yannis@cs.ucy.ac.cy

**Abstract.** This paper presents an argumentative version of the well known alternating offers negotiation protocol. The negotiation mechanism is based on an abstract preference based argumentation framework where both epistemic and practical arguments are taken into consideration in order to decide about different strategic issues. Such issues are the offer that is proposed at each round, acceptance or refusal of an offer, concession or withdrawal from the negotiation. The argumentation framework shows clearly how offers are linked to practical arguments that support them, as well as how the latter are influenced by epistemic arguments. Moreover it illustrates how agents' argumentative theories evolution, due to the exchange of arguments, influences the negotiation outcome. Finally, a generic algorithm that implements a concession based negotiation strategy is presented.

**Keywords:** Argumentation, Negotiation.

## 1 Introduction

Negotiation is the process of looking for an agreement between two or several agents on one or more issues. There exist three main approaches to negotiation, namely game theory (see e.g. [11]), heuristics (see e.g. [6]) and argumentation (see e.g. [10,13]).

In the last years there is a plethora of works on English on argumentation based negotiation (see e.g. [12,7,8]), testifying the increasing importance that is attached to the role of argumentation in negotiation. Although a precise and formal account of the added value of argumentation in negotiation is still missing, it is at least clear that exchanging arguments revealing (at least some of) the reasons for which a negotiator is proposing an offer may release several blocked situations. Such an example is a situation where the conflict is due to different perceptions of the world, which may have further repercussions on the behavior of a negotiator, including even parameters like his own preferences. Indeed, arguments received by the opponent on some issue might provide a piece of missing information to the proponent who could suddenly discover that the proposed offer is not optimal for himself, or that there is an objective constraint that forbids his opponent to accept his offer.

It is, therefore, evident that trying to "influence", in one way or another, the opponent's beliefs about the world may be a meaningful way to defend or attack an offer. This situation can be handled through the simultaneous consideration of both *practical* and *epistemic* arguments in the reasoning process and by deciding in which situation each type of argument must prevail. This may be part of the strategy of the agent. We remind that practical arguments support offers while epistemic arguments represent what the agent believes about the world.

The above intuitions define the perspective that is taken in this paper. To capture these intuition, we propose an original adaptation of the well known *alternating offers protocol* [12] in the argumentation context. Then, we adapt a reasoning mechanism combining practical and epistemic arguments proposed in [3] in the negotiation context and we exploit the possibilities that our *argumentative alternating offers protocol* provides for alternating practical and epistemic arguments depending on the evolution of the negotiation. Finally, we present a generic algorithm, which, building on the above reasoning mechanism and the possibilities that the argumentative alternating offers protocol provides, implements a parameterized concession based negotiation strategy. The algorithm is generic in the sense it can operate regardless of whether there is a time constraint or not (which is the case in this paper), or of the tactics (or heuristics) the agents might use in several situations where a choice has to be made (e.g. accept or reject an offer, choose the best offer to propose). Thus, it can be parameterized to capture the previous issues without further modification.

To the best of our knowledge, it is the first time that the way that epistemic arguments interfere with practical arguments in a negotiation process is presented along with a generic algorithm that incorporates this mechanism in the service of strategic considerations. This seriously differentiates our work from other important works in the domain such as [1,2,7,9], etc. A similar combination of epistemic and practical arguments is proposed in [4] but in a deliberation dialogue.

## 2 Negotiation Framework

The negotiation framework we propose is based on the abstract preference-based argumentation framework of [2].

We assume two agents,  $ag_i$  and  $ag_j$ ,  $i \neq j$ , that are involved in a bilateral negotiation over a set of offers (options)  $O = \{o_1, o_2, \dots, o_n\}$  which are identified from a logical language  $\mathcal{L}$ . We further assume that there is an option  $o_D \in O$  that represents disagreement. The options are mutually exclusive, which means that each agent can choose only one of them at once.

### 2.1 Arguments

From the language  $\mathcal{L}$  a set of arguments  $Args(\mathcal{L})$  are constructed. By argument we mean a *reason* for *believing* or *doing something*. We assume that an agent is aware of all the arguments of the set  $Args(\mathcal{L})$ . It encodes the fact that when an agent receives an argument from another agent, it can interpret it correctly and it can also compare it with its own arguments.

**Types of Arguments.** Unlike [2], we distinguish between *epistemic* and *practical* arguments, that are both taken into account, as in [3], in the reasoning mechanism used by the agents. Thus, we have:

1. *Practical arguments*  $A_p$  support offers (or decisions) by trying to justify those offers.
2. *Epistemic arguments*  $A_e$  represent what the agent believes about the world

In what follows, we are not interested in the construction of these arguments. We make the following assumptions:

- $Args(\mathcal{L}) = A_e \cup A_p$ ,
- $A_e \cap A_p = \emptyset$ ,
- Arguments structure is unknown.

Epistemic arguments are denoted by variables  $\alpha_1, \alpha_2, \dots$ , while practical arguments by variables  $\delta_1, \delta_2, \dots$ . When no distinction is necessary between arguments, we use variables  $a, b, c, \dots$

Let  $F$  be a function that maps each option to the arguments that support it, i.e.,  $\forall o \in O, F(o) \subseteq A_p$ . Each argument can support only one option, thus  $\forall o_y, o_z \in O, o_y \neq o_z, F(o_y) \cap F(o_z) = \emptyset$ . When  $\delta \in F(o)$ , we say that  $o$  is the conclusion of  $\delta$ , noted  $Conc(\delta) = o$ .

**Comparison Between Arguments.** As in [3], we assume three binary preference relations on arguments.

- $\succ_e$ : Partial preorder on the set  $A_e$ ,
- $\succ_p$ : Partial preorder on the set  $A_p$ ,
- $\succ_m$ : defined on the sets  $A_e$  and  $A_p$ , such that  $\forall \alpha \in A_e, \forall \delta \in A_p, (\alpha, \delta) \in \succ_m$  and  $(\delta, \alpha) \notin \succ_m$ . That means that any epistemic argument is stronger (preferred) than any practical argument ( $m$  stands for mixed relation).

In what follows  $\succ_x$  with  $x \in \{e, p, m\}$  denotes the strict relation associated with  $\succeq_x$ . It is defined as  $(a, b) \in \succ_x$  iff  $(a, b) \in \succeq_x$  and  $(b, a) \notin \succeq_x$ . Moreover when  $(a, b) \in \succeq_x$  and  $(b, a) \in \succeq_x$  we will say that the arguments  $a$  and  $b$  are *indifferent*, denoted by  $a \sim b$ .

**Conflict Between Arguments.** Conflicts between arguments in  $\mathcal{A} = A_p \cup A_e$  are captured by the binary relation  $R$  ([3]).

- $R_e$ : Represents the conflicts between arguments in  $A_e$ .
- $R_p$ : Represents the conflict between practical arguments, such that  $R_p = \{(\delta, \delta') \mid \delta, \delta' \in A_p, \delta \neq \delta' \text{ and } Conc(\delta) \neq Conc(\delta')\}$ . This relation is symmetric.
- $R_m$ : Represents the conflicts between epistemic and practical arguments s.t.  $(\alpha, \delta) \in R_m, \alpha \in A_e$  and  $\delta \in A_p$ .

Thus we have  $R = R_e \cup R_p \cup R_m$ .

We assume that practical arguments supporting different offers are in conflict. Thus for any two offers  $o_y, o_z$ ,  $\forall a \in F(o_y)$  and  $\forall a' \in F(o_z)$ , it holds that  $(a, a') \in R_p$  and  $(a', a) \in R_p$ .

**Attacks Between Arguments (Defeat).** Each preference relation  $\succeq_x$  (with  $x \in \{e, p, m\}$ ) is combined with the relation of conflict  $R_x$  (with  $x \in \{e, p, m\}$ ), to give a *defeat relation* between arguments, noted  $Def_x$  (with  $x \in \{e, p, m\}$ ).

**Definition 1.** (*Defeat*) Let  $A \subseteq \text{Args}(\mathcal{L})$  be a set of arguments and  $a, b \in A$ . Then  $(a, b) \in Def_x$  iff  $(a, b) \in R_x$ , and  $(b, a) \notin \succ_x$ .

We have  $Def_{global} = Def_e \cup Def_p \cup Def_m$ . In the following sections we will need two particular notions of defeat namely *rebuttal* and *undercutting*. For explaining those notions we will consider here a particular structure of arguments based on a propositional language  $\mathcal{L}'$  although our negotiation framework is independent of the structure of the arguments.  $\vdash$  stands for classical inference and  $\equiv$  for logical equivalence.

**Definition 2.** (*Argument Structure*) An argument is a pair  $a = (S, q)$  where  $q$  is a formula in  $\mathcal{L}'$  and  $S$  a set of formulae in  $\mathcal{L}'$  s.t.

- $S$  is consistent
- $S \vdash q$
- $S$  is a minimal set of propositions that satisfies the two previous conditions

Here  $S$  is called the support of the argument  $a$  and it is written  $S = \text{Support}(a)$  and  $q$  its conclusion and it is written  $q = \text{Conclusion}(a)$ .

**Definition 3.** (*Undercutting*) Let  $a$  and  $b$  be two arguments. Argument  $a$  undercuts  $b$  iff  $\exists p \in \text{Support}(b)$  s.t.  $p \equiv \neg \text{Conclusion}(a)$ .

**Definition 4.** (*Rebuttal*) Let  $a$  and  $b$  be two arguments. Argument  $a$  rebuts  $b$  iff  $\text{Conclusion}(a) \equiv \neg \text{Conclusion}(b)$ .

In the context of a negotiation, *practical* arguments *rebut* practical arguments, *epistemic* arguments *undercut* practical arguments, whereas *epistemic* arguments can both *undercut* and *rebut* other epistemic arguments. Recall that practical arguments cannot attack epistemic arguments.

## 2.2 Extensions of Arguments

In [5], different acceptability semantics have been introduced for computing the status of arguments. These are based on two basic concepts, *defence* and *conflict-freeness*, defined as follows:

**Definition 5.** (*Defence/Conflict-free*) Let  $\mathcal{T} = \langle A, Def \rangle$  be an argumentation system with  $A \subseteq \text{Args}(\mathcal{L})$ . Let  $A' \subseteq A$ .

- $A'$  is conflict free iff  $\nexists a, b \in A'$  s.t.  $(a, b) \in Def$ .
- $A'$  defends  $a \in A$  iff  $\forall b \in A$ , if  $(b, a) \in Def$ , then  $\exists c \in A'$  s.t.  $(c, b) \in Def$ .

**Definition 6.** (*Acceptability semantics*) Let  $\mathcal{T} = \langle A, Def \rangle$  be an argumentation system with  $A \subseteq \text{Args}(\mathcal{L})$  and  $A'$  a conflict free set of arguments.

- $A'$  is an admissible extension iff  $A'$  defends any element in  $A'$ .
- $A'$  is a preferred extension iff  $A'$  is a maximal(w.r.t set  $\subseteq$ ) admissible set.
- $A'$  is a stable extension iff it is a preferred extension that defeats any argument in  $A \setminus A'$ .

**Definition 7.** (*Argument status*) Let  $\mathcal{T} = \langle A, Def \rangle$  be an argumentation system with  $A \subseteq \text{Args}(\mathcal{L})$  and  $E_1, E_2, \dots, E_n$  its extensions under a given semantics. Let an argument  $a \in A$ .

- $a$  is skeptically accepted iff  $\forall E_q, 1 \leq q \leq n, a \in E_q$ .
- $a$  is credulously accepted iff  $\exists E_q, 1 \leq q \leq n, s.t a \in E_q$  and  $\exists E_w, 1 \leq w \leq n, s.t a \notin E_w$ .
- $a$  is rejected iff  $\nexists E_q, 1 \leq q \leq n, such that a \in E_q$ .

### 2.3 Negotiating Agents Theories

As in [2], we assume that each agent involved in a negotiation has a negotiation theory that contains arguments  $\mathcal{A}$  that can be exchanged during the negotiation. However, in our work we distinguish two types of arguments, i.e  $\mathcal{A} = A_p \cup A_e$ . This, as it will become evident in the following, has several effects on the reasoning process of the agents and consequently the negotiation process. Formally, a negotiation theory is defined as follows.

**Definition 8.** (*Negotiation theory*) Let  $O$  be a set of options,  $ag \in Ag$  an agent and  $Ag$  the set of negotiating agents. The negotiation theory  $\mathcal{T}^{ag}$  of agent  $ag$  is a tuple  $\mathcal{T}^{ag} = \langle \mathcal{A}^{ag}, F^{ag}, Def_{global}^{ag} \rangle$  where  $Def_{global}^{ag} = Def_e \cup Def_p \cup Def_m$  and  $\mathcal{A}^{ag} = A_p^{ag} \cup A_e^{ag}$  such that:

- $\mathcal{A}^{ag} \subseteq \text{Args}(\mathcal{L})$ . This set represents all the arguments that the agent can built from his beliefs and all the arguments that support each option in  $O$ .
- $F^{ag} : O \rightarrow 2^{A_p^{ag}}$  associates practical arguments to offers. We have

$$\bigcup_{1 \leq y \leq n} F^{ag}(o_y) = A_p^{ag}.$$

- $Def_{global}^{ag} \subseteq \mathcal{A}^{ag} \times \mathcal{A}^{ag}$

### 2.4 Offer Status and Preferences between Offers

In [3], five statuses are defined for the options/offers. In this work, we use only two of them. A *skeptical offer* which has a supporting argument that is skeptically accepted, and a *credulous offer* which has a supporting argument that is credulously accepted.

The *Effective Supporting Arguments* of an offer, defined formally below, are all arguments, either *skeptically* or *credulously* accepted, that support the offer.

**Definition 9.** (*Effective Supporting Arguments-ESA*) Let  $O$  be a set of offers,  $E_1, \dots, E_n$  the extensions under a given semantics of the theory  $\mathcal{T} = \langle \mathcal{A}, F, Def_{global} \rangle$  and  $o_y \in O$  an offer. Then the set of effective supporting arguments for offer  $o_y$  is  $ESA(o_y) = \{a \mid a \in F(o_y) \text{ and } a \in E_1 \cup \dots \cup E_n\}$ .

In simple words,  $ESA(o_y)$  is the set of arguments that support  $o_y$  and are included in at least one extension. The cardinalities of the ESA of the offers can be used to define a *preference relation* on these offers.

**Definition 10.** Let  $O$  be a set of offers,  $\mathcal{T} = \langle \mathcal{A}, F, Def_{global} \rangle$  a negotiation theory, and  $o_x, o_y \in O$ . Then  $o_x \succeq o_y$  iff  $\forall a \in ESA(o_x)$  and  $\forall b \in ESA(o_y)$  it holds that  $a \sim b$  (i.e. they are indifferent) and  $|ESA(o_x)| \geq |ESA(o_y)|$ .

Therefore,  $\succeq$  favors options that are supported by more arguments. Although this is a simple preference relation, and possibly more sophisticated methods for ranking offers exist, it suffices for the purposes of this work.

### 3 Argumentation-Based Alternating Offers Protocol

In [12], Rubinstein introduced the *Alternating Offers protocol* for bargaining between agents. It is a bilateral protocol between the *proposer* who initiates the process, and the *responder* who replies to the proposal. The proposer starts the negotiation process by presenting a *proposal* using a SUBMITPROPOSAL message. The responder can *accept* or *reject* the offer in its entirety by sending an ACCEPT or REJECT message as a reply. The responder can also propose a counter-offer by sending the COUNTER reply accompanied by the counter proposal. In this case, the proposer has the same options and therefore can accept, reject or reply with a counter proposal of its own. If one of the agents is satisfied with the current iteration of the proposal, he can send an ACCEPT message to the other. He can also signal his dissatisfaction and abort the negotiation session by sending a REJECT message. To seal the agreement, the other agent has to send a CONFIRM message and receive a CONFIRM-ACCEPTANCE message in reply.

The protocol, as described above, is generic, with no time limits and no central coordinator to manage the negotiations, and either of the parties can leave the process at any time.

In this work we adapt the classical alternating offers protocol to the case of argumentation-based negotiation. To do so we extend the concept of *round* used in the classical protocol to include, besides the classical *propose*, *accept* or *reject* messages, the possibility to *argue* in order to defend or attack an offer. In addition, *propose* and *argue* are accompanied by *supporting (practical or epistemic) arguments*.



### 3.1 Moves

Arguments and offers are conveyed through *dialogue moves* (or simply *moves*). A move is denoted by  $m_{r,g}$ , whereas  $r \geq 1$  identifies the round (and therefore the offer which is currently discussed), and  $g \geq 1$  the number (order) of the move in that round. In the argumentative alternating offers protocol the following moves are used. In all moves  $ag_i$  and  $ag_j$  are the participating agents and  $o_y \in O$ .

- *Propose*( $ag_i, ag_j, o_y, \delta$ ), where  $\delta \in F^{ag_i}(o_y)$ . This move allows agent  $ag_i$  to propose an offer  $o_y$  to agent  $ag_j$ , along with a practical argument  $\delta$  that supports it.
- *Argue*( $ag_i, ag_j, a, Target$ ), where  $a \in A^{ag_i}$  and *Target* is the move the argument of which is attacked by  $a$  or nil. This move allows agent  $ag_i$  to argue by defending his own offer  $o_y$  or to counter-attack an offer sent by  $ag_j$ . The arguments used in this move satisfy the following conditions
  - If *Target* = nil then  $a \in F^{ag_i}(o_y)$ , i.e.,  $a$  is a practical argument that support the offer  $o_y$ .
  - If *Target*  $\neq$  nil then  $a \in A_e^{ag_i}$  is an argument presented against the argument of *Target*. Thus, an agent can't present an argument against his own arguments.
- *Reject*( $ag_i, ag_j, o_y$ ). This move is sent by  $ag_i$  to inform  $ag_j$  that he has no arguments to present and he does not accept  $ag_j$ 's offer.
- *Nothing*( $ag_i, ag_j$ ). This move notifies  $ag_j$  that  $ag_i$  has no arguments to present and he either still considers his offer as a most preferred one for him (when he is the proposer), or believes that he has better options than the current offer (when he is the recipient of an offer sent by the other agent).
- *Accept*( $ag_i, ag_j, o_y$ ). This move is used by agent  $ag_i$  to notify that he accepts the offer  $o_y$  made by  $ag_j$ .
- *Agree*( $ag_i, ag_j$ ). This move means that  $ag_i$  now believes that his current offer is not optimal for himself and therefore accepts the arguments sent by  $ag_j$ . Agent  $ag_j$  starts a new round.
- *Withdraw*( $ag_i, ag_j$ ). This move indicates that agent  $ag_i$  withdraws from negotiation.
- *final*( $ag_i, ag_j$ ). This is a shorthand for *Propose*( $ag_i, ag_j, o_y, \emptyset$ ) and is used during a final round of the negotiation. Its use and semantics will become apparent in the following.

The following functions retrieve the parameters of the moves.

- *Performative*( $m_{r,g}$ ) returns one of *Propose, Argue, Nothing, Reject, Accept, Withdraw, Agree*.
- *Agent*( $m_{r,g}$ ) returns the agent who sent the move.
- *Offer*( $m_{r,g}$ ) returns the offer sent in the round  $r$ .
- *Argument*( $m_{r,g}$ ) returns the argument sent to the other agent.
- *Targ*( $m_{r,g}$ ) returns the target of the move.

Finally, the following hold.

- If  $Performative(m_{r,g})=Propose$  then  $Argument(m_{r,g}) \in A_p^{agi}$  arguments
- If  $Performative(m_{r,g})=Argue$  then  $Argument(m_{r,g}) \in A_e^{agi} \cup A_p^{agi}$

### 3.2 Round

A round takes place in alternating way between two agents  $P$  (the proposer of the offer) and  $R$  (the recipient of the offer). The agent proposing an offer may send moves with performative from  $\{Propose, Argue, Agree, Nothing, Withdraw\}$ , whereas the agent that receives an offer may send moves with performative from  $\{Argue, Reject, Accept, Nothing, Withdraw\}$ . A round is defined formally as follows.

**Definition 11.** (Round) *A round  $r$  between two agents  $P$  and  $R$  is a non empty sequence of moves  $m_{r,1}, \dots, m_{r,n}$ , such that:*

- $\forall i, k, i \neq k, \forall g, g', g \neq g' Offer(m_{i,g}) \neq Offer(m_{k,g'})$ .
- $\forall r, Agent(m_{r,g}) = P$  if  $Odd(g)$ , and  $Agent(m_{r,g}) = R$  if  $Even(g)$ .
- $\forall m_{r,g}$ , if  $Odd(g)$  then  $Performative(m_{r,g}) \in \{Propose, Argue, Agree, Nothing, Withdraw\}$ .
- $\forall m_{r,g}$ , if  $Even(g)$  then  $Performative(m_{r,g}) \in \{Argue, Reject, Accept, Nothing, Withdraw\}$ .
- $\forall r, Performative(m_{r,1}) \in \{Propose, Withdraw\}$ .
- $\forall r$ , if  $Performative(m_{r,g}) = Performative(m_{r,g+1}) = Withdraw$  then the dialogue ends with a disagreement.
- $\forall m_{r,g}$ , if  $Performative(m_{r,g})=Argue$  then:
  - If  $Targ(m_{r,g}) \neq nil$  then  $Targ(m_{r,g})=m_{r,g'}$  with  $g' < g$ ,  $Argument(m_{r,g}) \in Def_{global}^{Agent(m_{r,g})} Argument(m_{r,g'})$  and  $Agent(m_{r,g}) \neq Agent(m_{r,g'})$ . Here the agent sends an argument which attacks one presented previously by the other agent in the same round.
  - Else  $Agent(m_{r,g})=Agent(m_{r,1})$  and  $Argument(m_{r,g}) \in F^{Agent(m_{r,g})}(Offer(m_{r,1}))$ . Here the agent sends a new argument to support his offer.
- If  $Performative(m_{r,n}) = Accept$  then  $Offer(m_{r,1})$  is the outcome of the dialogue which terminates with agreement.

- If  $Performative(m_{r,n}) \in \{Agree, Reject\}$  then a new round  $r + 1$  starts with  $Agent(m_{r+1,1}) \neq Agent(m_{r,1})$  i.e. with the other agent as proposer.
- $\forall m_{r,g}$ , if  $Performative(m_{r,g}) = \text{Nothing}$  then  $Argument(m_{r,g}) = \emptyset$  and  $Offer(m_{r,1}) = \emptyset$ .
- $\forall m_{r,1}, m_{r,g'}, g' > 1$  if  $Offer(m_{r,1}) = Offer(m_{r,g'})$  then  $Agent(m_{r,1}) = Agent(m_{r,g'})$  and  $Argument(m_{r,1}) \neq Argument(m_{r,g'})$ . In our protocol, unlike [2], an agent can propose the same offer more than once during a round provided that he supports it with an argument not used before.

**Definition 12.** (*Argumentative alternating offers dialogue*) An argumentative alternating offers dialogue  $d$  between two agents  $P, R$  is a non-empty sequence of rounds  $d = \{r_1 \dots r_\lambda\}$  between  $P$  and  $R$ .

In the alternating offers protocol ([12]), two outcomes are possible: (a) no agreement (disagreement), or (b) an agreement in some round. In the argumentative protocol the situation is similar.

**Definition 13.** (*Outcome*) Let  $d = \{r_1, \dots, r_n\} \in D$  be an argumentative alternating offers dialogue where  $D$  is the set of all the dialogues built from the argumentative alternating offers protocol and  $r_n = \{m_{r_n,1}, \dots, m_{r_n,k}\}$  be the last round of  $d$ . If  $Performative(m_{r_n,k}) = \text{Accept}$  then  $Outcome(d) = Offer(m_{r_n,1})$  (*Agreement*). Else  $Outcome(d) = \text{nil}$  (*Disagreement*).

## 4 Negotiation Strategy

In this section, we present a strategy that can be used by the two agents involved in an argumentative alternating offers negotiation. The strategy is based on the theory of the agent  $T$ , his preference on the set of offers  $\succeq$ , and the alternating offers protocol as defined in the previous section.

In order to improve presentation, some of the parameters of the messages of the negotiation dialogue are omitted from the algorithms that follow. These are mainly agent and target move names, and are easily derivable from context.

The main procedure of the strategy is described by procedure `negotiate` ( $T, O, \text{outcome}$ ), depicted in Algorithm 4. It accepts as parameters the agent theory  $T$ , and the set of possible offers  $O$ , and returns an *outcome* that can be either an offer, when an agreement is reached, or *nil* when the negotiation fails. As noted before, the set  $O$  contains an option  $o_D$  representing the possibility that the agent leaves the negotiation without an agreement, and therefore remains in the same state that he was initially. Therefore, offers that lead to situations that are less desirable than his current state are less preferred by the agent. This option  $o_D$  corresponds to what in classical negotiation theory is referred to as *reservation value*.

```

Algorithm 1: Procedure compute-best( $T_{r,g}, O, O^{best}$ )
begin
  Compute the extensions  $E_1, E_2, \dots, E_n$  of  $T_{r,g}$ ;
  Compute  $O^{cand} = \{o \mid o \in O \text{ s.t. } \exists a \in \cup_{i=1}^n E_i \text{ and } a \in F(o)\}$ ;
  Compute  $O^{best} = \{o \mid o \in O^{cand}, o_D \not\triangleright o, \text{ and } \neg \exists o' \in O^{cand} \text{ s.t. } o' \triangleright o\}$ ;
  return  $O^{best}$ ;
end

```

```

Algorithm 2: Procedure proposal( $T_{r,g}, O, o, a$ )
begin
  Call compute-best( $T_{r,g}, O, O^{best}$ );
  if  $O^{best} = \emptyset$  then
     $o = nil$ ;  $a = nil$ ;
  else
    Select an offer  $o$  from  $O^{best}$  and  $a \in F(o)$  such that
     $a$  belongs to some extension of  $T_{r,g}$ ;
  end
  return  $o, a$ ;
end

```

```

Algorithm 3: Procedure check( $T_{r,g}, O, o, R, UA$ )
begin
  Call compute-best( $T_{r,g}, O, O^{best}$ );
  if  $o \in O^{best}$  then
    Send accept;
  else
    Compute  $\mathcal{E}^{best} = \{E \mid E \text{ is an extension of } T_{r,g} \text{ s.t. } \exists a \in E \text{ and } a \in F(o) \text{ for } o \in O^{best}\}$ ;
    if there is  $a \in E, E \in \mathcal{E}^{best}$  s.t.  $a$  is an epistemic
    argument that attacks some argument  $b \in R$ ,
    and  $(a, b) \notin UA$  then
      Send argue( $a$ );
       $UA = UA \cup \{(a, b)\}$ ; return  $UA$ ;
    else
      Send nothing;
    end
  end
   $g = g + 1$ ;
end

```

```

Algorithm 4: Procedure negotiate( $T, O, outcome$ )
begin
   $r = 1; g = 1; own = false; T_{1,1} = T;$ 
   $Received = \emptyset; Offered = \emptyset; UsedAtt = \emptyset;$ 
  if Agent proposes first then
    Call proposal( $T_{1,1}, O, o^{cur}, a^{curr}$ );
    Send Propose( $o^{cur}, a^{cur}$ );  $own = true;$ 
  end
  while true do
     $g = g + 1;$  Get  $m_{r,g};$ 
    Incorporate argument( $m_{r,g}$ ) into  $T_{r,g};$ 
    switch Performative( $m_{r,g}$ ) do
      case Argue
        Add argument( $m_{r,g}$ ) to Received;
        if  $own$  then
          Call defend( $T_{r,g}, O, o^{cur}, Received, UsedAtt$ );
        else
          Call check( $T_{r,g}, O, o^{cur}, Received, UsedAtt$ );
        end
      case Propose
        Add argument( $m_{r,g}$ ) to Received;
         $o^{cur} = Offer(m_{r,g});$ 
        Add  $o^{cur}$  to Offered;
         $r = r + 1; g = 1;$ 
        Call check( $T_{r,g}, O, o^{cur}, Received, UsedAtt$ );
      case Agree
        Call proposal( $T_{r,g}, O, o^{cur}, a^{curr}$  );
        if  $o^{cur} = nil$  then
          Send withdraw;  $g = g + 1;$ 
        else
          Send Propose( $o^{cur}, a^{cur}$ );
           $Received = \emptyset; UsedAtt = \emptyset;$ 
           $r = r + 1; g = 1;$ 
           $o = o^{cur}; own = true;$ 
        end
      case Nothing
        Call nothing-
        reply( $T_{r,g}, O, own, o^{cur}, Received, UsedAtt$ );
      case Reject
         $O = O - \{o^{cur}\}; own = false;$ 
        Remove from  $T_{r,g}$  all arguments of  $F(o^{cur})$ 
      case Withdraw
        Call withdrawal( $T_{r,g}, O, Offered, outcome$ );
        return outcome and exit;
      case Accept
         $outcome = o^{cur};$ 
        return outcome and exit;
      case Final
         $outcome = Offer(m_{r,g});$ 
        return outcome and exit;
    end
  end
end
end

```

One of the agents initiates the negotiation by sending a proposal via a *propose* message. This proposal is selected by procedure  $\text{proposal}(T_{r,g}, O, o, a)$  (Algorithm 2). This selection at some round  $r$  and step  $g$ , is based on the current theory of the agent  $T_{r,g}$ , and the current set of offers  $O$ . The offer  $o$  that is proposed must be supported by some argument  $a$  that belongs to some of the extensions of  $T_{r,g}$ . Among several possible such offers, the best wrt  $\succeq$  is selected. Note that an agent never proposes, accepts or defends an offer that is worse wrt  $\succeq$  than  $o_D$ , as any such deal is considered by the agent worse than no deal.

```

Algorithm 5: Procedure  $\text{defend}(T_{r,g}, O, o, R, UA)$ 
begin
  Call  $\text{compute-best}(T_{r,g}, O, O^{best})$ ;
  if  $o \notin O^{best}$  then
    Send agree;
  else
    Compute  $\mathcal{E}^{best} = \{E \mid E \text{ is an extension of } T_{r,g}$ 
      s.t.  $\exists a \in E \text{ and } a \in \mathcal{F}(o) \text{ for } o \in O^{best}\}$ ;
    Compute  $\mathcal{E}^o = \{E \mid E \in \mathcal{E}^{best} \text{ and } \exists a \in E \text{ s.t. } a \in F(o)\}$ ;
    if there is  $a \in E$ ,  $E \in \mathcal{E}^o$  s.t.  $a$  is an epistemic
      argument that attacks some argument  $b \in R$ ,
      and  $(a, b) \notin UA$  then
      Send argue( $a$ );
       $UA = UA \cup \{(a, b)\}$ ;
      return  $UA$ ;
    else if there is  $a \in E$ ,  $E \in \mathcal{E}^o$  s.t.  $a$  is
      a practical argument that has not been used before
      and  $a \in F(o^{cur})$  then
      Send argue( $a$ );
    else
      Send nothing;
  end
   $g = g + 1$ ;
end

```

Upon receiving a proposal in a move  $m_{r,g}$ , the agent incorporates the supporting argument in his theory, adds the argument to the set *Received* of arguments that have been sent by the other agent, and runs procedure  $\text{check}(T_{r,g}, O, o, \text{Received}, \text{UsedAtt})$  (Algorithm 3). If the proposed offer is one of the best (wrt  $\succeq$ ), he accepts the offer and the negotiation terminates. Otherwise, he attempts to find an epistemic argument  $a$  that belongs to one of the extensions of  $T_{r,g}$ , and counterattacks the argument supporting the offer. Note that  $a$  must not have been used before to attack the supporting argument of the other agent in the same round. This avoids loops in argumentation, and is achieved by recording the counterattacks in *UsedAtt*. If he is successful, he sends argument  $a$  with an *argue* to the other agent.

If he is unsuccessful, he is confronted with a situation where on the one hand he can not counterattack the proposal, but on the other hand there are offers that are more desirable than the proposal. In such a case he sends a *nothing* message to the other agent.

Therefore, the reply to a proposal can be any of *accept*, *argue*, or *nothing*. The first case is straightforward. Whenever an agent receives an *argue* during a round during which he is the proposer, he runs procedure  $\text{defend}(T_{r,g}, O, o, \text{Received}, \text{UsedAtt})$  (Algorithm 5). If in the light of the last argument sent by the other agent his proposal is not one of the most preferred for himself, he replies with *agree*. Otherwise, he tries to defend his proposal against the attack by attacking one of the arguments sent by his counter-party during the current round. If no such attack exists, another argument supporting his offer is sent in a *argue* message. If no such argument exists, a *nothing* message is sent. This signifies that the agent insists that his current offer is one of the most preferred for himself. Upon receiving *nothing* the other agent sends a *reject* message, and becomes the proposer in the new round. This task is carried out by the part of procedure  $\text{nothing-reply}(T_{r,g}, O, \text{own}, o, \text{Received}, \text{UsedAtt})$  (Algorithm 6), which runs when parameter *own* is false, meaning that the offer currently discussed has been proposed by the other agent.

```

Algorithm 6: Procedure  $\text{nothing-reply}(T_{r,g}, O, \text{own}, o, R, UA)$ 
begin
  if own then
    Compute  $\mathcal{E}^o = \{E \mid E \text{ is an extension of } T_{r,g} \text{ and } \exists a \in E \text{ s.t. } a \in F(o)\}$ ;
    if there is  $a \in E$ ,  $E \in \mathcal{E}^o$ ,  $a \in \mathcal{F}(o)$  and  $a$  has not
      been used before then
      Send argue( $a$ );
    else
      Send nothing;
    end
     $g = g + 1$ ;
  else
    Send reject;  $g = g + 1$ ;
    Call Proposal( $T_{r,g}, O, o^{cur}, a^{curr}$ );
    if  $o^{cur} = \text{nil}$  then
      Send withdraw;  $g = g + 1$ ;
    else
      Send Propose( $o^{cur}, a^{cur}$ );
       $R = \emptyset$ ;  $UA = \emptyset$ ;
       $r = r + 1$ ;  $g = 1$ ;
       $o = o^{cur}$ ;  $\text{own} = \text{true}$ ;
      return  $o, \text{own}, R, UA$ ;
    end
  end
end
end

```

If an agent receives a *nothing* message in a round where he is the proposer, he is in a situation where he can not defend the argument that supports his offer. Therefore, he needs to find some other argument to support it, which he sends in an *argue* message. If no such argument exists, a *nothing* message is sent.

If at some point one of the agent has no offers, he sends a *withdraw* message, signifying that he is willing to leave the negotiation. This triggers a final round of negotiation that is carried out by procedure  $\text{withdrawal}(T_{r,g}, O, Offered, out)$  (Algorithm 7) that selects one of the offers from the input set *Offered*. This set contains all the offers proposed during the negotiation by the agent who wishes to withdraw. The agent who receives the withdraw message, finds the best offer contained in *Offered*. If this is better than disagreement, he sends it in a *final* message. The negotiation terminates with agreement if this offer is still better than  $o_D$  for the other agent, otherwise it terminates with disagreement.

Algorithm 7: Procedure  $\text{withdrawal}(T_{r,g}, O, Offered, out)$

```

begin
  if  $\text{Performative}(m_{r,g-2}) = \text{withdraw}$  then
     $out = \text{nil}$ ;
  else
    Select  $o \in Offered$  s.t.  $o \triangleright o_D$  and  $\neg \exists o' \in Offered$  s.t.  $o' \triangleright o_D$ 
    and  $o' \triangleright o$ ;
    if  $o$  exists then
      Send  $\text{final}(o)$ ;  $out = o$ ;
    else
      Send  $\text{withdraw}$ ;  $out = \text{nil}$ ;
    end
  end
end
return  $out$ ;
end

```

It is worth noting that although the above algorithms implement a specific negotiation strategy, the overall process they describe is generic in the sense that it can easily be adapted to accommodate other strategies. Consider for instance the case where one of the agents receives a *reject* message to an offer he made in some previous move. In the current version of procedure *negotiate* he removes his offer and the other agent takes turn. Moreover, in the next round he will *concede*, by sending his next best offer. All these are strategic decisions that can easily be modified without altering in any way the working of the overall algorithms, and more importantly the role of argumentation in negotiation.

Moreover, the argumentative alternating offers protocol we propose has two useful properties. The first property is *soundness*. This property guarantees that any offer agreed by the two agents through the argumentative alternating offers protocol is better for both agents than the offer that corresponds to disagreement i.e.  $o_D$ . More formally:



**Proposition 1.** (*Soundness*) Let  $d = \{r_1, \dots, r_n\}$  be an argumentative alternating offers dialogue between two agents  $\alpha$  and  $\beta$ . If  $\text{Outcome}(d) = o$ ,  $o \neq \text{nil}$  then  $o \triangleright o_D^\alpha$  and  $o \triangleright o_D^\beta$ .

Another interesting property of the argumentative alternating offers protocol is that any negotiation dialogue produced through this protocol terminates.

**Proposition 2.**  $\forall d \in \mathcal{D}$  where  $\mathcal{D}$  is the set of all the dialogues built from the alternating offers protocol,  $d$  terminates.

## 5 Example

For illustrating our negotiation algorithm we will use a simple scenario where a buyer ( $ag_b$ ) and a seller ( $ag_s$ ) negotiate over the price of a product. The set of options is  $O = \{o_1, o_2, o_3, o_D\}$ , where  $o_1 = \text{high}$ ,  $o_2 = \text{medium}$ ,  $o_3 = \text{low}$  are referring to the price of the product, and  $o_D$  represents the options of not selling (buying) for the seller (buyer).

Assume that the seller prefers a high price to a medium price to not selling to a low price. Symmetrically, the buyer prefers a low price to a medium price to not buying to a high price. We also assume that we are in a high-season period, but the buyer agent is not aware of that before the negotiation. Both agents represent their knowledge in some propositional language  $\mathcal{R}$ .

Assume that the buyer has the following knowledge:

*regular\_customer*  
*regular\_customer*  $\rightarrow$  *discount*  
*discount*  $\rightarrow$  *buy2*, *buy2*  $\rightarrow$   $o_3$ .  
*high\_season*  $\rightarrow$   $o_1$   
*high\_season*  $\rightarrow$   $\neg$ *discount*  
*high\_season*  $\wedge$  *regular\_customer*  $\rightarrow$   $o_2$

From this knowledge base the agent can construct one practical argument  $\delta_1 = (\{\textit{regular\_customer}, \textit{regular\_customer} \rightarrow \textit{discount}, \textit{discount} \rightarrow \textit{buy2}, \textit{buy2} \rightarrow o_3\}, o_3)$  that supports  $o_3$ . Two epistemic arguments can also be constructed:  $\alpha_1 = (\{\textit{regular\_customer}, \textit{regular\_customer} \rightarrow \textit{discount}\}, \textit{discount})$ , and  $\alpha_2 = (\{\textit{regular\_customer}, \textit{regular\_customer} \rightarrow \textit{discount}, \textit{discount} \rightarrow \textit{buy2}\}, \textit{buy2})$ . We have therefore only one extension  $E = \{\delta_1, \alpha_1, \alpha_2\}$ . Thus, the option  $o_3$  is skeptical and  $o_2, o_1$  are rejected.

Assume now that the seller has the following knowledge:

*high\_season*  
*high\_season*  $\rightarrow$   $\neg$ *discount*, *high\_season*  $\rightarrow$   $o_1$ .  
*high\_season*  $\wedge$  *regular\_customer*  $\rightarrow$   $o_2$ .  
*sales\_season*  $\wedge$  *regular\_customer*  $\rightarrow$   $o_3$ .

The seller agent has one practical argument  $\delta_2 = (\{\textit{high\_season}, \textit{high\_season} \rightarrow o_1\}, o_1)$  which supports  $o_1$ . He has also one epistemic argument  $\alpha_3 = (\{\textit{high\_season}, \textit{high\_season} \rightarrow \neg \textit{discount}\}, \neg \textit{discount})$ .

We have therefore only one extension  $E = \{\delta_2, \alpha_3\}$ . Thus, the offer  $o_1$  is skeptical and  $o_2, o_3$  are rejected.

Supposing that  $ag_b$  begins the negotiation, the dialogue between  $ag_b$  and  $ag_s$  will be as follows:

$m_{1,1}$ :*Propose*( $ag_b, ag_s, o_3, \delta_1$ )  
 $m_{1,2}$ :*Argue*( $ag_s, ag_b, \alpha_3, (m_{1,1})$ )  
 $m_{1,3}$ :*Agree*( $ag_b, ag_s$ )  
 $m_{2,1}$ :*Propose*( $ag_s, ag_b, o_1, \delta_2$ )  
 $m_{2,2}$ :*Nothing*( $ag_b, ag_s$ )  
 $m_{2,3}$ :*Nothing*( $ag_s, ag_b$ )  
 $m_{2,4}$ :*Reject*( $ag_b, ag_s$ )  
 $m_{3,1}$ :*Propose*( $ag_b, ag_s, o_2, \delta_3$ )  
 $m_{3,2}$ :*Accept*( $ag_s, ag_b, o_2$ )

The buyer agent proposes first his optimal offer which is  $o_3$  with his supporting argument  $\delta_1$ . The seller agent updates his theory which now contains the arguments  $\delta_2, \delta_3 = (\{ \textit{high\_season, regular\_customer, high\_season} \wedge \textit{regular\_customer} \rightarrow o_2 \}, o_2), \alpha_3, \delta_1, \alpha_1$  and  $\alpha_2$ . For the defeat relation of the seller agent we have the following situation:

- $(\delta_2, \delta_3), (\delta_2, \delta_1), (\delta_3, \delta_1)$  because the conclusions of  $\delta_2, \delta_3$  and  $\delta_1$  are not the same (and therefore conflicting) and also because the preferences of the agent are  $\delta_2 \succ_p \delta_3, \delta_2 \succ_p \delta_1$  and  $\delta_3 \succ_p \delta_1$ .
- $(\alpha_3, \delta_1)$  because there is undercutting and  $\alpha_3 \succ_m \delta_1$ .
- $(\alpha_3, \alpha_2)$  and  $(\alpha_3, \alpha_1)$  because there is undercutting between  $\alpha_3$  and  $\alpha_2$ , rebuttal between  $\alpha_3$ , and  $\alpha_1$  and the preferences of the agent are  $\alpha_3 \succ_e \alpha_2$  and  $\alpha_3 \succ_e \alpha_1$ .

The theory of the seller agent has one extension,  $E = \{\alpha_3, \delta_2\}$ , and therefore the seller tries to defeat with  $\alpha_3$  the argument he received in the last move.

When the buyer receives move *argue*( $\alpha_3$ ), he first updates his theory. This theory now contains the arguments  $\alpha_1, \alpha_2, \alpha_3, \delta_1, \delta_2, \delta_3$ , whereas the defeat relation is as follows:

- $(\alpha_3, \alpha_2)$  and  $(\alpha_3, \alpha_1)$  because there is undercutting between  $\alpha_3$  and  $\alpha_2$ , rebuttal between  $\alpha_3$  and  $\alpha_1$ . The preferences of the agent in this context are now  $\alpha_3 \succ_e \alpha_2$  and  $\alpha_3 \succ_e \alpha_1$ .
- $(\delta_1, \delta_2), (\delta_1, \delta_3)$  and  $(\delta_3, \delta_2)$  because the conclusions of  $\delta_2, \delta_3$  and  $\delta_1$  are not the same and the preferences of the agent are  $\delta_1 \succ_p \delta_2, \delta_1 \succ_p \delta_3, \delta_3 \succ_p \delta_2$ .
- $(\alpha_3, \delta_1)$  because there is an undercutting attack and  $\alpha_3 \succ_m \delta_1$ .

The theory of the buyer after  $m_{1,2}$  has one extension,  $E = \{\alpha_3, \delta_3\}$ , and therefore  $o_3$  is not the best offer for him. Consequently, he agrees with the seller.

This initiates a new round, where the seller proposes the offer  $o_1$ , which is a skeptical conclusion of his theory. The buyer updates his theory with argument  $\delta_2$ , but there is no change. The offer  $o_2$  remains the best, and thus he needs to defeat the argument he received. But, none of the acceptable arguments defeats  $\delta_2$  and then he sends *nothing* to indicate that he does not accept the offer.

When the seller agent receives the *nothing* message, he attempts to find another argument which supports his offer  $o_1$ . As this attempt fails, he sends a *nothing* message to signify that he has not change his preference on offer  $o_1$  but he has no other argument for supporting it. The buyer ends the round with a *reject* and thus, a third round begins in which the buyer agent proposes his best offer  $o_2$  with  $\delta_3$  to support it. Here we assume that the seller agent is willing to concede, and therefore he accepts  $o_2$  because this is an acceptable offer for him. Thus the negotiation ends with an agreement.

## 6 Related Work and Conclusion

In the last years several works have appeared in the argumentation based negotiation literature. These works have focused on several aspects of negotiation such as the problem of decision making (see e.g. [17]), the study of specific types of negotiation such as interest based negotiation [9], whereas the work of [2] proposed a general framework for argumentation based negotiation where several interesting issues have been studied. These issues include the link between the status of the arguments and the offers they support, the definition of important concepts such as the concession and its impact on the evolution of the negotiation, etc. This work is the most relevant to ours. Nevertheless, there are some important differences. One of them is that in our paper we make a clear distinction between epistemic and practical arguments and, by adapting the work presented in [3], we show how epistemic arguments interfere with practical arguments in the definition of the acceptable offers. Then, we show how this reasoning mechanism can be used by the agents in the context of an original adaptation of the well known alternating offers protocol [12]. Another difference is that in our work we are interested in strategic issues. More precisely, we propose a generic algorithm that implements a strategy that can be used by both agents. This algorithm can be parameterized in different ways in order to capture, for example, different conditions of concession, different methods for ranking offers or different tactics for deciding whether withdrawing or making a concession. Our future work will address several open issues. One such issue is the study of several tactics for choosing the best offer to propose, especially in the context of time constraint negotiations. Another issue is the investigation of different methods for ranking the offers, whereas a third issue is that of the formal properties of the argumentative alternating offers protocol, apart from the ones of soundness and termination that we have already presented.

## References

1. Amgoud, L., Belabbes, S., Prade H.: Towards a formal framework for the search of a consensus between autonomous agents. In: 4th International Joint Conference on Autonomous Agents and Multi-Agents systems, pp. 537–543 (2005)
2. Amgoud, L., Dimopoulos, Y., Moraitis, P.: A Unified and General Framework for Argumentation-based Negotiation. In: 6th International Joint Conference on Autonomous Agents and Multiagent Systems, pp. 113–124 (2007)

3. Amgoud, L., Dimopoulos, Y., Moraitis, P.: Making Decisions through Preference-based Argumentation. In: 11 th International Conference on the Principles of Knowledge Representation and Reasoning, KR 2008, pp. 963–970 (2008)
4. Black, E., Atkinson, K.: Dialogues that account for different perspectives in collaborative argumentation. In: 8th International Conference on Autonomous Agents and Multi-Agents systems, pp. 867–874 (2009)
5. Dung, P. M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and  $n$ -person games. *Artificial Intelligence Journal* 77, 321–357 (1995)
6. Faratin, P., Sierra, C., Jennings, N.R.: Using similarity criteria to make issue trade-offs in automated negotiations. *Artificial Intelligence* 142(2), 205–237 (2002)
7. Kakas, A., Moraitis, P.: Adaptive Agent Negotiation via Argumentation. In: 5th International Joint Conference on Autonomous Agents and Multi-Agents systems, pp. 384–391 (2006)
8. Kraus, S., Sycara, K., Evenchik, A.: Reaching agreements through argumentation: a logical model and implementation. *Artificial Intelligence* 104, 1–69 (1998)
9. Rahwan, I., Pasquier, P., Sonenberg, L., Dignum, F.: On the benefits of exploiting underlying goals in argument-based negotiation. In: 22nd Conference on Artificial Intelligence, pp. 116–121 (2007)
10. Rahwan, I., Ramchurn, S.D., Jennings, N.R., McBurney, P., Parsons, S., Sonenberg, E.: Argumentation-based negotiation. *Knowledge Engineering Review* 18 (4), 343–375 (2003)
11. Rosenschein, J., Zlotkin, G.: *Rules of Encounter: Designing Conventions for Automated Negotiation Among Computers*. MIT Press, Cambridge, Massachusetts (1994)
12. Rubinstein, A.: Perfect equilibrium in a bargaining model. *Econometrica* 50(1), 97–109 (1982)
13. Sycara, K.: Persuasive argumentation in negotiation. *Theory and Decision* 28, 203–242 (1990)

# On a Computational Argumentation Framework for Agent Societies

Stella Heras, Vicente Botti, and Vicente Julián

Departamento de Sistemas Informáticos y computación  
Universitat Politècnica de València  
Camino de Vera s/n, 46022, Valencia, Spain  
sheras@dsic.upv.es

**Abstract.** In this paper, we analyse the requirements that argumentation frameworks should take into account to be applied in agent societies. Then, we propose a generic framework for the computational representation of argument information. It is able to represent different types of complex arguments in open multi-agent societies, where agents have social relations between them. In addition, we have formalised our framework by defining an argumentation framework based on it.

**ACM Categories:** Coherence and Coordination, Multi-Agent Systems.

**Keywords:** Agreement Technologies, Argumentation.

## 1 Introduction

A recent trend in Multi-Agent Systems (MAS) research is to broaden the applications of the paradigm to open MAS [20], where heterogeneous agents could enter into (or leave) the system, interact, form societies and adapt to changes in the environment. This and other paradigms for computing, such as grid computing or peer-to-peer technologies, have given rise to a new approach of computing as interaction [16]. This notion is used to define large complex systems in terms of the services that their entities or agents can offer and consume and consequently, in terms of the interactions between them. The high dynamism of these systems requires them to have a way of harmonising knowledge inconsistencies and reaching common agreements, for instance, when agents in an open MAS are faced with the goal of collaborating and solving a problem together.

Argumentation is probably the most natural way of harmonising conflicts. It provides a fruitful means of dealing with non-monotonic and defeasible reasoning. During the last decade, this important property has made many Artificial Intelligence (AI) researchers to pay attention on argumentation theory. In addition, research on argumentation is also at its peak in the Multi-Agent Systems (MAS) community, since it has been very successful to implement agents' internal and practical reasoning and to manage multi-agent dialogues [24]. Nowadays, argumentation is an active research area in AI and MAS [5].

However, most argumentation systems consider abstract notions of argument that are not intended for performing automated reasoning over them (automatic argument generation, selection and evaluation). In fact, the proposed computational argumentation frameworks take a narrow view of the argument structure [26]. On the other hand, most MAS whose agents have argumentation capabilities use ad-hoc and domain-dependent representations for arguments [30] [31]. Moreover, little work, if any, has been done to study the effect of the social relations between agents in the way that they argue and manage arguments. Commonly, the term *agent society* is used in the argumentation and AI literature as a synonym for an *agent organisation* [12] or a *group of agents* that play specific roles, follow some established interaction patterns and collaborate to reach global objectives [14] [19]. Nevertheless, the social context of agent societies (the social dependencies between agents and the effects of their membership to a group in the way that they can argue with other agents), is not analysed.

To our knowledge, no research is done to adapt argumentation frameworks to represent and manage arguments in agent societies taking into account their social context both in the representation of arguments and in the argument management process. Nevertheless, this social information plays an important role in the way agents can argue and learn from argumentation experiences. Depending on their social relations with other agents, an agent can accept arguments from a member of its society that it would never accept before acquiring social dependencies with this member. For instance, in a company a subordinate must sometimes accept arguments from his superior that go against his own ideas and that he would never accept without this power relation. Also, despite having no knowledge about an opponent agent, an agent could try to infer the potential willingness of the opponent to accept an argument by taking into account its social relation with the opponent, or even its social relation with similar agents in the past. These are major considerations that should be studied to apply argumentation techniques in real domains modelled by means of open MAS.

The purpose of this paper is twofold. On one hand, Section 2 analyses the requirements for an argumentation framework for agent societies and proposes a generic computational representation of arguments. This framework stresses the importance of the social dependencies between agents and the effects of their membership to a group in the way that they argue. On the other hand, in Section 3 we formalise this proposal by defining a computational argumentation framework (AF) for the design and implementation of argumentation dialogues in MAS. Our notion of argument relies on technological standards for argument and data interchange on the web. Hence, our argumentation framework can be adapted to work in multiple domains and distributed environments.

## 2 A Computational Model for Argument Representation in Agent Societies

In this section, we introduce the formal definition of the concepts that define our approach for agent societies. Then, we analyse the issues that have been considered to choose a suitable argumentation framework for agent societies. Taking

them into account, we propose a computational representation of arguments. Finally, an example is provided.

## 2.1 Society Model

In this work, we follow the approach of [10] and [2], who define an *agent society* in terms of a set of *agents* that play a set of *roles*, observe a set of *norms* and a set of *dependency relations* between roles and use a *communication language* to collaborate and reach the global objectives of the *group*. This definition is generic enough to fit most types of agent societies, such as social networks of agents or open agent organisations. Broadly speaking, it can be adapted to any open MAS where there are norms that regulate the behaviour of agents, roles that agents play, a common language that allow agents to interact defining a set of permitted locutions and a formal semantics for each of these elements. Moreover, the set of norms in open MAS define a *normative context* (covering both the set of norms defined by the system itself as well as the norms derived from agents' interactions) [8].

However, we consider that the values that individual agents or groups want to promote or demote and preference orders over them have also a crucial importance in the definition of an argumentation framework for agent societies. These values could explain the reasons that an agent has to give preference to certain beliefs, objectives, actions, etc. Also, dependency relations between roles could imply that an agent must change or violate its value preference order. For instance, an agent of higher hierarchy could impose their values to a subordinate or an agent could have to adopt a certain preference order over values to be accepted in a group. Therefore, we endorse the view of [21], [28] and [3], who stress the importance of the audience in determining whether an argument (e.g. for accepting or rejecting someone else's beliefs, objectives or action proposals) is persuasive or not. Thus, we have included in the above definition of agent society the notion of values and preference orders among them. Next, we provide a formal definition for the model of society that we have adopted:

**Definition 1 (Agent Society).** *An Agent society in a certain time  $t$  is defined as a tuple  $S_t = \langle Ag, Rl, D, G, N, V, Roles, Dependency, Group, val, Valpref_Q \rangle$  where:*

- $Ag = \{ag_1, ag_2, \dots, ag_I\}$  is the set of  $I$  agents of  $S_t$  in a certain time  $t$ .
- $Rl = \{rl_1, rl_2, \dots, rl_J\}$  is the set of  $J$  roles that have been defined in  $S_t$ .
- $D = \{d_1, d_2, \dots, d_K\}$  is the set of  $K$  possible dependency relations in  $S_t$ .
- $G = \{g_1, g_2, \dots, g_L\}$  is the set of groups that the agents of  $S_t$  form, where each  $g_l = \{a_1, a_2, \dots, a_M\}$ ,  $M \leq I$  consist of a set of agents  $a_i \in A$  of  $S_t$ .
- $N$  is the defined set of norms that affect the roles that agents play in  $S_t$ .
- $V = \{v_1, v_2, \dots, v_P\}$  is the set of  $P$  values predefined in  $S_t$ .
- $Roles : Ag \rightarrow 2^{Rl}$  is a function that assigns an agent its roles in  $S_t$ .
- $Dependency_{S_t} : \prec_D^{S_t} \subseteq Rl \times Rl$  defines a reflexive, transitive and asymmetric partial order relation over roles.
- $Group : Ag \rightarrow 2^G$  is a function that assigns an agent its groups in  $S_t$ .

- $val : Ag \rightarrow V$  is a function that assigns an agent its set of values.
- $Valpref_Q \subseteq V \times V$ , where  $Q = Ag \vee Q = G$ , defines a irreflexive, transitive and asymmetric preference relation  $<_Q^{S_t}$  over the values.

That is,  $\forall r_1, r_2, r_3 \in R, r_1 <_d^{S_t} r_2 <_d^{S_t} r_3$  implies that  $r_3$  has the highest rank with respect to the dependency relation  $d$  in  $S_t$ . Also,  $r_1 <_d^{S_t} r_2$  and  $r_2 <_d^{S_t} r_1$  implies that  $r_1$  and  $r_2$  have the same rank with respect to  $d$  in  $S_t$ . Finally,  $\forall v_1, v_2, v_3 \in V, Valpref_{ag_i} = v_1 <_{ag_i}^{S_t} v_2 <_{ag_i}^{S_t} v_3$  implies that agent  $ag_i$  prefers value  $v_3$  to  $v_2$  and value  $v_2$  to value  $v_1$  in  $S_t$ . Similarly,  $Valpref_{g_j} = v_1 <_{g_j}^{S_t} v_2 <_{g_j}^{S_t} v_3$  implies that group  $g_j$  prefers value  $v_3$  to  $v_2$  and value  $v_2$  to value  $v_1$  in  $S_t$ .

Once the concepts that we use to define agent societies are specified, the next section analyses the computational requirements for argument representation in these societies. Then, our approach for agent societies and the analysed requirements are used to propose a new computational representation for arguments.

## 2.2 Computational Requirements for Arguments in Agent Societies

An argumentation process is conceived as a reasoning model with several steps:

1. Building arguments (supporting or attacking conclusions) from knowledge bases.
2. Defining the strengths of those arguments by comparing them in conflict situations.
3. Evaluating the acceptability of arguments in view of the other arguments that are posed in the dialogue.
4. Defining the justified conclusions of the argumentation process.

The first step to design MAS whose agents are able to perform argumentation processes is to decide how agents represent arguments. According to the *interaction problem* defined in [7], “...representing knowledge for the purpose of solving some problem is strongly affected by the nature of the problem and the inference strategy to be applied to the problem...”. Therefore the way in which agents computationally represent arguments should ease the automatic performance of argumentation processes.

Most research effort on the computational representation of arguments is performed in the area of developing models for argument authoring and diagramming [23][27] (OVA<sup>1</sup>). However, these systems assume human users interacting with the software tool and are not conceived for performing agents’ automatic reasoning processes. Other research works where the computational modelling of arguments has been studied are those on case-based argumentation. From the first uses of argumentation in AI, arguments and cases are intertwined [29]. Case-based argumentation particularly reported successful applications in American common law [5], whose judicial standard orders that similar cases must be

<sup>1</sup> OVA at ARG:dundee: [www.arg.dundee.ac.uk](http://www.arg.dundee.ac.uk)



resolved with similar verdicts. In [4] a model of legal reasoning with cases is proposed. But, again, this model assumed human-computer interaction and cases were not thought to be only acceded by software agents. Case-Based Reasoning (CBR) systems [1] allow agents to learn from their experiences. In MAS, the research in case-based argumentation is quite recent with just a few proposals to date. These proposals are highly domain-specific or centralise the argumentation functionality in a *mediator* agent that manages the dialogue between the agents of the system [15].

As pointed out before, we focus on argumentation processes performed among a set of agents that belong to an agent society and must reach an agreement to solve a problem taking into account their social dependencies. Each agent builds its individual position in view of the problem (a solution for it). At this level of abstraction, we assume that this could be a generic problem of any type (e.g. a resource allocation problem, an agreed classification, a joint prediction, etc.) that could be characterised with a set of features. Thus, we assume that each agent has its individual knowledge resources to generate a potential solution. Also, agents have their own argumentation system to create arguments to support their positions and defeat the ones of other agents.

Taking into account the above issues, there are a set of requirements that a suitable framework to represent arguments in agent societies should met:

- be computationally tractable and designed to ease the performance of automatic reasoning processes over it.
- be rich enough to represent general and context dependent knowledge about the domain and social information about the agents' dependency relations or the agents' group.
- be generic enough to represent different types of arguments.
- comply with the technological standards of data and argument interchange on the web.

These requirements suggest that an argumentation framework for agent societies should be easily interpreted by machines and have highly expressive formal semantics to define complex concepts and relations over them. Thus, we propose a Knowledge-Intensive (KI) case-based argumentation framework [9], which allows automatic reasoning with semantic knowledge in addition to the syntactic properties of cases. Reasoning with cases is specially suitable when there is a weak (or even unknown) domain theory, but acquiring examples encountered in practice is easy. Most argumentation systems produce arguments by applying a set of inference rules. In open MAS the domain is highly dynamic and the set of rules that model it is difficult to specify in advance. However, tracking the arguments that agents put forward in argumentation processes could be relatively simple. Other important problem with rule-based systems arises when the knowledge-base of rules must be updated (e.g. adding a new rule). Updates imply to check the knowledge-base for conflicting or redundant rules. Case-based systems are in most cases easier to maintain than rule-based systems and hence, more suitable for being applied in dynamic domains.

In the following section, we present the framework proposed accordingly to the above requirements. This framework is also conceived to allow agents to improve their argumentation skills and be able to evaluate the persuasive power of arguments for specific audiences in view of their previous argumentation experiences.

### 2.3 Case-Based Model for Argument Representation

In open multi-agent argumentation systems the arguments that an agent generates to support its position can conflict with arguments of other agents and these conflicts are solved by means of argumentation dialogues between them. To allow agents to take the maximum profit from previous argumentation experiences, the structure that agents use to store information related to their argumentation experiences must be able to represent knowledge about individual arguments and also about the argumentation dialogues where arguments were posed. Therefore, agents that implement our argumentation framework have an individual case-based argumentation system with the following knowledge resources:

- Domain-cases: a set of cases that store information about problems that were solved in the past. The structure and concrete feature set of these cases depends on the specific application domain, but at least, they have a minimum set of features that represent the problem and the solution applied to it.
- Argument-cases: a set of cases that store information about arguments that the agent posed in the past and the results that were obtained by putting forward them in a previous argumentation dialogue<sup>2</sup>.
- Dialogue graphs: a set of directed graphs that link argument-cases and represent previous argumentation dialogues. Nodes represent arguments and arrows between nodes represent attack relations.
- Ontology of Argumentation Schemes: an ontology that encodes the set of argumentation schemes that agents can use to produce arguments. These schemes are stereotyped patterns of reasoning [32] that can be used to create presumptive arguments from a set of premises that characterise the problem to solve. In addition, argumentation schemes have a set of critical questions, which represent attacks to the conclusion drawn from the scheme.

The argument-cases are the main structure that we use to implement our framework and computationally represent arguments in agent societies. Argument-cases have two main objectives: 1) they can be used by agents as knowledge resource to generate new arguments in view of past argumentation experiences and 2) they can be used to store new argumentation knowledge that agents gain in each dialogue, improving the agents' argumentation skills. Due to space restrictions, we focus here on explaining this knowledge resource. Table 1 shows an

---

<sup>2</sup> Note that argument-cases and arguments are not the same, but the former are knowledge structures that store information about previous arguments (and maybe represent a generalisation of several arguments).

**Table 1.** Structure of an Argument Case

<b>PROBLEM</b>	Domain Context	Premises = {Volume, Price, etc.}	
	Social Context	Proponent	ID = F2
			Role = Farmer
			Norms = $N_{F2}$
			$ValPref_{F2} = [EC < SO]$
		Opponent	ID = BA
			Role = Basin Administrator
			Norms = $N_{BA}$
		Group	$ValPref_{BA} = \emptyset$
			ID = RB
		Norms = $N_{RB}$	
		$ValPref_{RB} = [SO < EC]$	
		Dependency Relation = Power	
<b>SOLUTION</b>	Argument Type = Inductive		
	Conclusion = $F2tr$ (F2 wins the water-right transfer)		
	Acceptability State = Acceptable		
	Received Attacks	Critical Questions = $\emptyset$	
		Distinguish Case = $\emptyset$	
	Counter Examples = {C1}		
<b>JUSTIFICATION</b>	Cases = {C2}		
	Schemes = $\emptyset$		
	Associated Dialogue Graph		

example of the structure of a specific argument-case (explained in the example of Section 2.4). As it is usual in CBR systems, the argument-cases have three main parts: the description of the *problem* that the case represents, the *solution* applied to this problem and the *justification* why this particular solution was applied. An argument-case stores the information about a previous argument that an agent posed in certain step of a dialogue with other agents.

**Problem:** The problem description stores the *premises* of the argument, which represent the context of the domain where the argument was put forward. In addition, if we want to store an argument and use it to generate a persuasive argument in the future, the features that characterise the audience of the previous argument (the social context) must also be kept.

For the definition of the social context of arguments, we follow our model of society presented in Section 2.1. Therefore, we store in the argument-case the social information about the *proponent* of the argument, the *opponent* to which the argument is addressed, the *group* to which both agents belong and the dependency relation established between the roles that these agents play. For the sake of simplicity, in what follows we assume that in each step of the dialogue, one proponent agent generates an argument and sends it to one opponent agent that belongs to its same group. However, either the proponent or the opponent's features could represent information about agents that act as representatives of a group and any agent can belong to different groups at the same time.

Thus, the proponent and opponent's features represent information about the agent that generated the argument and the agent that received it respectively.

Concretely, for each agent the argument-case stores a unique *ID* that identifies it in the system and the *role* that the agent was playing when it sent or received the argument (e.g. farmer and basin administrator, do not confuse with the role of proponent and opponent from the argumentation perspective). In addition, a reference to the set of norms that governed the behaviour of the agents at this step of the dialogue is also stored, since the normative context of agents could force or forbid them to accept certain facts and the arguments that support them (e.g. a norm could invalidate a dependency relation or a value preference order). Moreover, if known, we also store the preferences of each agent over the pre-defined set of general values in the system (e.g. security, solidarity, economy, etc.). As pointed out before, these preferences ( $ValPref_{F_2}$  and  $ValPref_{BA}$ ) affect the persuasive power of the proponent's argument over the opponent's behaviour.

Regarding the group features, the argument-case stores the unique identifier *ID* of the agents' group, the set of *norms* that regulates the behaviour of the group members at this moment, since changes can occur due to norm emergence, and the preference order ( $ValPref_{RB}$ ) about the *social values*<sup>3</sup> of the group. Finally, the dependency relation between the proponent's and the opponent's roles is also stored. To date, we define the possible dependency relations between roles as in [10]:

- *Power*: when an agent has to accept a request from other agent because of some pre-defined domination relationship between them (e.g. in a society  $S_t$  that manages the water of a river basin,  $Farmer <_{Power}^{S_t} Basin Administrator$ , since farmers must comply with the laws announced by the basin administrator).
- *Authorisation*: when an agent has committed itself to other agent for a certain service and a request from the latter leads to an obligation when the conditions are met (e.g. in the society  $S_t$ ,  $Farmer_i <_{Authorisation}^{S_t} Farmer_j$ , if  $Farmer_j$  has contracted a service that offers  $Farmer_i$ ).
- *Charity*: when an agent is willing to answer a request from other agent without being obliged to do so (e.g. in the society  $S_t$ , by default  $Farmer_i <_{Charity}^{S_t} Farmer_j$  and  $Farmer_j <_{Charity}^{S_t} Farmer_i$ ).

**Solution:** In the solution part, the *argument type* that defines the method by which the conclusion of the argument was drawn and this *conclusion* itself are stored. By default, we do not assume that agents have a pre-defined set of rules to infer deductive arguments from premises, which is difficult to maintain in open MAS. In our framework, agents have the following ways of generating new arguments:

- *Presumptive arguments*: by using the premises that describe the problem to solve and an argumentation scheme whose premises match them.

---

<sup>3</sup> We use the term social values to refer to those values that are agreed by (or commanded to) the members of a society as the common values that this society should promote (e.g. justice and solidarity in an ideal society) or demote.

- *Inductive arguments*: by using similar argument-cases and/or domain-cases stored in the case-bases of the system.
- *Mixed arguments*: by using premises, cases and argumentation schemes.

Moreover, the argument-case stores the information about the *acceptability state* of the argument at the end of the dialogue. This feature shows if the argument was deemed *acceptable*, *unacceptable* or *undecided* in view of the other arguments that were put forward during the dialogue (see Section 3 for details). Regardless of the final acceptability state of the argument, the argument-case also stores the information about the possible *attacks* that the argument received. These attacks could represent the justification for an argument to be deemed unacceptable or else reinforce the persuasive power of an argument that, despite being attacked, was finally accepted. Argument-cases can store different types of attacks, depending on the type of argument that they represent:

- For presumptive arguments: *critical questions* associated with the scheme.
- For inductive arguments [4]: either
  - Premises which value in the context where the argument was posed was different (or non-existent) than the value that it took in the cases used to generate the argument (*distinguish the case*) or
  - Cases which premises also match the premises of the context where the argument was posed, but which conclusion is different than the conclusion of the case(s) used to generate the argument (*counter-examples*).
- For mixed arguments: any of the above attacks.

**Justification:** The justification part of the argument-case stores the information about the knowledge resources that were used to generate the argument represented by the argument-case (e.g. the set argumentation schemes in presumptive arguments, the set of cases in inductive arguments and both in mixed arguments). In addition, each argument-case has associated a dialogue-graph that represents the dialogue where the argument was posed. This graph can be used later to develop dialogue strategies. The same dialogue graph can be associated with several argument-cases.

Following a CBR methodology, the knowledge resources of the agents' case-based argumentation system allow them to automatically generate, select and evaluate arguments. However, the complete argument management process (how agents generate, select and evaluate arguments by using the knowledge resources of their argumentation systems) is out of the scope of this paper. Also, the framework presented is flexible enough to represent different types of arguments and their associated information, but the value of some features on argument-cases and domain-cases could remain unspecified in specific domains. For instance, in some open MAS, the preferences over values of other agents could not be previously known. However, agents could try to infer the unknown features by using CBR adaptation techniques [17].

## 2.4 Example

To exemplify our framework, let us propose a simple scenario of an open MAS that represents a water market [6], where agents are users of a river basin that can buy or sell their water-rights to other agents. A water-right is a contract with the basin administration organism that specifies the rights that a user has over the water of the basin (e.g. the maximum volume that he can spend, the price that he must pay for the water or the district where it is settled<sup>4</sup>). In this setting, suppose that two agents that play the role of farmers, F1 and F2, are arguing with a basin administrator, BA, to decide over a water-right transfer agreement that will grant an offered water-right to a farmer. Then, the premises of the domain context would store data about the water-right transfer offer and other domain-dependent data about the current problem. All agents belong to the same group (the river basin RB) whose behaviour is controlled by certain set of norms  $N_{RB}$ , its value preference order promotes economy over solidarity (SO<EC) and commands a dependency relation of charity (C) between two farmers and power relation (P) between a basin administrator and a farmer. Also, F1 prefers economy over solidarity (SO<EC) and has a set of norms  $N_{F1}$ , F2 prefers solidarity over economy (EC<SO) and has a set of norms  $N_{F2}$  and by default, BA has the value preference order of the basin (SO<EC) and a set of norms  $N_{BA}$ .

Suppose that F1 has a domain-case C1 that represents a previous water-right transfer agreement that granted a similar water-right to a farmer whose land was adjacent to the district associated with the current water-right offer. Thus, F1 would put forward an argument to BA, generated by using C1.

A1: I should be the beneficiary of the transfer because my land is adjacent to the owner's land.

Here, we suppose that the closer the lands the cheaper the transfers between them and then, this argument would promote economy. However, F2 has a domain-case C2 that represents a previous water-right transfer agreement that granted a similar water-right to a farmer whose land needed an urgent irrigation to save the crop due to a drought. Thus, F2 would put forward the following argument to BA, generated by using C2.

A2: I should be the beneficiary of the transfer because there is a drought and my land is almost dry.

In this argument, we assume that crops are lost in dry lands and helping people to avoid losing crops promotes solidarity. In addition, suppose that as basin administrator, BA knows that there is a drought in the basin, which is a new premise that should be considered. Also, its ontology of argumentation schemes includes an *Argument for An Exceptional Case* scheme [32] S1 stating that the value preference order of the basin can be waived in case of drought and changed for EC<SO. Therefore, BA could generate an argument by using S1 and certain

---

<sup>4</sup> Following the Spanish Water Law, a water-right is always associated to a district.

domain-case C3 that granted a similar water-right transfer to a farmer whose land was dry in a drought to promote solidarity.

A3: There is a drought in the basin and dry lands must be irrigated first.

Table 1 shows the argument-case that F2 could store for A2 at the end of the dialogue, including the attacks received and the knowledge resources that support the argument. The dialogue graph of this argument-case would point to the node that represents it in the whole dialogue (represented with several argument-cases interlinked). Assuming that in our open MAS all agents can receive the arguments posed by the agents of their group, A1 and A2 will attack each other. In addition, A3 will attack A1, which do not takes into account the exceptional case of drought in the basin. Also, assuming that in this society S the administrator BA has a power dependency relation over any farmer ( $\text{Farmer} <_{\text{Power}}^S \text{Basin Administrator}$ ), F1 would have to accept the attack that defeats its argument A1 and withdraw it. If the dialogue ends here, the water-right transfer would be granted to F2.

Recall that argument-cases store the social information about roles, values, norms, etc. Therefore, agents can use this information when they are faced with the task of selecting a case from a set of possible cases to support their arguments. For instance, suppose that BA has also found a domain-case C4 that turned down a similar water-right transfer to a farmer whose land was dry in a drought. To decide which C3 or C4 is most suitable to draw a conclusion for the current problem, BA can check its arguments case-base. Then, suppose that BA finds the argument-case that represents the argument that ended the past dialogue that motivated the creation of the domain-case C4 by turning down the transfer. However, the social information about the group does not match with the current one. Thus, BA could infer that in those situation, the farmer was member of a different group where the irrigation of dry lands does not take priority in the case of drought and hence, C4 could not be cited in the current situation.

Up to this point, we have specified our approach for agent societies, analysed the requirements that a suitable argumentation framework for these type of societies should met and proposed our framework. In next section, we formalise this framework.

### 3 Case-Based Argumentation Framework for Agent Societies

Following our case-based computational representation of arguments, we have designed a formal AF as an instantiation of Dung's AF [11]. The main advantages that this framework contributes over other existent AFs are: 1) the ability to represent social information in arguments; 2) the possibility of automatically managing arguments in agent societies; 3) the improvement of the agents' argumentation skills; and 4) the easy interoperability with other frameworks that follow the argument and data interchange web standards. Next, the elements of the AF (according to Prakken's AF elements [22]) are specified.

### 3.1 The Notion of Argument: Case-Based Arguments

We have adopted the Argument Interchange Format (AIF) [33] view of arguments as a set of interlinked premiss-illative-conclusion sequences. The notion of argument is determined by our KI case-based framework to represent arguments. In our framework agents can generate arguments from previous cases (domain-cases and argument-cases), from argumentation schemes or from both. However, note that the fact that a proponent agent use one or several knowledge resources to generate an argument does not imply that it has to show all this information to its opponent. The argument-cases of the agents' argumentation systems and the structure of the actual arguments that are interchanged between agents is not the same. Thus, arguments that agents interchange are tuples of the form:

**Definition 2 (Argument).**  $Arg = \{\phi, \langle S \rangle\}$ , where  $\phi$  is the conclusion of the argument and  $\langle S \rangle$  is a set of elements that support it.

This support set can consist of different elements, depending on the argument purpose. On one hand, if the argument provides a potential solution for a problem (e.g. who should be the beneficiary of the transfer), the support set is the set of features (premises) that describe the problem to solve and optionally, any knowledge resource used by the proponent to generate the argument (domain-cases, argument-cases, argumentation schemes or elements of them). On the other hand, if the argument attacks the argument of an opponent, the support set can also include any of the allowed attacks in our framework (critical questions, distinguishing premises or counter-examples). Then, the support set consists of the following tuple of sets of support elements<sup>5</sup>:

**Definition 3 (Support Set).**  $S = \langle \{Premises\}, \{DomainCases\}, \{ArgumentCases\}, \{ArgumentationSchemes\}, \{CriticalQuestions\}, \{DistinguishingPremises\}, \{CounterExamples\} \rangle$

For instance, assuming that  $\sim$  stands for the logical negation and the set of  $n$  premises is defined as  $Pre = \{pre_1, \dots, pre_n\}$ , in our example we have that:

$$\begin{aligned} A1 &= \{F1tr, \langle Pre, \{C1\}, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle\} \\ A2 &= \{F2tr, \langle Pre, \{C2\}, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle\} \\ A3 &= \{\sim C1, \langle Pre \cup \{Drought\}, \emptyset, \emptyset, \{S1\}, \emptyset, \{Drought\}, \{C3\} \rangle\} \end{aligned}$$

where  $F1tr$  and  $F2tr$  mean that the transfer should be granted to the farmers F1 or F2 respectively and  $\sim C1$  means that this case cannot be applied in this context, due to the new distinguishing premise  $\{Drought\}$  and the counter-example  $C3$ .

### 3.2 The Logical Language

The logical language represents argumentation concepts and possible relations among them. In our framework, these concepts are represented in the form of KI

<sup>5</sup> This representation is only used for illustrative purposes and efficiency considerations about the implementation are obviated.



cases and argumentation schemes. Therefore, the logical language of the AF is defined in terms of the vocabulary to represent these resources. In this section, we focus on the definition of the logical language to represent cases. To represent schemes, we use the AIF ontology proposed in [25].

The vocabulary of cases is defined by using an ontology inspired by the approach proposed in [9] and the AIF ontology. We have selected the Ontology Web Language OWL-DL [6] as the formal logics to represent the vocabulary of cases. This variant is based on Description Logics (DL) and guarantees computational completeness and decidability. Thus, it allows for automatic description logic reasoning over argument-cases and domain-cases. In addition, it facilitates the interoperability with other systems. Next, we provide a partial view of the top levels of the ontology [7] for the AF proposed.

In the top level of abstraction, the terminological part of the ontology distinguishes between three disjoint concepts: *Case*, which is the basic structure to store the argumentation knowledge of agents; *CaseComp*, which represent the usual parts that cases have in CBR systems; and *CaseAtt*, which are the specific attributes that make up each component:

$$\begin{aligned} Case &\sqsubseteq Thing & Case &\sqsubseteq \neg CaseComp \\ CaseComp &\sqsubseteq Thing & CaseComp &\sqsubseteq \neg CaseAtt \\ CaseAtt &\sqsubseteq Thing & CaseAtt &\sqsubseteq \neg Case \end{aligned}$$

As pointed out before, there are two disjoint types of cases:

$$\begin{aligned} ArgumentCase &\sqsubseteq Case & DomainCase &\sqsubseteq Case \\ ArgumentCase &\sqsubseteq \neg DomainCase \end{aligned}$$

Both argument-cases and domain-cases have the three possible types of components that usual cases of CBR systems have: the description of the state of the world when the case was stored (*Problem*); the solution of the case (*Conclusion*); and the explanation of the process that gave rise to this conclusion (*Justification*):

$$\begin{aligned} Problem &\sqsubseteq CaseComp & Conclusion &\sqsubseteq CaseComp \\ Justification &\sqsubseteq CaseComp \\ Case &\sqsubseteq \forall hasProblem.Problem \\ Case &\sqsubseteq \forall hasConclusion.Conclusion \\ Case &\sqsubseteq \forall hasJustification.Justification \end{aligned}$$

Case components are composed of one or more attributes:

$$CaseComp \sqsubseteq \geq 1 hasAttribute.CaseAtt$$

For instance, the attributes of the solution description of an argument-case are presented below. The cardinality of the possible attacks that an argument-case can receive is not specified, since the case could not have been attacked.

$$\begin{aligned} ArgumentType &\sqsubseteq CaseAtt \\ Solution &\sqsubseteq = 1 hasArgumentType.ArgumentType \\ Conclusion &\sqsubseteq CaseAtt \\ Solution &\sqsubseteq = 1 hasConclusion.Conclusion \\ AcceptabilityState &\sqsubseteq CaseAtt \\ Solution &\sqsubseteq = 1 hasAcceptabilityState.AcceptabilityState \\ ReceivedAttacks &\sqsubseteq CaseAtt \end{aligned}$$

<sup>6</sup> <http://www.w3.org/TR/owl-guide/>

<sup>7</sup> The complete specification of the ontology is out of the scope of this paper.

*CriticalQuestions*  $\sqsubseteq$  *ReceivedAttacks*  
*DistinguishingPremises*  $\sqsubseteq$  *ReceivedAttacks*  
*CounterExamples*  $\sqsubseteq$  *ReceivedAttacks*

In addition, some additional properties about the concepts of the ontology can also be defined. For instance, instances of argument-cases can have a unique identifier, a creation date or a date for the last time that the case was used, which could be used to determine if a case is outdated and should be removed from the case-base<sup>8</sup>. For simplicity, these elements are not shown in Table 11.

*Case*  $\sqsubseteq$  *1identifier*  
 $T \sqsubseteq \forall identifier.ID \quad T \sqsubseteq \forall identifier^-.Case$   
*Case*  $\sqsubseteq$  *1creationDate*  
 $T \sqsubseteq \forall creationDate.Date \quad T \sqsubseteq \forall creationDate^-.Case$   
*Case*  $\sqsubseteq$  *1lastUsed*  
 $T \sqsubseteq \forall lastUsed.Date \quad T \sqsubseteq \forall lastUsed^-.Case$

### 3.3 The Concept of Conflict between Arguments

The concept of conflict between arguments defines in which way arguments can attack each other. There are two typical attacks studied in argumentation: *rebut* and *undercut*. In an abstract definition, rebuttals occur when two arguments have contradictory conclusions. Similarly, an argument undercuts other argument if its conclusion is inconsistent with one of the elements of the support set of the latter argument or its associated conclusion. This section shows how our AF instantiates these two attacks. Taking into account the possible elements of the support set, rebut and undercut attacks can be formally defined as follows. Let  $Arg_1 = \{\phi_1, \langle S_1 \rangle\}$  and  $Arg_2 = \{\phi_2, \langle S_2 \rangle\}$  be two different arguments, where  $S_1 = \langle \{Premises\}_1, \dots, \{CounterExamples\}_1 \rangle$ ,  $S_2 = \langle \{Premises\}_2, \dots, \{CounterExamples\}_2 \rangle$ ,  $\sim$  stands for the logical negation,  $\Rightarrow$  stands for the logical implication and  $conc(x)$  is a function that returns the conclusion of the formula  $x$ . Then:

**Definition 4 (Rebut).**  $Arg_1$  rebuts  $Arg_2$  iff  $\phi_1 = \sim\phi_2$  and  $\{Premises\}_1 \supseteq \{Premises\}_2$

That is, if  $Arg_1$  supports a different conclusions for a problem description that includes the problem description of  $Arg_2$ . Assuming  $F1tr = \sim F2tr$  and vice-versa, in our example,  $A1$  and  $A2$  rebut each other.

**Definition 5 (Undercut).**  $Arg_1$  undercuts  $Arg_2$  if

- 1)  $\phi_1 = \sim conc(as_k) /$   
 $\exists cq \in \{CriticalQuestions\}_1 \wedge \exists as_k \in \{ArgumentationSchemes\}_2 \wedge$   
 $cq \Rightarrow \sim conc(as_k),$  or
- 2)  $\phi_1 = dp /$   
 $(\exists dp \in \{DistinguishingPremises\}_1 \wedge \exists pre_k \in \{Premises\}_2 \wedge dp = \sim pre_k) \vee$   
 $(dp \notin \{Premises\}_2),$  or
- 3)  $\phi_1 = ce /$

<sup>8</sup> In DL, the range of a property  $C$  is specified as  $T \sqsubseteq \forall R.C$  and its domain as  $T \sqsubseteq \forall R^-.C$ .

$$\begin{aligned}
 & (\exists ce \in \{CounterExamples\}_1 \wedge \exists dc_k \in \{DomainCases\}_2 \\
 & \wedge conc(ce) = \sim conc(dc_k)) \vee \\
 & (\exists ce \in \{CounterExamples\}_1 \wedge \\
 & \exists ac_k \in \{ArgumentCases\}_2 \wedge conc(ce) = \sim conc(ac_k))
 \end{aligned}$$

That is, if the conclusion drawn from  $Arg_1$  makes one of the elements of the support set of  $Arg_2$  or its conclusion non-applicable in the current context of the argumentation dialogue. In our example,  $A3$  undercuts  $A1$ , since its conclusion makes  $C1$  non-applicable due to the counter-example  $C3$  and the distinguishing premise  $\{Drought\}$ , which is not considered in the premises that describe the previous problem that is represented by  $C1$  and made  $F1$  to infer  $A1$  from it.

### 3.4 The Notion of Defeat between Arguments

Once possible conflicts between argument have been defined, the next step in the formal specification of an AF is to define the defeat relation between a pair of arguments. This comparison must not be misunderstood as a strategical function to determine with which argument an argumentation dialogue can be won [22]. A function like this must also consider other factors, such as other arguments put forward in the dialogue or agents' profiles. Therefore, it only tells us something about the relation between two arguments. Hence, the relation of defeat between two arguments is defined in our AF as follows. Let  $Arg_1 = \{\phi_1, \langle S_1 \rangle\}$  and  $Arg_2 = \{\phi_2, \langle S_2 \rangle\}$  be two conflicting arguments. Then:

**Definition 6 (Defeat).**  $Arg_1$  defeats  $Arg_2$  if  $Arg_1$  rebuts  $Arg_2$  and  $Arg_2$  does not undercut  $Arg_1$ , or else  $Arg_1$  undercuts  $Arg_2$

The first type of defeat poses a stronger attack on an argument, directly attacking its conclusion. In addition, an argument can strictly defeat other argument.

**Definition 7 (Strict Defeat).**  $Arg_1$  strictly defeats  $Arg_2$  if  $Arg_1$  defeats  $Arg_2$  and  $Arg_2$  does not defeat  $Arg_1$

In our example,  $A1$  and  $A2$  defeat each other and  $A3$  strictly defeats  $A1$ .

### 3.5 The Acceptability State of Arguments

The acceptability state of arguments determines their status on the basis of their interaction. Only comparing pairs of arguments is not enough to decide if their conclusions are acceptable, since defeating arguments can also be defeated by other arguments. Taking into account the underlying domain theory of a dialectical system, arguments can be considered *acceptable*, *unacceptable* and *undecided* [11]. However, the acquisition of new information in further steps of the dialogue could change the acceptability state of arguments.

Therefore, to decide the acceptability state of arguments a proof theory that takes into account the dialogical nature of the argumentation process is necessary. To evaluate the acceptability of arguments by using a dialogue game is a common approach. Dialogue games are interactions between two or more players, where each one moves by posing statements in accordance with a set or

predefined rules [18]. In our AF, the acceptability state of arguments could be decided by using a dialogue game and storing in the argument-case associated to each argument its acceptability state when the dialogue ends. However, the definition of this game is out of the scope of this paper.

## 4 Discussion

In this paper, we have presented a computational framework to represent arguments in agent societies. This framework takes into account the social dependencies between agents and the effects of their membership to a group in the way that they can argue. However, although the framework is flexible enough to store complex knowledge about arguments and dialogues, the value of some case features could not be specified or known in some domains. For instance, the proponent of an argument obviously knows its own preferences over its set of values, probably knows the preferences of its group but, in a real open MAS, we cannot assume that it also knows the value preferences of its opponent. However, the proponent can know the value preferences of the opponent's group (if both belong to the same) or have some previous knowledge about the value preferences of similar agents playing the same role that the opponent is playing now. The same could happen when agents belong to different groups. Thus, the group features could be unknown, but the proponent could try to use its experience with other agents of the opponent's group and infer these features.

In addition, the argumentation framework was inspired by the standard for argument interchange on the web and hence, an argumentation system based on it can interact with other systems that comply with the standard. Elements of cases are specified by using an ontologic case representation language. This means that agents that implement our case-based framework for argument representation and management could argue with agents with other models of reasoning. Each element of the knowledge structures of the argumentation framework proposed can be translated to a concept of the AIF ontology [23] or an ontology for CBR systems based on [9]. For instance, domain premises can be translated into AIF *Premise Descriptions Forms* and premise values into *Premise I-Nodes*, value preferences can instantiate *Preference-Application-Nodes S-Nodes* and argument types *Presumptive Rule-of-Inference Schemes*. Even temporal propositions, agents, roles and norms can be described with OWL ontologies, as proposed in [13]. Although agents in open MAS are heterogeneous, by sharing these ontologies they can *understand* the arguments interchanged in the system.

Moreover, a formal argumentation framework has been presented. This framework is aimed at providing agents with the ability of having argumentation dialogues with other agents in agent societies, with a weak or unknown domain theory. Moreover, the KI case-based approach used for representing argumentation related information allows agents to apply CBR techniques to learn from the experience and improve their argumentation skills. Current work is focused on the development of the necessary CBR algorithms to generate, select and evaluate arguments from domain-cases, argument-cases and argumentation schemes.

## Acknowledgment

This work is supported by the Spanish government grants CONSOLIDER INGENIO 2010 CSD2007-00022, TIN2008-04446 and TIN2009-13839-C03-01 and by the GVA project PROMETEO 2008/051.

## References

1. Aamodt, A., Plaza, E.: Case-based reasoning: foundational issues, methodological variations and system approaches. *AI Communications* 7, no. 1, 39–59 (1994)
2. Artikis, A., Sergot, M., Pitt, J.: Specifying norm-governed computational societies. *ACM Transactions on Computational Logic* 10(1) (2009)
3. Bench-Capon, T., Atkinson, K.: *Argumentation in Artificial Intelligence*, chap. Abstract argumentation and values, pp. 45–64 (2009)
4. Bench-Capon, T., Sartor, G.: A model of legal reasoning with cases incorporating theories and values. *Artificial Intelligence* 150(1-2), 97–143 (2003)
5. Bench-Capon, T., Dunne, P.: *Argumentation in artificial intelligence*. *Artificial Intelligence* 171(10-15), 619–938 (2007)
6. Botti, V., Garrido, A., Giret, A., Noriega, P.: Managing water demand as a regulated open mas. In: *Workshop on Coordination, Organization, Institutions and Norms in agent systems in on-line communities, COIN-09*. vol. 494, pp. 1–10 (2009)
7. Bylander, T.C., Chandrasekaran, B.: Generic tasks in knowledge-based reasoning: The right level of abstraction for knowledge acquisition. *International Journal of Man-Machine Studies* 26(2), 231–243 (1987)
8. Criado, N., Argente, E., Botti, V.: A Normative Model For Open Agent Organizations. In: *International Conference on Artificial Intelligence, ICAI-09* (2009)
9. Diaz-Agudo, B., Gonzalez-Calero, P.A.: *Ontologies: A Handbook of Principles, Concepts and Applications in Information Systems*, chap. An Ontological Approach to Develop Knowledge Intensive CBR Systems, pp. 173–214 (2007)
10. Dignum, V.: PhD Dissertation: A model for organizational interaction: based on agents, founded in logic. Ph.D. thesis (2003)
11. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming, and n -person games. *Artificial Intelligence* 77, 321–357 (1995)
12. Ferber, J., Gutknecht, O., Michel, F.: From Agents to Organizations: an Organizational View of Multi-Agent Systems. In: *Agent-Oriented Software Engineering VI*. LNCS, vol. 2935, pp. 214–230. Springer-Verlag (2004)
13. Fornara, N., Colombetti, M.: Ontology and Time Evolution of Obligations and Prohibitions using Semantic Web Technology. In: *Workshop on Declarative Agent Languages and Technologies, DALT-09* (2009)
14. Gaertner, D., Rodriguez, J.A., Toni, F.: Agreeing on institutional goals for multi-agent societies. In: *5th International Workshop on Coordination, Organizations, Institutions, and Norms in agent systems, COIN-08* (2008)
15. Heras, S., Botti, V., Julian, V.: Challenges for a CBR framework for argumentation in open MAS. *Knowledge Engineering Review* 24(4), 327–352 (2009)
16. Luck, M., McBurney, P.: Computing as interaction: agent and agreement technologies. In: *IEEE International Conference on Distributed Human-Machine Systems* (2008)

17. López de Mántaras, R., McSherry, D., Bridge, D., Leake, D., Smyth, B., Craw, S., Faltings, B., Maher, M.L., Cox, M., Forbus, K., Keane, M., Watson, I.: Retrieval, Reuse, Revision, and Retention in CBR. *The Knowledge Engineering Review* 20(3), 215–240 (2006)
18. McBurney, P., Parsons, S.: Dialogue games in multi-agent systems. *Informal Logic. Special Issue on Applications of Argumentation in Computer Science* 22(3), 257–274 (2002)
19. Oliva, E., McBurney, P., Omicini, A.: Co-argumentation artifact for agent societies. In: 5th International Workshop on Argumentation in Multi-Agent Systems, ArgMAS-08 (2008)
20. Ossowski, S., Julian, V., Bajo, J., Billhardt, H., Botti, V., Corchado, J.M.: Open issues in open mas: An abstract architecture proposal. vol. 2, pp. 151–160 (2007)
21. Perelman, C., Olbrechts-Tyteca, L.: *The New Rhetoric: A Treatise on Argumentation* (1969)
22. Prakken, H., Sartor, G.: A dialectical model of assessing conflicting arguments in legal reasoning. *Artificial Intelligence and Law* 4, 331–368 (1996)
23. Rahwan, I., Zablith, F., Reed, C.: Laying the foundations for a world wide argument web. *Artificial Intelligence* 171(10-15), 897–921 (2007)
24. Rahwan, I.: Argumentation in multi-agent systems. *Autonomous Agents and Multiagent Systems*, Guest Editorial 11(2), 115–125 (2006)
25. Rahwan, I., Banihashemi, B.: Arguments in OWL: a progress report. pp. 297–310 (2008)
26. Reed, C., Grasso, F.: Recent advances in computational models of natural argument. *International Journal of Intelligent Systems* 22, 1–15 (2007)
27. Rowe, G., Reed, C.: Diagramming the argument interchange format. In: Conference on Computational Models of Argument, COMMA-08. pp. 348–359 (2008)
28. Searle, J.: *Rationality in Action* (2001)
29. Skalak, D., Rissland, E.: Arguments and cases: An inevitable intertwining. *Artificial Intelligence and Law* 1(1), 3–44 (1992)
30. Soh, L.K., Tsatsoulis, C.: A real-time negotiation model and a multi-agent sensor network implementation. *Autonomous Agents and Multi-Agent Systems* 11(3), 215–271 (2005)
31. Tolchinsky, P., Cortés, U., Modgil, S., Caballero, F., López-Navidad, A.: Increasing human-organ transplant availability: Argumentation-based agent deliberation. *IEEE Intelligent Systems* 21(6), 30–37 (2006)
32. Walton, D., Reed, C., Macagno, F.: *Argumentation Schemes* (2008)
33. Willmott, S., Vreeswijk, G., Chesñevar, C., South, M., McGinnis, J., Modgil, S., Rahwan, I., Reed, C., Simari, G.: Towards an argument interchange format for Multi-Agent Systems. In: 3rd International Workshop on Argumentation in Multi-Agent Systems, ArgMAS-06. pp. 17–34 (2006)

# Towards a Dialectical Approach for Conversational Agents in Selling Situations

Maxime Morge, Sameh Abdel-Naby, and Bruno Beaufile

Laboratoire d'Informatique Fondamentale de Lille  
Université Lille 1  
Bât M3 - F-59655 Villeneuve d'Ascq  
{maxime.morge,sameh.abdel-naby,bruno.beaufils}@lifl.fr

**Abstract.** The use of virtual agents to intelligently interface with online customers of e-commerce businesses is remarkably increasing. Most of these virtual agents are designed to assist online customers while searching for information related to a specific product or service, while few agents are intended for promoting and selling a product or a service. Within the later type, our aim is to provide proactive agents that recommend a specific item and justify this recommendation to a customer based on his purchases history and his needs. In this paper, we propose a dialectical argumentation approach that would allow virtual agents that have sales goals to trigger persuasions with e-commerce's customers. Then, we illustrate the proposed idea through its integration with an example from real-life.

## Categories and Subject Descriptors

I.2.11 [Artificial Intelligence:] Intelligent Agents.

**General Terms:** Algorithms.

**Keywords:** Argumentation, E-commerce, Agents, Language Processors.

## 1 Introduction

Within the last twelve years, precisely from 1998 wherein the dot-coms' boom first made an impact, e-commerce has succeeded to pursue a massive number of shoppers to change their idea of buying [1]. Several existing businesses have taken an advantage of this boom by adding a virtual presence to their physical one by means of an e-commerce website, these companies are now called *brick and mortar* businesses (e.g., Barnes & Noble). Additionally, new companies that exist only through the web, called *bricks and clicks* businesses, have also appeared (e.g., Amazon). Although the online presence of companies is cost-efficient, yet the lack of a persuading salesman affects the transformation ratio (sales vs. visits).

Apart from the *Business's* reaction to the boom, in Computer Science, several research efforts were made to study, analyze, and better shape the processes of assisting customers while being present in an e-commerce space [2,3]. In Artificial

Intelligence, a considerable amount of the research conducted in the area of Software Agents [4] focus on the enhancement and the proper provision of online Embodied Conversational Agents (ECAs) [5].

Whether these agents sell, assist, or just recommend, it is now clear that such autonomous agents are capable of engaging in verbal and non-verbal dialogues with e-commerce's customers. However, the ability of these agents to transform an ordinary visitor of an e-commerce who needs assistance to an actual buyer is yet of no notable weight. For an overview of the issues encountering the development of virtual sales agents refer to [6].

Since most of the currently available ECAs for e-commerce are designed to ask questions and wait for answers, one of the major challenges of concerned scholars is related to the reversibility of the current dialogue schemas. Meaning, to reach a proper ECA proactivity / sales attitude, the questions an agent should ask to collect sales data should be placed in nowadays agents' answers, and the vice-versa. Consequently, for the ECAs of existing literature; the current design approach of agents' answers generation mechanism must be adjusted for a conversational agent who is in the process of asking questions too (proactive) and not just giving answers.

In this paper, we propose the use of dialectical argumentation technologies as a step on the way to increase the sales-oriented negotiation skills of software agents in the business-to-consumer (B2C) segment of e-commerce. For this purpose, we suggest the exploitation of existing argumentation tools, such as those found in [7,8,9]. Using these tools we intend to build a sales-driven dialogue system that is capable of leading a virtual seller agent to influence the decision of a potential buyer in an e-commerce setting. Then, we illustrate the proposed idea through its integration with an example from real-life.

This paper is organized as follows. In section [2] we give an overview of the existing dialogue systems while pointing out their limitations. In section [3] we adopt a different approach for dialogue management based upon argumentation. Section [4] illustrates this approach using an intuitive scenario. Section [5] briefly describes the CSO language processor on which our dialogue system is based. The rest of the paper overviews the dialectical argumentation technology we consider. Section [6] outlines the dialogue-game protocol we use. Section [7] presents our realization of the dialogue strategy. We then conclude this paper by discussing some of the related work and, providing a summary of our future work.

## 2 Dialogue Systems

A dialogue system is a computer system that is capable of interacting with humans using the language they understand - *natural language*. Similar to that we can find TRAINS-93 [10], Collagen [11] and Artemis Agent Technology [12], which are mixed-initiative dialogue systems for collaborative problem solving. These dialogue systems can respond to initiatives made by users and, they also take initiatives themselves, which is required to support a selling process.

TRAINS-93 [10], Collagen [11] and Artemis Agent Technology [12] are adopting the same approach of focusing on the dialogue modelling itself besides the



dialogue management that is based on intentions recognition. For example, out of the following utterance of a user, "I want to purchase a quilt", there can be three possible interpretations:

1. It can be a direct report of a need;
2. It can be a statement of a goal that a user is pursuing independently;
3. It can be a proposal to adopt this joint goal.

Particularly, the discourse structure considered by Collagen in [11] is based on a comprehensive axiomatization of SharedPlans [13], while TRAINS-93 and Artemis Agent Technology are based upon a BDI approach [14]. The semantics of utterances is specified with the help of a first order modal logic language using operators as Beliefs, Desires and Intentions. The notions of persistent goal is a composite mental attitude which is defined from the previous operators in order to formalize the intention expressed by utterances. According to the semantic language of FIPA-ACL [15] adopted by the Artemis Agent Technology, an agent  $i$  has  $p$  as a persistent goal, if  $i$  has  $p$  as a goal and is self-committed toward this goal until  $i$  comes to believe that the goal is achieved or, this goal is unachievable. Here, an intention is defined as a persistent goal imposing the agent to act, which accordingly generates a planning process.

The process of inferring intentions from actions is needed to constraint and reduce the amount of communications exchanged. Also, it is worth noticing here that it is hard to incorporate this process into practical computer systems due to the complexities encountered while facilitating natural intractability. Therefore, it is then required to develop a heuristic mechanism for software agents in a collaborative setting.

For this purpose, dialogue systems are required to recognize the intention of the user and reason about it. The implementation of this theory is problematic due to its computational complexity [16]. Moreover, the specification of the semantics for the speech acts in terms of mental states is not adapted for resolving the conflicts which can appear during a selling process. For instance, an information that is received by a virtual seller agent must be adopted even if this information is contradictory with its beliefs. Those are the reasons why we consider an alternative approach based upon dialectical argumentation.

### 3 Dialectical Approach

Our approach for dialogue modelling considers the exchange of utterances as an argumentation process regulated by some normative rules that we call *dialogue-game protocol*. Our approach is inspired by the notion of dialectical system that Charles L. Hamblin introduced in [17]. A **dialectical system** is a family of regulated dialogue, (i.e., a system through which a set of participants communicate in accordance with some rules).

From this perspective, Walton and Krabbe in [18] define a **dialogue** as a coherent and structured sequence of utterances aiming at moving from an initial state to reach the goals of the participants. These are the dialogue's *goals* that

**Table 1.** Systemic overview of dialogue categories

Initial situation → Goal ↓	Conflict	Open problem	Ignorance of a participant
Stable agreement i.e., Resolution	persuasion	enquiry	information seeking
Practical settlement i.e., Decision	negotiation	deliberation	∅

can be shared by the participants or they can be also each of the participants' individual goals. Based on this definition, Walton and Krabbe have distinguished between five main categories of dialogues depending on the initial situation and goals. These categories are: information seeking, persuasion, negotiation, enquiry and deliberation [18].

Table 1 represents the analysis grid for dialogues proposed by Walton and Krabbe. An **information seeking** appears when a participant aims at catching knowledge from its interlocutor. The goal is to spread knowledge. In a **persuasion** dialogue, the initial situation is disagreement, (i.e., a conflict of opinion). The goal consists of solving the conflict by verbal means. In a **negotiation** dialogue, the initial situation is a conflict of interest mixed with a need for collaboration. The goal consists of a deal, i.e. an agreement attracting all participants to maximizing their gains. An **enquiry** dialogue aims at establishing (or demonstrating) the truth of a predicate. This one must answer to an open question and a stable agreement emerges. Each participant aims at extending their knowledge. A **deliberation**, as an enquiry, begins with an open problem rather than a conflict. The discussion is about the means and ends of a future action. It is worth noticing that, in real world, the nature of dialogues can be mixed. A dialogue can be composed of different sub-dialogues with different natures as we will see in our scenario.

## 4 Dialogue: Phases and Purposes

In this section, we explain the different phases of the overall online sales process that we are attempting to tackle in our research. Within these phases, we expect our virtual agent to rely on a specific language processor - *explained further ahead* - to handle online one-to-one conversations, related misspelling, and the use of diverse languages. Since the existing language processor is already capable of handling what is known to us as After-Sales, (i.e., assisting online users while searching for problems' answers), we then became extra interested to increase the salesability of this agent.

- **BEFORE-SALE:** in this phase we distinguish between two different processes that are possibly interleaved: a) the process of needs identification and, b) the process of product selection.

The **Needs Identification** can be performed with the help of an information seeking dialogue shifting from an initial asymmetric situation to a final one where both of the players share the user requirements.

The **Product Selection** allows the participants to constraint and to reduce the amount of communication by considering only relevant products later in the selling process. This task, in overall, also supports the information seeking dialogue where the virtual seller agent asks discriminatory questions in order to narrow its focus into a single product.

Both of these dialogues can be interleaved. The aim of the virtual seller agent here is to spread information about the products, while the aim of the user is expected to be the spreading of information about his needs.

- **SALE:** here, the aim for all dialogues' parties is to bargain over their interests and, eventually, "make a deal". For this purpose, the participants play a role in a **negotiation** dialogue. The simplest dialogue is: the virtual seller agent makes an offer and the user accepts or refuses this proposal.

If there is no single product corresponding to the user needs, then the participants attempt to maximize their benefits by conceding some aspects while insisting on others. If no product is matching user's needs, the user's high-ranked features of the products, (i.e., top priority conditions), are then altered to a lower ranked ones. On the other hand, if more than one product is corresponding to user's needs, the virtual seller agent picks the ones with the highest gross margin.

Later to that, the virtual seller agent can suggest - *after a single sale* - additional sub-items or an offer which invokes more of the user needs. These later approaches are well-known marketing techniques, (i.e., cross-selling and up-selling). In both cases, the new items aim at persuading the customer.

- **AFTER-SALE:** The initial condition of this phase is not coming out of a sales conflict, however, it comes out of the user's unawareness of the problem related to the concerned product. Therefore, participants here aim at extending their knowledge about each foreseen encounter and, eventually, decide about their subsequent action, (e.g., product replacement). As a result, users here play in an enquiry dialogue or a deliberation fashion.

The dialogue in Fig 1 consists of a specific case run which embraces the before-sale and sale phases between a user and a virtual seller agent (VSA). While the utterances #3 and #5 are related to the *Needs Identification*, the utterances #1,#6 and #8 are of the *Product Selection*. The rest of utterances are directly related to the sales process. The utterance #9 is a single-sale offer, the utterance #11 captures an up-selling offer and the utterance #13 is a cross-selling proposal.

It is also worth highlighting here that relying on argumentation will make the virtual seller agent more credible and, then, persuasive, since its appearance or involvement will yield from the fulfillment of a set of predefined conditions, (e.g., user X is wandering in the quilt section).

The After-Sale phase, (i.e., assistance), is out of the scope of this paper since the existing query answering systems, such as the CSO Language Processor we

1. **VSA:** Your moves through our website indicates that you might be interested in Quilts, do you need help to select one?
2. **User:** Nice.
3. **VSA:** Are you allergic?
4. **User:** Yes.
5. **VSA:** Do you sweat?
6. **User:** Sure. Is there any organic quilts ?
7. **VSA:** Definitely, but which dimensions are you looking for?
8. **User:** 200x200. Is there any one within a budget of €350?
9. **VSA:** Yes, it cost €239 since the dimensions are 200x200, it is nonallergenic and organic.
10. **User:** I buy it.
11. **VSA:** Since it is still within your budget, I may also suggest you to buy a “four seasons” one that cost €279?
12. **User:** Good idea.
13. **VSA:** What about a quilt cover that cost €99 only? still in your budget.
14. **User:** No, thanks.

**Fig. 1.** Virtual Seller Agent Vs. User: Before-Sale & Sale scenarios

present in the next section, have been practically proven to be efficient for such tasks.

## 5 CSO Language Processor

The CSO Language Processor [9] is the technology upon which our dialogue system is built.

The CSO Language Processor provided by Artificial Solutions<sup>1</sup> allows to produce virtual dialogue assistants such as Anna<sup>2</sup> and numerous other systems in more than 20 different languages. Elbot, which has won the Loebner Prize<sup>3</sup> in 2008, is built upon this technology.

This technology manages users’ sessions, handles misspellings and it also contains a language dependent preprocessing feature. In accordance with the dialogue state, it selects and carries out the best dialogue move. Additionally, this technology is able to interact with a back-end system, (e.g., databases), to hand out answer document for requesting application/front end and to write log files for analysis.

The inputs of the language processor are the user queries, (i.e., the user’s identity and his text inputs). After the identification of the session, the inputs are divided in sentences and words and the spelling is corrected. Another phase is carried out wherein an interpretation of the inputs is made: an answer retrieval for each sentences of the user’s inputs based on some **interaction rules in a knowledge base**. Finally, the answer is selected and generated by replacing some template variables.

<sup>1</sup> <http://www.artificial-solutions.com>

<sup>2</sup> <http://www.ikea.com>

<sup>3</sup> <http://www.loebner.net/Prizef/loebner-prize.html>

The interaction rules combine the meaning of the user's inputs and the dialogue state to define the conditions under which a dialogue move may be uttered. A given move can only be performed if the conditions are completely fulfilled.

However, the core of the CSO Language Processor is an inference engine that implements forward-chaining and so reactionary. Therefore, the language processor only makes it possible to respond to a user's queries and not to initiate or lead a sales-driven conversations. consequently, in order for us to make CSO proactive and enable it to go through sales-driven encounters, we introduce in the next section a formal framework for possible sale-driven dialogue management that can be adapted by virtual agents.

## 6 Dialectical System

A dialogue is a social interaction amongst self-interested parties intended to reach a common goal. In this section, we present how our game-based social model [8] handles the forseen conversation between a user and a virtual seller agent (VSA).

A dialectical system is a formal system that regulate persuasion dialogue, (See [19] for an overview). According to the game metaphor for social interactions, the parties are players which utter moves according to social rules.

**Definition 1 (Dialectical system).** *Let us consider  $\mathcal{L}$  a common object language and  $\mathcal{ACL}$  a common agent communication language. A **dialectical system** is a tuple  $DS = \langle P, \Omega_M, H, T, \mathit{proto}, Z \rangle$  where:*

- $P$  is a set of participants called players;
- $\Omega_M \subseteq \mathcal{ACL}$  is a set of well-formed moves;
- $H$  is a set of histories, the sequences of well-formed moves s.t. the speaker of a move is determined at each stage by the turn-taking function  $T$  and the moves agree with the dialogue-game protocol  $\mathit{proto}$ ;
- $T: H \rightarrow P$  is the turn-taking function;
- $\mathit{proto}: H \rightarrow 2^{\Omega_M}$  is the function determining the legal moves which are allowed to expand an history;
- $Z$  is the set of dialogues, i.e. the terminal histories.

Here,  $DS$  reflects the formalization of social interactions between players uttering moves during a dialogue. Each dialogue is a maximally long sequence of moves. Later to that, we specify informally the elements of  $DS$  for bilateral negotiation and information-seeking.

In our scenario, there are two players: the initiator `init` and, the responder `resp`, which utter moves each in turn. Since we address a proactive dialogue system, we consider the initiator to be a VSA. The **syntax** of moves is in conformance with a common **agent communication language**,  $\mathcal{ACL}$ . A move at time  $t$ : has an identifier, `mvt`; is uttered by a speaker (`spt ∈ P`) and the speech act is composed of a locution `loct` and a content `contentt`.

The possible locutions are: `question`, `assert`, `unknow`, `introduce`, `request`, `accept` and `reject`. The content consists of all instances of the following

1. VSA: `question(is(product, quilt) because search(user, product))`.
2. User: `assert(is(product, quilt))`.
3. VSA: `question(is(user, allergic))`.
4. User: `assert(is(user, allergic))`.
5. VSA: `question(is(user, sweat))`.
6. User: `assert(is(user, sweat))`.  
`question(is(product, organic))`.
7. VSA: `assert(is(product, organic))`.  
`question(dimension(product, 200, 200))`.
8. User: `assert(dimension(product, 200, 200))`.  
`question(budget(product, 350))`.
9. VSA: `introduce(is(product, quilt) because budget(product, 239) and is(product, nonallergenic) and is(product, organic))`.
10. User: `accept(is(product, quilt))`.
11. VSA: `introduce(is(product, quilt) because budget(product, 279) and is(product, nonallergenic) and is(product, fourseasons))`.
12. User: `accept(is(product, quilt))`.
13. VSA: `introduce(is(product, quilt)) and is(product, quiltcovers) because budget(product, 333.90) and is(product, nonallergenic) and is(product, fourseasons))`.
14. User: `reject(is(product, quiltcovers))`.

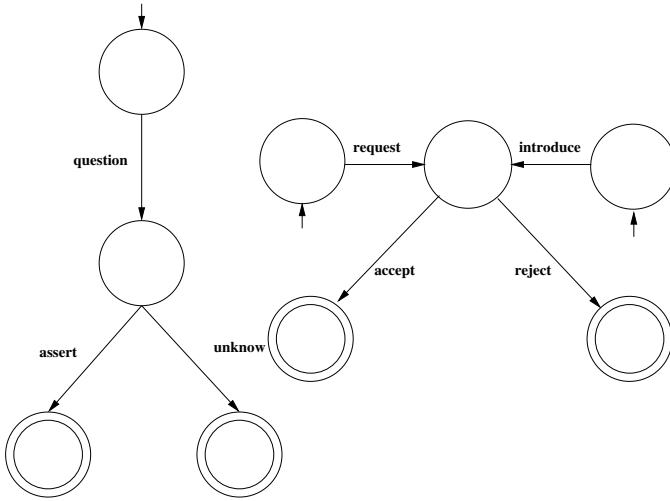
**Fig. 2.** A Possible Scenario Formalization

schemata " $S$  (because  $S'$ )" where  $S$  (eventually  $S'$ ) is a set of sentences in the common object language,  $\mathcal{L}$ . Actually, natural language utterances are interpreted/generated by the language dependent preprocessing of CSO (See Section 5). A move is an abstract representation of natural language utterances.

The dialogue in Fig 2 depicts a possible formalization of the natural language dialogue of Fig 1. It is worth noticing here that each utterance can contain more than one move.

In Fig. 3, we present our dialogue-game protocols by means of a deterministic finite-state automaton. An information-seeking dialogue begins with a question. The legal responding speech acts are **assert** and **unknown**. Two possible cases can occur: i) the dialogue is a failure if it is closed by an **unknown**; ii) the dialogue is a success if it is closed by an **assert**. A negotiation dialogue either begins with an offer from the VSA through the speech act **introduce** or the offer is suggested by the user through the speech act **request**. The legal responding speech acts are **accept** and **reject**. Here, the possibly occurring cases are: i) the dialogue is a failure if it is closed by a **reject**; ii) the dialogue is a success if it is closed by an **accept**.

The strategy interfaces with the dialogue-game protocol through the condition mechanism of utterances for a move. For example, at a certain point in the dialogue the VSA is able to send **introduce** or **question**. The choice of which locution and which content to send is depending on the VSA's strategy. For instance, the VSA is **benevolent** in the dialogue represented in Fig 2 since he first attempts to identify the dialogue's party needs, he continues with the product



**Fig. 3.** Dialogue-game protocol for information-seeking (on the left), and negotiation (on the right)

selection phase and then it terminates with the sale dialogue. An **aggressive** agent would consider the sale prior to anything whether the before-sale tasks have been performed or not.

## 7 Arguing over Utterances

In this section, we present how our computational model of argumentation for decision making [7] handles the dialogue strategy in order to generate and evaluate utterances.

In our framework, the knowledge is represented by a logical theory built upon an underlying logic-based language. In this language we distinguish between several different categories of predicate symbols. We use *goals* to represent the possible objectives of the decision making process (e.g. the dialogue to perform), *decisions* an agent can adopt (e.g. the move to utter) and a set of predicate symbols for *beliefs* (e.g. the previous utterance).

Assumptions here are required to carry on the reasoning process with incomplete knowledge, (e.g. some information about user’s needs are missing), and we need to express *preferences* between different goals (e.g. some dialogues are prior depending on the agent’s strategy). Finally, we allow the representation of explicit *incompatibilities* between goals, decisions and beliefs.

**Definition 2 (Decision framework).** A decision framework is a tuple  $DF = \langle \mathcal{DL}, \mathcal{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P} \rangle$ , where:

- $\mathcal{DL} = \mathcal{G} \cup \mathcal{D} \cup \mathcal{B}$  is a set of predicate symbols called the **decision language**, where we distinguish between goals ( $\mathcal{G}$ ), decisions ( $\mathcal{D}$ ) and beliefs ( $\mathcal{B}$ );

- $\mathcal{Asm}$  is a set of atomic formulae built upon predicates in  $\mathcal{DL}$  called **assumptions**;
- $\mathcal{I}$  is the **incompatibility relation**, i.e. a binary relation over atomic formulae in  $\mathcal{G}$ ,  $\mathcal{B}$  or  $\mathcal{D}$ . We require  $\mathcal{I}$  to be asymmetric;
- $\mathcal{T}$  is a logic **theory** built upon  $\mathcal{DL}$ ; statements in  $\mathcal{T}$  are clauses, each of them has a distinguished name;
- $\mathcal{P} \subseteq \mathcal{G} \times \mathcal{G}$  is the **priority** relation, namely a transitive, irreflexive and asymmetric relation over atomic formulae in  $\mathcal{G}$ .

In our framework, we consider multiple objectives which may or not be fulfilled by a set of decisions under certain circumstances. Additionally, we explicitly distinguish *assumable* (respectively *non-assumable*) literals which can (respectively cannot) be assumed to hold, as long as there is no evidence to the contrary. Decisions as well as some beliefs can be assumed. In this way, DF can model the incompleteness of knowledge.

The most natural way to represent conflicts in our object language is through-out some forms of logical negation. We consider two types of negation, as usual, (e.g., in extended logic programming), namely *strong negation*  $\neg$  (also called *explicit* or *classical negation*), and *weak negation*  $\sim$ , also called *negation as failure*. As a consequence we will distinguish between strong literals, i.e. atomic formula possibly preceded by strong negation, and weak literals, i.e. literals of the form  $\sim L$ , where  $L$  is a strong literal. The intuitive meaning of a strong literal  $\neg L$  is "L is definitely not the case", while  $\sim L$  intuitively means "There is no evidence that L is the case". The set  $\mathcal{I}$  of incompatibilities contains some *default* incompatibilities related to negation on the one hand, and to the nature of decision predicates on the other hand. Indeed, given an atom  $A$ , we have  $A \mathcal{I} \neg A$  as well as  $\neg A \mathcal{I} A$ . Moreover,  $L \mathcal{I} \sim L$ , whatever  $L$  is, representing the intuition that  $L$  is evidence to the contrary of  $\sim L$ . Notice, however, that we do not have  $\sim L \mathcal{I} L$ , as in the spirit of weak negation.

Other default incompatibilities are related to decisions, since different alternatives for the same decision predicate are incompatible with one another. Hence,  $D(a_1) \mathcal{I} D(a_2)$  and  $D(a_2) \mathcal{I} D(a_1)$ ,  $D$  being a decision predicate in  $\mathcal{D}$ , and  $a_1$  and  $a_2$  being different constants representing different<sup>4</sup> alternatives for  $D$ . Depending on the particular decision problem being represented by the framework,  $\mathcal{I}$  may contain further non-default incompatibilities. For instance, we may have  $g \mathcal{I} g'$ , where  $g, g'$  are different goals.

To summarize, the incompatibility relation captures the conflicts, either default or domain dependent, amongst decisions, beliefs and goals. The incompatibility relation can be easily lifted to set of sentences. We say that two sets of sentences  $\Phi_1$  and  $\Phi_2$  are *incompatible* (still denoted by  $\Phi_1 \mathcal{I} \Phi_2$ ) if there is a sentence  $\phi_1$  in  $\Phi_1$  and a sentence  $\phi_2$  in  $\Phi_2$  such that  $\phi_1 \mathcal{I} \phi_2$ .

A theory gathers the statements about the decision problem.

**Definition 3 (Theory).** A theory  $\mathcal{T}$  is an extended logic program, i.e a finite set of rules  $R: L_0 \leftarrow L_1, \dots, L_j, \sim L_{j+1}, \dots, \sim L_n$  with  $n \geq 0$ , each  $L_i$  (with

<sup>4</sup> Notice that in general a decision can be addressed by more than two alternatives.



$i \geq 0$ ) being a strong literal in  $\mathcal{L}$ .  $R$ , called the unique name of the rule, is an atomic formula of  $\mathcal{L}$ . All variables occurring in a rule are implicitly universally quantified over the whole rule. A rule with variables is a scheme standing for all its ground instances.

To simplify, we assume that names of rules are neither in the bodies nor in the head of the rules thus avoiding self-reference problems. We assume that the elements in the body of rules are independent. Besides, we suppose the decisions do not influence the beliefs and the decisions have no side effects.

In order to evaluate the relative importance of goals, we consider the *priority* relation  $\mathcal{P}$  over the goals in  $\mathcal{G}$ , which is transitive, irreflexive and asymmetric.  $G_1 \mathcal{P} G_2$  can be read " $G_1$  has priority over  $G_2$ ". There is no priority between  $G_1$  and  $G_2$ , either because  $G_1$  and  $G_2$  are *ex æquo* (denoted  $G_1 \simeq G_2$ ), or because  $G_1$  and  $G_2$  are not comparable.

We consider the dialogue formalized in Section 6. The generation and the evaluation of utterances by the VSA are captured by a *decision framework*  $DF = \langle \mathcal{DL}, \mathcal{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P} \rangle$  where:

- the decision language  $\mathcal{DL}$  distinguishes,
  - a set of **goals**  $\mathcal{G}$ . This set of literals identifies various motivations for driving the possible dialogues, negotiation (`negotiating(product)`) and information-seeking ones for product selection (`infoseeking(product)`) or need identification (`infoseeking(user)`),
  - a set of **decisions**  $\mathcal{D}$ . This set of literals identifies the possible utterances (e.g. `send(question(is(user, allergic)))`),
  - a set of **beliefs**, i.e. a set of literals identifying various situations identifying the possible queries of the user, (e.g. `receive(question(is(product, nonallergenic)))`, behavior through the website (e.g. `search(user, quilt)`) or the knowledge about the product/needs information (e.g. `is(user, allergic)`);
- the set of assumptions  $\mathcal{Asm}$  contains the possible decisions and the missing information about the user, (e.g. `~ is(user, allergic)`), or the product, (e.g. `~ is(product, nonallergenic)`);
- the incompatibility relation  $\mathcal{I}$  is trivially defined. For instance, `send(x) I send(y)`, with  $x \neq y$   
`infoseeking(topic1) I infoseeking(topic2)`, with  $topic_1 \neq topic_2$   
`negotiating(topic1) I infoseeking(topic2)` whatever  $topic_1$  and  $topic_2$  are
- the theory  $\mathcal{T}$  contains the rules in Table 2.
- If the VSA is benevolent, then the priority is defined such that:  
`infoseeking(user) P infoseeking(product)` and  
`infoseeking(product) P negotiating(product)`.  
 If the VSA is aggressive, then the priority is defined such that:  
`negotiating(product) P infoseeking(product)` and  
`infoseeking(product) P infoseeking(user)`.

Our formalization allows to capture the incomplete representation of a decision problem with assumable beliefs. It provides a knowledge base on top of which

**Table 2.** The rules of a Virtual Seller Agent (VSA)

$r_{11} : \text{infoseeking}(\text{user})$	$\leftarrow \text{send}(\text{question}(\text{is}(\text{user}, \text{allergic}))),$ $\sim \text{is}(\text{user}, \text{allergic}), \text{is}(\text{product}, \text{quilt}), \sim \text{receive}(x)$
$r_{12} : \text{infoseeking}(\text{user})$	$\leftarrow \text{send}(\text{question}(\text{is}(\text{user}, \text{sweat}))),$ $\sim \text{is}(\text{user}, \text{sweat}), \text{is}(\text{product}, \text{quilt}), \sim \text{receive}(x)$
$r_{21} : \text{infoseeking}(\text{product})$	$\leftarrow \text{send}(\text{question}(\text{is}(\text{product}, \text{quilt}))),$ $\text{search}(\text{user}, \text{quilt}), \sim \text{is}(\text{product}, \text{quilt})$
$r_{22} : \text{infoseeking}(\text{product})$	$\leftarrow \text{send}(\text{question}(\text{is}(\text{product}, \text{nonallergenic}))),$ $\sim \text{is}(\text{product}, \text{nonallergenic}), \sim \text{receive}(x)$
$r_{23} : \text{infoseeking}(\text{product})$	$\leftarrow \text{send}(\text{question}(\text{is}(\text{product}, \text{organic}))),$ $\sim \text{is}(\text{product}, \text{organic}), \sim \text{receive}(x)$
$r_{24} : \text{infoseeking}(\text{product})$	$\leftarrow \text{send}(\text{question}(\text{dimension}(\text{product}, x, y))),$ $\sim \text{dimension}(\text{product}, x, y), \sim \text{receive}(z)$
$r_{25} : \text{infoseeking}(\text{product})$	$\leftarrow \text{send}(\text{question}(\text{budget}(\text{product}, x))),$ $\sim \text{budget}(\text{product}, x), \sim \text{receive}(y)$
$r_{26} : \text{infoseeking}(\text{product})$	$\leftarrow \text{send}(\text{assert}(\text{is}(x, y))), \text{receive}(\text{question}(\text{is}(x, y))), \text{is}(x, y)$
$r_{27} : \text{infoseeking}(\text{product})$	$\leftarrow \text{send}(\text{assert}(\neg \text{is}(x, y))), \text{receive}(\text{question}(\text{is}(x, y))), \neg \text{is}(x, y)$
$r_{28} : \text{infoseeking}(\text{product})$	$\leftarrow \text{send}(\text{unknow}(\text{is}(x, y))), \text{receive}(\text{question}(\text{is}(x, y))), \sim \text{is}(x, y)$
$r_{29} : \text{negotiating}(\text{product})$	$\leftarrow \text{send}(\text{introduce}(\text{product})), \text{budget}(\text{product}, y)$
$r_{31} : \text{budget}(\text{product}, 350)$	$\leftarrow \text{is}(\text{product}, \text{nonallergenic}),$ $\text{is}(\text{product}, \text{organic}), \text{dimension}(\text{product}, 200, 200)$
$r_{32} : \text{is}(\text{product}, \text{nonallergenic})$	$\leftarrow \text{is}(\text{user}, \text{allergic})$
$r_{33} : \text{is}(\text{product}, \text{organic})$	$\leftarrow \text{is}(\text{user}, \text{sweat})$

arguments are built in order to reach decisions. We adopt here a tree-like structure for arguments.

**Definition 4 (Argument).** Let  $DF = \langle \mathcal{DL}, \text{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P}, \mathcal{RV} \rangle$  be a decision framework. An **argument**  $\bar{a}$  deducing the **conclusion**  $c \in \mathcal{DL}$  (denoted  $\text{conc}(\bar{a})$ ) supported by a set of **assumptions**  $A$  in  $\text{Asm}$  (denoted  $\text{asm}(\bar{a})$ ) is a tree where the root is  $c$  and each node is a sentence of  $\mathcal{DL}$ . For each node :

- if the node is a leaf, then it is either an assumption in  $A$  or  $\top$ <sup>5</sup>;
- if the node is not a leaf and it is  $\alpha \in \mathcal{DL}$ , then there is an inference rule  $\alpha \leftarrow \alpha_1, \dots, \alpha_n$  in  $\mathcal{T}$  and,
  - either  $n = 0$  and  $\top$  is its only child,
  - or  $n > 0$  and the node has  $n$  children,  $\alpha_1, \dots, \alpha_n$ .

The sentences of  $\bar{a}$  (denoted  $\text{sent}(\bar{a})$ ) is the set of literals of  $\mathcal{DL}$  in the bodies/heads of the rules including the assumptions of  $\bar{a}$ . We write  $\bar{a} : A \vdash \alpha$  to denote an argument  $\bar{a}$  such that  $\text{conc}(\bar{a}) = \alpha$  and  $\text{asm}(\bar{a}) = A$ . The set of arguments built upon  $DF$  is denoted by  $\mathcal{A}(DF)$ .

Arguments are built by reasoning backwards. Additionally, arguments interact with one another, and consequently, we reach to define the following attack relation.

**Definition 5 (Attack relation).** Let  $DF = \langle \mathcal{DL}, \text{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P} \rangle$  be a decision framework, and  $\bar{a}, \bar{b} \in \mathcal{A}(DF)$  be two arguments.  $\bar{a}$  **attacks**  $\bar{b}$  iff  $\text{sent}(\bar{a}) \mathcal{I} \text{sent}(\bar{b})$ .

<sup>5</sup>  $\top$  denotes the unconditionally true statement.

This relation encompasses both the direct (often called *rebuttal*) attack due to the incompatibility of the conclusions, and the indirect (often called *undermining*) attack, (i.e., directed to a "subconclusion").

Since the goals promoted by arguments have different priorities, the arguments interact with one another. For this purpose, we define the strength relation between concurrent arguments. Arguments are *concurrent* if their conclusions are identical or incompatible.

**Definition 6 (Strength relation).** *Let  $DF = \langle \mathcal{DL}, \text{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P} \rangle$  be a decision framework and  $\bar{a}_1, \bar{a}_2 \in \mathcal{A}(DF)$  be two arguments which are concurrent.  $\bar{a}_1$  is stronger than  $\bar{a}_2$  (denoted  $\bar{a}_1 \mathcal{P} \bar{a}_2$ ) iff  $\text{conc}(\bar{a}_1) = g_1 \in \mathcal{G}$ ,  $\text{conc}(\bar{a}_2) = g_2 \in \mathcal{G}$  and  $g_1 \mathcal{P} g_2$ .*

Due to the definition of  $\mathcal{P}$  over  $\mathcal{T}$ , the relation  $\mathcal{P}$  is transitive, irreflexive and asymmetric over  $\mathcal{A}(DF)$ .

The attack relation and the strength relation can be combined to adopt Dung's calculus of opposition as in [20]. We distinguish between one argument attacking another, and that attack succeeding due to the strength of arguments.

**Definition 7 (Defeat relation).** *Let  $DF = \langle \mathcal{DL}, \text{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P} \rangle$  be a decision framework and  $\bar{a}$  and  $\bar{b}$  be two structured arguments.  $\bar{a}$  defeats  $\bar{b}$  iff:*

1.  $\bar{a}$  attacks  $\bar{b}$ ;
2. and it is not the case that  $\bar{b} \mathcal{P} \bar{a}$ .

Similarly, we say that a set  $S$  of structured arguments defeats a structured argument  $\bar{a}$  if  $\bar{a}$  is defeated by one argument in  $S$ .

Let us consider this example:

*Example 1 (Defeat relation).* Let us consider the situation after the second move in the dialogue represented in Fig. 11.

The argument  $\bar{a}$  concludes `infoseeking(user)` since the VSA can ask to the user if he is allergic, (i.e. `question(is(user,allergic))`), the VSA is not yet aware about it, (i.e.  $\sim \text{is}(\text{user}, \text{allergic})$ ), the user is looking for a quilt, (i.e. `is(product,quilt)`), and the user did not query the VSA, (i.e.  $\sim \text{receive}(x)$ ). The argument  $\bar{b}$  concludes `infoseeking(product)` since the VSA can ask to the user if the product must be nonallergenic, (i.e. `send(question(is(product,nonallergenic))`), the VSA is not yet aware about it (i.e.  $\sim \text{is}(\text{product}, \text{nonallergenic})$ ) and the user did not query the VSA ( $\sim \text{receive}(x)$ ). While  $\bar{a}$  is built upon  $r_{11}$ ,  $\bar{b}$  is built upon  $r_{22}$ . Since these arguments suppose different decisions, they attack each others. If the VSA is benevolent, it is not the case that `infoseeking(product)Pinfoseeking(user)` and so  $\bar{a}$  defeats  $\bar{b}$ . If the VSA is aggressive, it is not the case that `infoseeking(user)Pinfoseeking(product)` and so  $\bar{a}$  defeats  $\bar{b}$ .

In our argumentation-based approach for dialogue strategy, arguments motivate decisions and they can also be defeated by other arguments. Formally, our argumentation framework (AF for short) is defined as follows.

**Definition 8 (AF).** Let  $DF = \langle \mathcal{DL}, \mathcal{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P} \rangle$  be a decision framework. The argumentation framework for decision making built upon  $DF$  is a pair  $AF = \langle \mathcal{A}(DF), \mathbf{defeats} \rangle$  where  $\mathcal{A}(DF)$  is the finite set of arguments built upon  $DF$  as defined in Definition 8, and  $\mathbf{defeats} \subseteq \mathcal{A}(DF) \times \mathcal{A}(DF)$  is the binary relation over  $\mathcal{A}(DF)$  as defined in Definition 7.

We adapt Dung’s extension-based semantics in order to analyze whenever a set of arguments can be considered as subjectively justified with respect to the agent’s priority.

**Definition 9 (Semantics).** Let  $DF = \langle \mathcal{DL}, \mathcal{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P} \rangle$  be a decision framework and  $AF = \langle \mathcal{A}(DF), \mathbf{defeats} \rangle$  be our argumentation framework for decision making. For  $S \subseteq \mathcal{A}(DF)$  a set of arguments, we say that:

- $S$  is conflict-free iff  $\forall \bar{a}, \bar{b} \in S$  it is not the case that  $\bar{a}$  defeats  $\bar{b}$ ;
- $S$  is admissible iff  $S$  is conflict-free and  $S$  defeats every argument  $\bar{a}$  such that  $\bar{a}$  defeats some argument in  $S$ ;

Here, we only consider admissibility but other Dung’s extension-based semantics [21] can easily be adapted.

Formally, given an argument  $\bar{a}$ , let

$$\mathbf{dec}(\bar{a}) = \{D(a) \in \mathbf{asm}(\bar{a}) \mid D \text{ is a decision predicate}\}$$

be the set of decisions supported by the argument  $\bar{a}$ .

The decisions are *suggested* to reach a goal if they are supported by an argument concluding this goal and this argument is a member of an admissible set of arguments.

**Definition 10 (Credulous decisions).** Let  $DF = \langle \mathcal{DL}, \mathcal{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P} \rangle$  be a decision framework,  $g \in \mathcal{G}$  be a goal and  $D \subseteq \mathcal{D}$  be a set of decisions. The decisions  $D$  **credulously argue for**  $g$  iff there exists an argument  $\bar{a}$  in an admissible set of arguments such that  $\mathbf{conc}(\bar{a}) = g$  and  $\mathbf{dec}(\bar{a}) = D$ . We denote  $\mathbf{val}_c(D)$  the set of goals in  $\mathcal{G}$  for which the set of decisions  $D$  credulously argues.

It is worth noticing here that the decisions that credulously argue for a goal cannot contain mutual exclusive alternatives for the same decision predicate. This is due to the fact that an admissible set of arguments is conflict-free.

If we consider the arguments  $\bar{a}$  and  $\bar{b}$  supporting the decisions  $D(a)$  and  $D(b)$  respectively where  $a$  and  $b$  are mutually exclusive alternatives, we have  $D(a) \mathcal{I} D(b)$  and  $D(a) \mathcal{I} D(b)$  and so, either  $\bar{a}$  **defeats**  $\bar{b}$  or  $\bar{b}$  **defeats**  $\bar{a}$  or both of them depending on the strength of these arguments.

**Proposition 1 (Mutual exclusive alternatives).** Let  $DF = \langle \mathcal{DL}, \mathcal{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P} \rangle$  be a decision framework,  $g \in \mathcal{G}$  be a goal and  $AF = \langle \mathcal{A}(DF), \mathbf{defeats} \rangle$  be the argumentation framework for decision making built upon  $DF$ . If  $S$  be an admissible set of arguments such that, for some  $\bar{a} \in S$ ,  $g = \mathbf{conc}(\bar{a})$  and  $D(a) \in \mathbf{asm}(\bar{a})$ , then  $D(b) \in \mathbf{asm}(\bar{a})$  iff  $a = b$ .

However, it is worth highlighting here the fact that mutual exclusive decisions can be suggested for the same goal through different admissible set of arguments. This case reflects the credulous nature of our semantics.

**Definition 11 (Skeptical decisions).** Let  $DF = \langle \mathcal{DL}, \mathcal{Psm}, \mathcal{I}, \mathcal{T}, \mathcal{P}, \mathcal{RV} \rangle$  be a decision framework,  $g \in \mathcal{G}$  be a goal and  $D \subseteq \mathcal{D}$  be a set of decisions. The decisions  $D$  **skeptically argue for**  $g$  iff for all admissible set of arguments  $S$  such that for some arguments  $\bar{a}$  in  $S$   $\text{conc}(\bar{a}) = g$ , then  $\text{dec}(\bar{a}) = D$ . We denote  $\text{val}_s(D)$  the set of goals in  $\mathcal{G}$  for which the set of decisions  $D$  skeptically argues.

Due to the uncertainties, some decisions satisfy goals for sure if they skeptically argue for them, or some decisions can possibly satisfy goals if they credulously argue for them. While the first case is required for convincing a risk-averse agent, the second case is enough to convince a risk-taking agent. Since some ultimate choices amongst various justified sets of alternatives are not always possible, we will consider in this paper only risk-taking agents.

Since agents can consider multiple objectives which may not be fulfilled all together by a set of non-conflicting decisions, high-ranked goals must be preferred to low-ranked goals.

**Definition 12 (Preferences).** Let  $DF = \langle \mathcal{DL}, \mathcal{Asm}, \mathcal{I}, \mathcal{T}, \mathcal{P}, \mathcal{RV} \rangle$  be a decision framework. We consider  $\mathcal{G}, \mathcal{G}'$  two set of goals in  $\mathcal{G}$  and  $\mathcal{D}, \mathcal{D}'$  two set of decisions in  $\mathcal{D}$ .  $\mathcal{G}$  is **preferred to**  $\mathcal{G}'$  (denoted  $\mathcal{G} \mathcal{P} \mathcal{G}'$ ) iff

1.  $\mathcal{G} \supseteq \mathcal{G}'$ , and
2.  $\forall g \in \mathcal{G} \setminus \mathcal{G}'$  there is no  $g' \in \mathcal{G}'$  such that  $g' \mathcal{P} g$ .

$\mathcal{D}$  is **preferred to**  $\mathcal{D}'$  (denoted  $\mathcal{D} \mathcal{P} \mathcal{D}'$ ) iff  $\text{val}_c(\mathcal{D}) \mathcal{P} \text{val}_c(\mathcal{D}')$ .

Formally, let

$\mathcal{AD} = \{D \mid D \subseteq \mathcal{D} \text{ such that } \forall D' \subseteq \mathcal{D} \text{ it is not the case that } \text{val}_c(D') \mathcal{P} \text{val}_c(D)\}$

be the decisions which can be accepted by the agent. Additionally, let

$\mathcal{AG} = \{G \mid G \subseteq \mathcal{G} \text{ such that } G = \text{val}_c(D)\}$

be the goals which can be reached by the agent.

Let us consider now the VSA's decision problem after the second move.

*Example 2 (Semantics).* The argument  $\bar{a}$  (respectively  $\bar{b}$ ) (described in Example 1), concludes  $\text{infoseeking}(\text{user})$  (respectively  $\text{infoseeking}(\text{product})$ ). Actually, the decisions  $\{\text{send}(\text{question}(\text{is}(\text{user}, \text{allergic})))\}$  credulously argue for  $\text{infoseeking}(\text{user})$  and the decisions  $\{\text{send}(\text{question}(\text{is}(\text{product}, \text{nonallergenic})))\}$  credulously argue for  $\text{infoseeking}(\text{product})$ . If the VSA is benevolent, then  $\{\text{send}(\text{question}(\text{is}(\text{user}, \text{allergic})))\}$  is an acceptable set of decisions. If the VSA is aggressive, then  $\{\text{send}(\text{question}(\text{is}(\text{product}, \text{nonallergenic})))\}$  is an acceptable set of decisions.

## 8 Related Works

Amgoud & Prade in [22] are presenting a general and abstract argumentation framework for multi criteria decision making. This framework captures the mental states (goals, beliefs and preferences) of the decision makers. Therefore, in their framework the arguments are prescribing actions to reach goals if these actions are feasible while certain circumstances are true. These arguments - *that eventually conflict* - are balanced according to their strengths. The argumentation framework we proposed earlier in this paper is conforming with this approach while being more specific and concrete.

The argumentation-based decision making process envisaged in [22] is divided into different steps where the arguments are successively constructed, weighted, confronted and evaluated. However, the computations we proposed earlier in this paper go through the construction of arguments, the construction of counterarguments, the evaluation of the generated arguments and the relaxation of preferences for making concessions. It is also worth noticing here that: a) the model we propose is unique in making it finally possible to concede, b) our argumentation-based decision process suggest some decisions even if low-ranked goals cannot be reached.

Bench-Capon & Prakken formalize in [23] defeasible argumentation for practical reasoning. As in [22], they select the best course of actions by confronting and evaluating arguments. Bench-Capon & Prakken focus on the abductive nature of practical reasoning which is directly modelled within in our framework.

Kakas & Moraitis propose in [24] an argumentation-based framework for decision making of autonomous agents. For this purpose, the knowledge of the agent is split and localized in different modules representing different capabilities. Whereas [24] is committed to one argumentation semantics, we can still deploy our framework for a number of such semantics by relying on assumption-based argumentation.

Finally, to the best of our knowledge, few implementation of argumentation over actions exist. CaSAPI<sup>6</sup> [25] and DeLP<sup>7</sup> [26] are restricted to the theoretical reasoning. GORGIAS<sup>8</sup> [27] implements an argumentation-based framework to support the decision making of an agent within a modular architecture.

## 9 Conclusions

In this paper, we have presented a dialogue management system that applies argumentation for generating and evaluating utterances. The agent start the conversation with the prior task which can consist of the need identification, the product selection or the negotiation depending on its strategy. During the dialogue, a proactive agent can query the user. Additionally, it can introduce a

<sup>6</sup> <http://www.doc.ic.ac.uk/~dg00/casapi.html>

<sup>7</sup> <http://lidia.cs.uns.edu.ar/DeLP>

<sup>8</sup> <http://www.cs.ucy.ac.cy/~nkd/gorgias/>

product to sell and justify this choice depending on the information collected in the previous steps.

In order for us to implement an agent's reasoning method we are considering MARGO<sup>9</sup> (A Multiattribute ARGumentation framework for Opinion explanation), which is an argumentation-based mechanism for decision-making [7]. We are currently rewriting MARGO in Java so that issues related to improving its performance, (i.e., the response time), and its scalability, (i.e., the number of rules which can be managed), are better tackled. This work is required to provide an industrial application rather than a research prototype. Besides, we need to interface this argumentation-based engine with the CSO Artificial Solutions' Language Processor in order to build conversational agents which are proactive in different selling situations.

Although the negotiation dialogue model we proposed allows single-sellings through the exchange of proposals and counter-proposals. However, we are currently working on an extension that will address cross-selling and up-selling.

## Acknowledgements

This work is supported by the Ubiquitous Virtual Seller (VUU) project that was initiated by the Competitvity Institute on Trading Industries (PICOM).

## References

1. Hof, R., Green, H., Himmelstein, L.: Now it's YOUR WEB. *BusinessWeek*, 68–75 (1998)
2. Poong, Y., Zaman, K.-U., Talha, M.: E-commerce today and tomorrow: a truly generalized and active framework for the definition of electronic commerce. In: *Proc. of the 8th International Conference on Electronic Commerce (ICEC)*, pp. 553–557. ACM, Fredericton (2006)
3. Palopoli, L., Rosaci, D., Ursino, D.: Agents' roles in B2C e-commerce. *AI Communications* 19, 95–126 (2006)
4. Wooldridge, M., Jennings, N.: Intelligent agents: Theory and practice. *Knowledge Engineering Review* 10, 115–152 (1995)
5. Isbister, K., Doyle, P.: Design and evaluation of embodied conversational agents: A proposed taxonomy. In: *Proc. of the First International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, Budapest, Hungary (2002)
6. Rist, T., André, E., Baldes, S., Gebhard, P., Klesen, M., Rist, P., Schmitt, M.: A review of the development of embodied presentation agents and their application fields. In: *Life-Like Characters – Tools, Affective Functions, and Applications*, pp. 377–404. Springer (2003)
7. Morge, M.: The hedgehog and the fox. In: Rahwan, I., Parsons, S., Reed, C. (eds.) *Argumentation in Multi-Agent Systems*. LNCS (LNAI), vol. 4946, pp. 114–131. Springer, Heidelberg (2008)
8. Morge, M., Mancarella, P.: Assumption-based argumentation for the minimal concession strategy. In: McBurney, P., Rahwan, I., Parsons, S., Maudet, N. (eds.) *ArgMAS 2009*. LNCS, vol. 6057, pp. 114–133. Springer, Heidelberg (2010)

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

9. Roberts, F., Gülsdorff, B.: Techniques of dialogue simulation. In: Pelachaud, C., Martin, J.-C., André, E., Chollet, G., Karpouzis, K., Pelé, D. (eds.) IVA 2007. LNCS (LNAI), vol. 4722, pp. 420–421. Springer, Heidelberg (2007)
10. Ferguson, J.A.G.: Mixed-initiative systems for collaborative problem solving. *AI Magazine* 28, 23–32 (2007)
11. Rich, C., Sidner, C.L., Lesh, N.: COLLAGEN applying collaborative discourse theory to human-computer interaction. *AI Magazine* 22, 15–25 (2001)
12. Sadek, D.: Artemis Rational Dialogue Agent Technology: An Overview. In: Multi-Agent Programming, pp. 217–225. Springer, Heidelberg (2005)
13. Grosz, B.J., Sidner, C.L.: Plans for Discourse. In: Cohen, Morgan, Pollack (eds.) *Intentions and Plans in Communication and Discourse*. MIT press, Cambridge (1990)
14. Rao, A.S., Georgeff, M.P.: Modeling rational agents within a BDI-architecture. In: Proc. of the 2nd International Conference on Principles of Knowledge Representation and Reasoning (KR), pp. 473–484. Morgan Kaufmann publishers Inc., San Mateo (1991)
15. C, F.T.: FIPA ACL Communicative Act Library Specification. Component, Foundation for Intelligent Physical Agents (2002), <http://fipa.org/specs/fipa00037/>
16. Breiter, P.: La communication orale coopérative : contribution à la modélisation et à la mise en œuvre d'un agent rationnel dialoguant. PhD thesis, Université de Paris Nord (1992)
17. Hamblin, C.L.: *Fallacies*. Methuen (1970)
18. Walton, D., Krabbe, E.: *Commitment in Dialogue*. SUNY Press (1995)
19. Prakken, H.: Formal systems for persuasion dialogue. *The Knowledge Engineering Review* 21, 163–188 (2006)
20. Amgoud, L., Cayrol, C.: On the acceptability of arguments in preference-based argumentation. In: Proc. of UAI, pp. 1–7. Morgan Kaufmann, Wisconsin (1998)
21. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artif. Intell.* 77, 321–357 (1995)
22. Amgoud, L., Prade, H.: Using arguments for making and explaining decisions. *Artificial Intelligence Journal* 173(3-4), 413–436 (2009)
23. Bench-Capon, T., Prakken, H.: Justifying actions by accruing arguments. In: Proc. of the 1st International Conference on Computational Models of Argument, pp. 247–258. IOS Press, Amsterdam (2006)
24. Kakas, A., Moraitis, P.: Argumentative-based decision-making for autonomous agents. In: Proc. of AAMAS, pp. 883–890. ACM Press, New York (2003)
25. Gartner, D., Toni, F.: CaSAPI: a system for credulous and sceptical argumentation. In: Proc. of ArgNMR, pp. 80–95 (2007)
26. García, A.J., Simari, G.R.: Defeasible logic programming: an argumentative approach. *Theory and Practice of Logic Programming* 4(2), 95–138 (2004)
27. Demetriou, N., Kakas, A.C.: Argumentation with abduction. In: Proc. of the 4th Panhellenic Symposium on Logic (2003)



# Reasoning about Trust Using Argumentation: A Position Paper

Simon Parsons<sup>1,2</sup>, Peter McBurney<sup>3</sup>, and Elizabeth Sklar<sup>1,2</sup>

<sup>1</sup> Department of Computer & Information Science, Brooklyn College,  
City University of New York, 2900 Bedford Avenue, Brooklyn, NY 11210 USA  
{sklar, parsons}@sci.brooklyn.cuny.edu

<sup>2</sup> Department of Computer Science, Graduate Center  
City University of New York, 365 Fifth Avenue, New York, NY 10016, USA

<sup>3</sup> Department of Computer Science, University of Liverpool,  
Ashton Building, Ashton Street, Liverpool, L69 3BX, United Kingdom  
mcburney@liverpool.ac.uk

**Abstract.** Trust is a mechanism for managing the uncertainty about autonomous entities and the information they store, and so can play an important role in any decentralized system. As a result, trust has been widely studied in multiagent systems and related fields such as the semantic web. Managing information about trust involves inference with uncertain information, decision making, and dealing with commitments and the provenance of information, all areas to which systems of argumentation have been applied. Here we discuss the application of argumentation to reasoning about trust, identifying some of the components that an argumentation-based system for reasoning about trust would need to contain and sketching the work that would be required to provide such a system.

## 1 Introduction

Trust is a mechanism for managing the uncertainty about autonomous entities and the information they store. As a result trust can play an important role in any decentralized system. As computer systems have become increasingly distributed, and control in those systems has become more decentralized, trust has steadily become more important in computer science. Trust is an especially important issue from the perspective of autonomous agents and multiagent systems. The premise behind the multiagent systems field is that of developing software agents that will work in the interests of their owners, carrying out their owners' wishes while interacting with other entities. In such interactions, agents will have to reason about the amount that they should trust those other entities, whether they are trusting those entities to carry out some task, or whether they are trusting those entities to not misuse crucial information.

This paper argues that systems of argumentation have an important role to play in reasoning about trust. We start in Section 2 by briefly reviewing work that defines important aspects of trust and giving an extended example which illustrates some of these aspects. Section 3 then briefly reviews some of the work on reasoning about trust and identifies some of the characteristics of any effective system for dealing with trust information. Building on this discussion, Section 4 then argues that systems of argumentation can handle trust and sketches a specific system of argumentation for doing this. Section 5 concludes.

## 2 Trust

As a number of authors have pointed out, trust is a concept that is both complex and rather difficult to pin down precisely. As a result, there are a number of different definitions in the literature. To pick a few specific examples, Sztompka [27] (cited in [7]) suggests that:

Trust is a bet about the future contingent actions of others.

while Mcknight and Chervany [21], drawing on a range of existing definitions, define trust as:

Trust is the extent to which one party is willing to depend on something or somebody in a given situation with a feeling of relative security, even though negative consequences are possible.

and Gambetta [4] states:

Trust is the subjective probability by which an individual, A, expects that another individual, B, performs a given action on which its welfare depends.

While these definitions differ, there are clearly some common elements. There is a degree of uncertainty associated with trust — whether expressed as a subjective probability, as a bet (which, of course, can be expressed as a subjective probability [11]), or as a “feeling of security”. Trust is tied up with the relationships between individuals. Trust is related to the actions of individuals and how those actions affect others.

It is also pointed out in a number of places that there are different kinds of trust, what Jøsang et al. [12] call “trust scopes”. For example, [9] identifies the following types of trust:

1. *Provision trust*: the trust that exists between the user of a service or resource, and the provider of that resource.
2. *Access trust*: the trust that exists between the owner of a resource and those that are accessing those resources.
3. *Delegation trust*: the trust that exists between an individual who delegates responsibility for some action or decision and the individual to which that action or decision is delegated.
4. *Identity trust*: trust that an individual is who they claim to be.
5. *Context trust*: trust that an individual has in the existence of sufficient infrastructure to support whatever activities that individual is engaged in.

We illustrate some of these different types of trust with the following example.

Alice is planning a picnic for a group of friends. She asks around amongst some of her acquaintances for ideas about where to hold the picnic. Bob suggests a park a little way outside of the city where he goes quite regularly (provision trust, relating to information) — he says it is quiet and easy to get to. Carol says she has never been to the park herself, but has heard that the bugs are terrible (provision trust, relating to information).

Alice decides that the picnic will be a potluck<sup>1</sup>. Alice asks David to bring potato salad (delegation trust) and Eric says he will bring bread from the bakery near his house (provision trust, relating to a good). Fran offers to bake a cake (provision trust, relating to a good). Carol says she will make her famous barbecue chicken, cooking it on the public barbecues that Alice believes are provided by the park (context trust).

The picnic is scheduled for midday. George arranges to pick up Alice from her house at 10am in order to drive her to the park (Alice doesn't have a car). Harry, who can borrow a minivan (access trust), offers to collect several people from their homes and stop on the way to buy a case of beer. Iain, who is going to ride with George, says he'll bring a soccer ball so they can all play after lunch. John asks if he can bring a friend of a friend, Keith, whom John has never met, but whom John knows will be visiting the city and is unoccupied that day (identity trust).

As Alice makes the arrangements, she is obviously trusting a lot of people to make sure that the plan comes together in ways that are rather distinct.

Bob and Carol are providing information. To decide whether to go to the park, Alice has to factor in the trustworthiness of that information. She has to take into account how reliable Bob and Carol are as information providers, not least because the information that they have given here is contradictory. Alice might judge that what she knows about Bob (that he goes to the park often) makes him more trustworthy than Carol in this regard (though in other contexts, such as when deciding what film to see, she might value Carol's opinion more), and the fact that Carol is relying on information from yet another person might strengthen this feeling (or, equally, make Alice value Carol's opinion about the park less).

The trust involved in handling the information from Bob and Carol seems to be somewhat different to the handling of trust when considering the makeup of the meal. Here Alice has to balance not the reliability of the information that people provide, but the *commitments* they are making, the extent to which Carol, David, Eric, Fran, George, Harry and Iain will do what they say they will do. Carol may be a terribly unreliable source of information about parks, and thus untrustworthy in that regard, but a superb provider of barbecued chicken, and one who has never failed to bring that chicken to a potluck when she says that she will. In contrast, Alice may know that Fran saying she will bake a cake means very little. She is just as likely to bake cookies, or realise late the night before the picnic that she has no flour and will have to bring a green salad instead (thus ruining the meal). David, on the other hand, is quite likely not to make potato salad; but if he doesn't, he can be relied upon to substitute it with some close approximation, a pasta salad for example.

In other words, an individual can be an untrustworthy source of information, but a trustworthy provider of services, or indeed an untrustworthy provider of services but a very reliable information source (it is perfectly possible that Fran only ever provides correct information despite her food-related flakiness) — there are different dimensions of trust for different services that are provided (here, information and food items).

---

<sup>1</sup> "Pot luck" means that all the guests are expected to bring something that will contribute to the meal, typically an item of food or a beverage.

We distinguish this by talking of the *context* of trust. Similarly, the failure of an individual to fulfill their commitments is not necessarily binary — how they fail can be important.

There are also other aspects to the failure of a commitment. Actions have time and location components. If George is a few minutes late picking Alice up, it may not affect the picnic. If he is an hour late, that might be catastrophic. If he has the wrong address, then even if he arrives at that (wrong) location at 10am, the success of the picnic is in danger. And if Harry can't find his way to the park, there won't be any soccer after lunch even if he successfully collected everyone and bought the beer just as he said he would. However, as long as he arrives while the picnic is going on, then his passengers have a chance to enjoy themselves, though the later he arrives, the less chance that they will have a good time.

### 3 Reasoning about Trust

As discussed above, a key aspect of trust is that it stems from the relationship between individuals or groups of individuals. This means that it is a relative notion — Alice and Bob may have different views about Carol's trustworthiness — and thus that *provenance* is important in reasoning about trust [6]. A situation that often arises is one where it is necessary to combine different people's information about trust and when this is done, it is important to know where information about trust is coming from.

In this context, Jøsang et al. [12] distinguish between *functional* trust, the trust in an individual to carry out some task, and *referral* trust, the trust in an individual's recommendation. Thus, in our example, Alice's reasoning about George's offer of a lift, and Carol's offer to bring chicken are *functional* trust — Alice is thinking about George's reliability as a provider of lifts and Carol's reliability as a provider of chicken. However, if Alice were to ask Carol for a recommendation for a good butcher, then Alice would base her assessment of Carol's answer on her (Alice's) assessment of Carol's ability to make good recommendations, an instance of referral trust, while what Carol expresses about her butcher is another instance of functional trust.

As [12] points out (in terms of our example), the fact that Carol trusts her butcher to supply good meat is not necessarily a reason for Alice to do the same, and it certainly isn't a reason for Alice to trust the butcher in any more general context (to do a good job of painting Alice's house, for example). However, under certain circumstances — and in particular when the trust context is the same, as it is when Alice is considering the use of Carol's butcher as a provider of meat [14] — it is reasonable to consider trust to be transitive. Thus Alice can consider combining her direct assessment of Carol's referral trustworthiness in the food domain, with Carol's direct assessment of her butcher's functional trustworthiness to derive an *indirect* functional assessment of the butcher.

Given this transitivity, the notion of a *trust network* then makes sense. If Alice can estimate the referral trustworthiness of her friends, and they can do the same for their friends, then Alice can make judgements about recommendations she receives not just

<sup>2</sup> Depending on the butcher, of course, even this might be too broad a trust context — perhaps the butcher provides excellent chicken and beef, but can only supply indifferent pork and his game is never hung for long enough.

from her friends, but also from the friends of her friends, and their friends and so on. The question is, what is a reasonable way to represent this computationally, taking into account that each of these friends trusts their friends to different degrees.

At the moment there is no definitive answer to the question<sup>3</sup>. As the definitions of trust cited above suggest, one way to model trust is to use some form of subjective probability — Alice’s degree of trust in Bob’s park recommendation is a measure of her belief that she will like the park since Bob says that he likes the park. *EigenTrust* [15] is a mechanism, derived for use in peer-to-peer networks, for establishing a global trust rating that estimates how much any individual should trust another. While such a global rating, based as it is on performance, is reasonable for peer-to-peer systems, it has been argued [16] that in the kind of social networks we are discussing here, it is necessary to capture the fact that, for example, Alice and Bob can have very different estimations of Carol’s trustworthiness (and, as we have argued, that they will have different ratings for Carol’s trustworthiness in different contexts).

Subjective logic [13] is a formalism for capturing exactly this aspect of trust, and for inferring the degree of trust existing between two nodes in a trust network. Based on the Dempster-Shafer theory of evidence [26] it computes a measure that is a generalisation of probability, distinguishing belief in the reliability of an individual, disbelief in the reliability, and the potential belief that has not yet been determined one way or another (termed the “uncertainty”). Singh and colleagues [10,28] provide extensions of the approach, the former looking at how best to update the measure of trust one individual has in another depending on their experience of interactions. Thus Alice may have her high regard for Carol’s food-related recommendations damaged by a bad experience with a supplier that Carol recommends. Subjective logic is not the only approach to handling this problem. For example, Katz and Golbeck [16] describe an algorithm called TidalTrust for establishing the trust between a *source* node (representing the individual doing the trusting) and a *sink* node (representing the individual being trusted). Later work by Kuter and Golbeck provides the SUNNY algorithm [18] which is reported to outperform TidalTrust on a benchmark database of trust information.

## 4 Argumentation and Trust

The Trust field, including sample literature discussed above, gives us methodologies for *computing* trust, while the Argumentation field can give us methodologies for *reasoning* about trust. In short, we believe that argumentation can provide a mechanism for handling many of the aspects that we need to capture about trust, as we discuss at some length in this section.

### 4.1 Argumentation in General

As we have discussed above, there are two major aspects that need to be handled by any representation of trust — we need to handle measures of trust, and we need to handle

---

<sup>3</sup> Indeed, there is not even complete agreement on the question of what patterns of inference of new trust relations are reasonable.

the provenance of trust information. Both of these are provided by several existing argumentation systems.

Some approaches to argumentation, for example abstract approaches such as that of Dung [3] and its derivatives, treat arguments as atomic objects. As a result, they say little or nothing about the internal structure of the argument and have no mechanism to represent the source of the information from which the argument is constructed. Such systems can represent the relationship between arguments (“ $a$  attacks  $b$ ”, and “ $b$  attacks  $c$ ”), but cannot represent *why* this is the case. As a result, such systems cannot capture the fact that  $a$  attacks  $b$  because  $b$  is based on information from source  $s$ , and there is evidence that source  $s$  is not trustworthy.

There are, however, a number of existing systems that extend [3] with more detailed information about the argument. One system is that of Amgoud [1], where an argument is taken to be a pair  $(H, h)$ ,  $h$  being a formula, the *conclusion* of the argument, and  $H$  being a set of formulae known as the *grounds* or *support* of the argument. Conclusion and support are related. In particular, [1] requires that  $H$  be a minimal consistent set of formulae such that  $H \vdash h$  in the language in which  $h$  and  $H$  are expressed. This means of representing the support is rather restricted. It presents the support as a bag of formulae with no indication as to how they are used in the construction of the argument, and without recording any of the intermediate steps. It is easy enough to see if another argument *rebuts*  $(H, h)$ , meaning that the conclusion of this second argument is the negation of  $h$ , and it is also quite simple to establish if the conclusion of the second argument contradicts any of the grounds in  $H$  (which in some systems of argumentation is known as *undercutting*). However, other forms of relationship are harder to establish. For example, in some cases it is interesting to know if an argument contradicts any of the intermediate steps in the chain of inferences between  $H$  and  $h$ .

Since the information about the steps in the argument can be useful, some systems of argumentation, for example [5] and [23], record more detail about the proof of  $h$  from  $H$  as part of the grounds. Some, including the system [20] which we will discuss in more detail below, go as far as to record the proof rules used in deriving  $h$  from  $H$ , permitting the notion of “attack” to include not only the intermediate conclusions but also the means by which they were derived.

Another problem with Dung’s argumentation system from the perspective of reasoning about trust is that it has no explicit means to represent degrees of trust. In [3] the important question is whether, given all the arguments that are known, a specific argument should be considered to hold. While one could construct a system for reasoning about trust in this way — the critical point, after all, is often whether someone’s argument is trustworthy or not — the prevalence of numerical measures of trust in the literature leads us to want to represent these.

Systems like that of Amgoud [1] provide one means of handling such measures, allowing formulae to have preference values attached to them. The values propagate to arguments and are taken into consideration when reasoning about the relationship between arguments (roughly speaking, strong arguments shrug off the attacks of weaker arguments). This approach seems a little too restrictive for dealing with trust, but there are systems that are more flexible. One example is the work of Oren et al. [22], which allows formulae and arguments to be weighted with the belief values used by Jøsang’s

subjective logic [13]. A more abstract approach is that of Fox [17] where values to represent belief in formulae are picked from some suitable *dictionary* of values, and propagated in a suitable way through the proof rules that are used to construct arguments. Arguments are then triples of conclusion, support, and value, and such systems are close to the notion of a *labelled deductive system* [2] (though they pre-date labelled deductive systems by some years).

## 4.2 A Suitable Argumentation System

Having given a high level description of how argumentation can help in handling a number of the aspects of reasoning about trust, we give a more detailed example of using a specific system of argumentation. The system we describe is the system  $TL$  that we introduced in [20], notable because it explicitly represents the rules of inference employed in constructing arguments in the support of the argument (which then makes it possible to dispute the application of those rules).

We start with a set of atomic propositions including  $\top$  and  $\perp$ , the ever true and ever false propositions. The set of well-formed formulae (*wffs*), labeled  $\mathcal{L}$ , is comprised of the set of atomic propositions closed under the connectives  $\{\neg, \rightarrow, \wedge, \vee\}$ .  $\mathcal{L}$  may then be used to create a database  $\Delta$  whose elements are 4-tuples:

$$(\theta : G : R : \tilde{d})$$

in which each element  $\theta$  is a formulae,  $G$  is the derivation of that formula,  $R$  is the sequence of rules of inference used in the derivation, and  $\tilde{d}$  is a suitable measure.

In more detail,  $\theta$  is a *wff* from  $\mathcal{L}$ ,  $G = (\theta_0, \theta_1, \dots, \theta_{n-1})$  is an ordered sequence of *wffs*, with  $n \geq 1$ , and  $R = (\vdash_1, \vdash_2, \dots, \vdash_n)$  is an ordered sequence of inference rules, such that:

$$\theta_0 \vdash_1 \theta_1 \vdash_2 \theta_2 \dots \theta_{n-1} \vdash_n \theta$$

In other words, each element  $\theta_k \in G$  is derived from the preceding element  $\theta_{k-1}$  as a result of the application of the  $k$ -th rule of inference,  $\vdash_k$ , ( $k = 1, \dots, n-1$ ). The rules of inference in any such sequence may be non-distinct. Thus  $G$  and  $R$  together provide an explicit representation of the way that  $\theta$  was inferred.

The element  $\tilde{d} = (d_1, d_2, \dots, d_n)$  is an ordered sequence of elements from some dictionary  $\mathcal{D}$ . For reasoning about trust, these elements could be a numerical measure of trust, or some linguistic term that indicates the trust in the relevant inference, for example:

*{very reliable, reliable, no opinion, somewhat unreliable, very unreliable}*

We also permit *wffs*  $\theta \in \mathcal{L}$  to be elements of  $\Delta$ , by including tuples of the form  $(\theta : \emptyset : \emptyset : \emptyset)$ , where each  $\emptyset$  indicates a null term. (Such tuples represent information that has not been derived — basic premises may take this form.) Note that the assignment of labels may be context-dependent, i.e., the  $d_i$  assigned to  $\vdash_i$  may also depend on  $\theta_{i-1}$ . This is the case for statistical inference, where the  $p$ -value depends on characteristics of the sample from which the inference is made, such as its size.

With this formal system, we can take a database  $\Delta$  and use the consequence relation  $\vdash_{TCR}$  defined in Figure 1 to build arguments for propositions of interest. This

$$\begin{array}{c}
\text{Ax} \frac{(\theta : G : R : \tilde{d}) \in \Delta}{\Delta \vdash_{TCR} (\theta : G : R : \tilde{d})} \\
\wedge\text{-I} \frac{\Delta \vdash_{TCR} (\theta : G : R : \tilde{d}) \text{ and } \Delta \vdash_{TCR} (\phi : H : S : \tilde{e})}{\Delta \vdash_{TCR} (\theta \wedge \phi : G \otimes H \otimes (\theta \wedge \phi) : R \otimes S \otimes (\vdash_{\wedge\text{-I}}) : \tilde{d} \otimes \tilde{e} \otimes (d_{\wedge\text{-I}}))} \\
\wedge\text{-E1} \frac{\Delta \vdash_{TCR} (\theta \wedge \phi : G : R : \tilde{d})}{\Delta \vdash_{TCR} (\theta : G \otimes (\theta) : R \otimes (\vdash_{\wedge\text{-E1}}) : \tilde{d} \otimes (d_{\wedge\text{-E1}}))} \\
\wedge\text{-E2} \frac{\Delta \vdash_{TCR} (\theta \wedge \phi : G : R : \tilde{d})}{\Delta \vdash_{TCR} (\phi : G \otimes (\phi) : R \otimes (\vdash_{\wedge\text{-E2}}) : \tilde{d} \otimes (d_{\wedge\text{-E2}}))} \\
\vee\text{-I1} \frac{\Delta \vdash_{TCR} (\theta : G : R : \tilde{d})}{\Delta \vdash_{TCR} (\theta \vee \phi : G \otimes (\theta \vee \phi) : R \otimes (\vdash_{\vee\text{-I1}}) : \tilde{d} \otimes (d_{\vee\text{-I1}}))} \\
\vee\text{-I2} \frac{\Delta \vdash_{TCR} (\phi : H : S : \tilde{e})}{\Delta \vdash_{TCR} (\theta \vee \phi : H \otimes (\theta \vee \phi) : S \otimes (\vdash_{\vee\text{-I2}}) : \tilde{e} \otimes (e_{\vee\text{-I2}}))} \\
\vee\text{-E} \frac{\Delta \vdash_{TCR} (\theta \vee \phi : G : R : \tilde{d}) \text{ and } \Delta, (\theta : \emptyset : \emptyset : \emptyset) \vdash_{TCR} (\gamma : H : S : \tilde{e}) \text{ and } \Delta, (\phi : \emptyset : \emptyset : \emptyset) \vdash_{TCR} (\gamma : J : T : \tilde{f})}{\Delta \vdash_{TCR} (\gamma : G \otimes H \otimes J \otimes (\gamma) : R \otimes S \otimes T \otimes (\vdash_{\vee\text{-E}}) : \tilde{d} \otimes \tilde{e} \otimes \tilde{f} \otimes (d_{\vee\text{-E}}))} \\
\neg\text{-I} \frac{\Delta, (\theta : \emptyset : \emptyset : \emptyset) \vdash_{TCR} (\perp : G : R : \tilde{d})}{\Delta \vdash_{TCR} (\neg\theta : G \otimes (\neg\theta) : R \otimes (\vdash_{\neg\text{-I}}) : \tilde{d} \otimes (d_{\neg\text{-I}}))} \\
\neg\text{-E} \frac{\Delta \vdash_{TCR} (\theta : G : R : \tilde{d}) \text{ and } \Delta \vdash_{TCR} (\neg\theta : H : S : \tilde{e})}{\Delta \vdash_{TCR} (\perp : G \otimes H \otimes (\perp) : R \otimes S \otimes (\vdash_{\neg\text{-E}}) : \tilde{d} \otimes \tilde{e} \otimes (d_{\neg\text{-E}}))} \\
\neg\neg\text{-E} \frac{\Delta \vdash_{TCR} (\neg\neg\theta : G : R : \tilde{d})}{\Delta \vdash_{TCR} (\theta : G \otimes (\theta) : R \otimes (\vdash_{\neg\neg\text{-E}}) : \tilde{d} \otimes (d_{\neg\neg\text{-E}}))} \\
\rightarrow\text{-I} \frac{\Delta, (\theta : \emptyset : \emptyset : \emptyset) \vdash_{TCR} (\phi : G : R : \tilde{d})}{\Delta \vdash_{TCR} (\theta \rightarrow \phi : G \otimes (\theta \rightarrow \phi) : R \otimes (\vdash_{\rightarrow\text{-I}}) : \tilde{d} \otimes (d_{\rightarrow\text{-I}}))} \\
\rightarrow\text{-E} \frac{\Delta \vdash_{TCR} (\theta : G : R : \tilde{d}) \text{ and } \Delta \vdash_{TCR} (\theta \rightarrow \phi : H : S : \tilde{e})}{\Delta \vdash_{TCR} (\phi : G \otimes H \otimes (\phi) : R \otimes S \otimes (\vdash_{\rightarrow\text{-E}}) : \tilde{d} \otimes \tilde{e} \otimes (d_{\rightarrow\text{-E}}))}
\end{array}$$

**Fig. 1.** The TL Consequence Relation



consequence relation is defined in terms of rules for building new arguments from old. The rules are written in a style similar to standard Gentzen proof rules, with the antecedents of the rule above the horizontal line and the consequent below. In Figure 11 we use the notation  $G \otimes H$  to refer to that ordered sequence created from appending the elements of sequence  $H$  after the elements of sequence  $G$ , each in their respective order. The rules are as follows:

Ax The rule Ax says that if the tuple  $(\theta : G : R : \tilde{d})$  is in the database, then it is possible to build the argument  $(\theta : G : R : \tilde{d})$  from the database. The rule thus allows the construction of arguments from database items.

$\wedge$ -I The rule  $\wedge$ -I says that if the arguments  $(\theta : G : R : \tilde{d})$  and  $(\phi : H : S : \tilde{e})$  may be built from the database, then an argument for  $\theta \wedge \phi$  may also be built. The rule thus shows how to introduce arguments about conjunctions; using it requires an inference of the form:  $\theta, \phi \vdash (\theta \wedge \phi)$ , which we denote

$$\vdash_{\wedge\text{-I}}$$

in Figure 11. This inference is then assigned a value of  $d_{\wedge\text{-I}}$ .

$\wedge$ -E The rule  $\wedge$ -E1 says that if it is possible to build an argument for  $\theta \wedge \phi$  from the database, then it is also possible to build an argument for  $\theta$ . Thus the rule allows the elimination of one conjunct from an argument, and its use requires an inference of the form:  $\theta \wedge \phi \vdash \theta$ .  $\wedge$ -E2 allows the elimination of the other disjunct.

$\vee$ -I The rule  $\vee$ -I1 allows the introduction of a disjunction from the left disjunct and the rule  $\vee$ -I2 allows the introduction of a disjunction from the right disjunct.

$\vee$ -E The rule  $\vee$ -E allows the elimination of a disjunction and its replacement by tuple when that tuple is a TL-consequence of each disjunct.

$\neg$ -I The rule  $\neg$ -I allows the introduction of negation.

$\neg$ -E The rule  $\neg$ -E allows the derivation of  $\perp$ , the ever-false proposition, from a contradiction.

$\neg\neg$ -E The rule  $\neg\neg$ -E allows the elimination of a double negation, and thus permits the assertion of the Law of the Excluded Middle (LEM).

$\rightarrow$ -I The rule  $\rightarrow$ -I says that if on adding a tuple  $(\theta : \emptyset : \emptyset : \emptyset)$  to a database, where  $\theta \in \mathcal{L}$ , it is possible to conclude  $\phi$ , then there is an argument for  $\theta \rightarrow \phi$ . The rule thus allows the introduction of  $\rightarrow$  into arguments.

$\rightarrow$ -E The rule  $\rightarrow$ -E says that from an argument for  $\theta$  and an argument for  $\theta \rightarrow \phi$  it is possible to build an argument for  $\phi$ . The rule thus allows the elimination of  $\rightarrow$  from arguments and is analogous to MP in standard propositional logic.

This is an intentionally abstract formalism — syntactically complete, but without a specified semantics. The idea is that to capture a specific domain, we have to identify a suitable dictionary from which to construct the  $\tilde{d}$  and that this set of values will determine the mechanism by which we can compute an overall value from the sequence of  $d_i$ . For example, if one wanted to use Jøsang's subjective logic, then the mechanism for combining the  $d_i$ 's would be taken from [13]. If one wanted to quantify trust using probability, then the combination rules would be those dictated by probability theory

(for example using [29]). If one wanted to use the dictionary mentioned above (“very reliable” and so on) then it would be necessary to determine the right way to combine these values across all the inference rules in Figure 1.

Even without specifying these mechanisms, it should be clear that whatever means we use to quantify trust in combination with  $TL$ , the formalism can both capture trust values and the precise source of information used. It is also possible to go further. The fact that  $TL$  includes explicit reference to different forms of inference allows us to capture the fact that inferences may differ depending on the source of the information on which they are based — we might want to make different inferences depending on whether the source was something we have direct experience of, or something that comes from a trusted source, or, indeed, something that comes from an untrusted source.

### 4.3 Extensions

The previous sections have argued that systems of argumentation can provide the core functionality required to reason about trust. Here we discuss how systems of argumentation, especially the system  $TL$  sketched above, can provide additional mechanisms that are important in dealing with trust.

First, argumentation systems explicitly allow the representation of different points of view. The system  $TL$  we have sketched above provides us with the rules for constructing arguments, and it does not limit the number of arguments that one can construct for a specific conclusion. Thus, the database  $\Delta$  may contain information that represents a number of different assessments of the trustworthiness of, for example, a source of information. This might be done through the inclusion of a number of tuples  $(\theta : G : R : \tilde{d})$  with different  $G$ s, representing different views of the sources, and different  $\tilde{d}$ s representing different assessments of trustworthiness. These pieces of information could then be used to make different inferences, with any potential choice between conclusions being made on the basis of the relevant  $\tilde{d}$  values.

That is one, fairly simple, way to represent different viewpoints. Another would be to have different argumentation systems represent the views of different individuals, and to use the mechanisms of argumentation-based dialogue (like those discussed in [25,8]) to explore the differences in the views of trust and to attempt to resolve them. In such a combination, the individual argumentation systems can be constructed using  $TL$ , and would then reason about trust based on a single viewpoint. The interaction between different viewpoints is then captured by the dialogue mechanisms of [25,8], enabling a rational discourse about trust issues.

Another important aspect of reasoning about trust, addressed in [10] for example, is the need for an individual to be able to revise the trust they have in another based on experience. Revision of beliefs is not a subject that has been widely considered within the argumentation community, but [24] suggests some approaches to the subject. We plan to examine how these can be implemented on top of  $TL$  giving us a means to represent the case in which one individual revises its view of a source as a result of considering information provided by another individual. In addition, [19] looks at how to use statistics on past performance to build arguments about trust, and combining this work with  $TL$  is a way to augment it with the ability to infer numerical degrees of trust.

## 5 Conclusion

This paper has presented the case for using argumentation as a mechanism for reasoning about trust. Starting from some of the many views of trust expressed in the literature, we extracted the major features that need to be represented, discussed formalisms for handling trust, and then suggested how argumentation could be used for reasoning about trust. We sketched in some detail how a specific system of argumentation, *TL*, could be used in this way and identified some additional argumentation-based mechanisms that could be of use in dealing with trust.

## Acknowledgement

Research was sponsored by the Army Research Laboratory and was accomplished under Cooperative Agreement Number W911NF-09-2-0053. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

## References

1. Amgoud, L., Cayrol, C.: A reasoning model based on the production of acceptable arguments. *Annals of Mathematics and Artificial Intelligence* 34(1-3), 197–215 (2002)
2. Chesñevar, C., Simari, G.: Modelling inference in argumentation through labelled deduction: Formalization and logical properties. *Logica Universalis* 1(1), 93–124 (2007)
3. Dung, P.M.: On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence* 77(2), 321–358 (1995)
4. Gambetta, D.: Can we trust them? In: Gambetta, D. (ed.) *Trust: Making and breaking cooperative relations*, pp. 213–238. Blackwell, Oxford (1990)
5. Garcia, A.J., Simari, G.R.: Defeasible logic programming: an argumentative approach. *Theory and Practice of Logic Programming* 4(2), 95–138 (2004)
6. Golbeck, J.: Combining provenance with trust in social networks for semantic web content filtering. In: *Proceedings of the International Provenance and Annotation Workshop*, Chicago, Illinois (May 2006)
7. Golbeck, J., Halaschek-Wiener, C.: Trust-based revision for expressive web syndication. *The Logic Journal of the IGPL* (to appear)
8. Gordon, T.F.: The pleadings game: An exercise in computational dialectics. *Artificial Intelligence and Law* 2(4), 239–292 (1994)
9. Grandison, T., Sloman, M.: A survey of trust in internet applications. *IEEE Communications Surveys and Tutorials* 4(4), 2–16 (2000)
10. Hang, C.-W., Wang, Y., Singh, M.P.: An adaptive probabilistic trust model and its evaluation. In: *Proceedings of the 7th International Conference on Autonomous Agents and Multiagent Systems* (2008)
11. Jaynes, E.T.: *Probability Theory: The Logic of Science*. Cambridge University Press, Cambridge (2003)

12. Jøsang, A., Gray, E., Kinatered, M.: Simplification and analysis of transitive trust networks. *Web Intelligence and Agent Systems* 4(2), 139–161 (2006)
13. Jøsang, A., Hayward, R., Pope, S.: Trust network analysis with subjective logic. In: *Proceedings of the 29th Australasian Computer Society Conference* (January 2006)
14. Jøsang, A., Ismail, R., Boyd, C.: A survey of trust and reputation systems for online service provision. *Decision Support Systems* 43(2), 618–644 (2007)
15. Kamvar, S.D., Schlosser, M.T., Garcia-Molina, H.: The Eigentrust algorithm for reputation management in P2P networks. In: *Proceedings of the 12th World Wide Web Conference* (May 2004)
16. Katz, Y., Golbeck, J.: Social network-based trust in prioritized default logic. In: *Proceedings of the 21st National Conference on Artificial Intelligence* (2006)
17. Krause, P., Ambler, S., Elvang-Gøransson, M., Fox, J.: A logic of argumentation for reasoning under uncertainty. *Computational Intelligence* 11(1), 113–131 (1995)
18. Kuter, Y., Golbeck, J.: Sunny: A new algorithm for trust inference in social networks using probabilistic confidence models. In: *Proceedings of the 22nd National Conference on Artificial Intelligence* (2007)
19. Matt, P.-A., Morge, M., Toni, F.: Combining statistics and arguments to compute trust. In: van der Hoek, W., Kaminka, G., Lespérance, Y., Luck, M., Sen, S. (eds.) *Proceedings of the 9th International Conference on Autonomous Agents and Multiagents Systems*, Toronto, Canada, pp. 209–216 (May 2010)
20. McBurney, P., Parsons, S.: Tenacious tortoises: A formalism for argument over rules of inference. In: *Proceedings of the ECAI Workshop on Computational Dialectics* (2000)
21. McKnight, D.H., Chervany, N.L.: The meanings of trust. Working Paper 96-04, Carlson School of Management, University of Minnesota (1996)
22. Oren, N., Norman, T., Preece, A.: Subjective logic and arguing with evidence. *Artificial Intelligence* 171(10-15), 838–854 (2007)
23. Parsons, S., Sierra, C., Jennings, N.R.: Agents that reason and negotiate by arguing. *Journal of Logic and Computation* 8(3), 261–292 (1998)
24. Parsons, S., Sklar, E.: How agents alter their beliefs after an argumentation-based dialogue. In: Parsons, S., Maudet, N., Moraitis, P., Rahwan, I. (eds.) *ArgMAS 2005. LNCS (LNAI)*, vol. 4049, pp. 297–312. Springer, Heidelberg (2006)
25. Prakken, H.: Coherence and flexibility in dialogue games for argumentation. *Journal of Logic and Computation* 15(6), 1009–1040 (2005)
26. Shafer, G.: *A Mathematical Theory of Evidence*. Princeton University Press, Princeton (1976)
27. Sztompka, P.: *Trust: A Sociological Theory*. Cambridge University Press, Cambridge (1999)
28. Wang, Y., Singh, M.P.: Trust representation and aggregation in a distributed agent system. In: *Proceedings of the 21st National Conference on Artificial Intelligence* (2006)
29. Xiang, Y., Jia, N.: Modeling causal reinforcement and undermining for CPT elicitation. *IEEE Transactions on Knowledge and Data Engineering* 19(12), 1708–1718 (2007)

# An Argument-Based Multi-agent System for Information Integration

Marcela Capobianco and Guillermo R. Simari

Artificial Intelligence Research and Development Laboratory  
Department of Computer Science and Engineering  
Universidad Nacional del Sur – Av. Alem 1253, (8000) Bahía Blanca ARGENTINA  
{mc,grs}@cs.uns.edu.ar

**Abstract.** In this paper we address the problem of obtaining a consolidated view of the knowledge that a community of information agents possesses in the form of private, possibly large, databases. Each agent in the community has independent sources of information and each database could contain information that is potentially inconsistent and incomplete, both by itself and/or in conjunction with some of the others. These characteristics make the consolidation difficult by traditional means. The idea of obtaining a single view is to provide a way of querying the resulting knowledge in a skeptical manner, *i.e.*, receiving one answer that reflects the perception of the information community.

Agents using the proposed system will be able to access multiple sources of knowledge represented in the form of deductive databases as if they were accessing a single one. One application of this schema is a novel architecture for decision-support systems (DSS) that will combine database technologies, specifically federated databases, which will cast as information agents, with an argumentation-based framework.

**Categories and Subjects Descriptors:** I.2.4 [Artificial Intelligence]: Knowledge Representation Formalisms and Methods—Representation languages, Representations (procedural and rule-based); H.1.0 [Models and Principles]: Systems and Information Theory—General systems theory.

**General Terms:** Algorithms, Design, Performance.

**Keywords:** Argumentation, Knowledge representation, Design languages for agent systems.

## 1 Introduction

Information systems, and the capability to obtain answers from them, play a key role in our society. In particular, data intensive applications are in constant demand and there is need for computing environments with much more intelligent capabilities than those present in today's Data-base Management Systems (DBMS). The expected requirements for these systems change every day: they constantly become more complex and more advanced features are demanded from them.

In this paper we address the problem of obtaining a consolidated view of the knowledge that a community of information agents possess in the form of private, possibly large, databases. Each agent in the community has independent sources of information and each database could contain information that is potentially inconsistent and incomplete, both by itself and/or in conjunction with some of the others. These characteristics make the consolidation difficult by traditional means. The idea of obtaining a single view is to provide a way of querying the resulting knowledge in a skeptical manner, *i.e.*, receiving one answer that reflects the perception of the information community.

Agents using the proposed system will be able to access multiple sources of knowledge represented in the form of deductive databases as if they were accessing a single one. One application of this schema is a novel architecture for decision-support systems (DSS) that will combine database technologies, specifically federated databases, which we will cast as information agents, with an argumentation-based framework.

Recently, there has been much progress in developing efficient techniques to store and retrieve data, and many satisfactory solutions have been found for the associated problems. However it remains an open problem how to understand and interpret large amounts of information. To do this we need specific formalisms that can perform complicated inferences, obtain the appropriate conclusions, and justify their results. We claim that these formalisms should also be able to access seamlessly databases distributed over a network.

In the field of deductive databases there has been a continued effort to produce an answer to this problem. Deductive databases store explicit and implicit information; explicit information is stored in the manner of a relational database and implicit information is recorded in the form of rules that enable inferences based on the stored data. These systems combine techniques and tools from relational databases with rule based formalisms. Hence, they are capable of handling large amounts of information and perform some sort of reasoning based on it. Nevertheless, these systems have certain limitations and shortcomings for knowledge representation and modeling commonsense reasoning, especially for managing incomplete and potentially contradictory information, as argued by several authors [17,23,16].

Argumentation frameworks [9,19,14] are an excellent starting point for building intelligent systems with interesting reasoning abilities. Research in argumentation has provided important results while striving to obtain tools for commonsense reasoning, and this prompted a new set of argument-based applications in diverse areas where knowledge representation issues play a major role [10,5,7].

We believe that deductive databases can be combined with argumentation formalisms to obtain interactive systems able to reason with large databases, even in the presence of incomplete and potentially contradictory information. This can be a significant advantage with respect to systems based on logic programming,

such as datalog, that cannot deal with contradictory information.<sup>1</sup> In particular, this could be useful in contexts where information is obtained from multiple databases, and these databases may be contradictory among themselves.

The multi-agent system introduced here virtually integrates different databases into a common view; in that manner users of this system can query multiple databases as if they were a single one. This schema can be applied to obtain a novel system architecture for decision-support systems (DSS) that combines database technologies, specifically federated databases [18], with an argumentation based framework.

In our proposal we consider that information is obtained from a number of different heterogeneous database systems, each represented by a particular agent. The reasoning mechanisms, based on an argumentative engine, use this information to construct a consolidated global view of the database. This task is performed by the reasoning agent, that is based on a set of rules expressed in a particular argumentation framework. This agent can deal with incomplete and contradictory information and can also be personalized for any particular DSS in a relatively simple way.

We have also considered that one of the design objectives of interactive systems is that they can respond in a timely manner to users' queries. So far the main objection to the use of argumentation in interactive systems is their computational complexity. In previous work [6] we have addressed the issue of optimizing argumentation systems, where the optimization technique consisted in using a precompiled knowledge component as a tool to allow significant speed-ups in the inference process. We also apply this technique in our reasoning agent.

To understand the advantages of the proposed reasoning mechanism used in our multiagent system, consider a set of databases used by the employers responsible of drug administration, sales, and delivery in a given hospital. These databases contains information regarding drugs, patients, known allergies, and addictions. Suppose a deductive database system in the style of datalog is used to query this information to derive certain relations. In this setting, there is a rule establishing that a drug should be given to a patient if the patient has a prescription for this drug signed by a physician. There could also be a rule saying that the drug should not be sold if the prescription is signed by the patient. In this case, if Dr. Gregory House enters the clinic with a prescription signed by himself to get Vicodin, the employers could query the system to see if the drug should or should not be sold. If a traditional deductive database is used, in the style of datalog or another logic programming based system, this would give raise to a contradiction and the system would not be able to recommend a course of action. Using our system, an argument can be built backing that the medication should be sold, given that there is a prescription signed by a doctor that warrants it. However, an argument for not selling the drug could also be built considering that the doctor and the patient are the same person. Our

---

<sup>1</sup> Some extensions of datalog handle negation using CWA (see [8]), but these approaches do not allow negation in the head of the rules in the style of extended logic programming.

argument-based framework can then compare both arguments, decide that the second argument is preferred, and answer the query saying that the drug should no be sold. In addition, it can also explain the reasons that back its answer. Note that this kind of behavior cannot be obtained in Datalog-like systems.

The rest of this article is organized as follows. Section 2 sums up related work, Section 3 contains our proposal for the system architecture, and section 4 formally describes the reasoning module, a key component of this architecture. Section 5 presents a key optimization for the argumentation-based reasoning process, and Section 6 shows a realistic example of the system's mechanics. Finally, Section 7 states the conclusions.

## 2 Integrating DBMS and Reasoning Systems

Previous work on the integration of databases and reasoning systems has almost been restricted to coupling Prolog interpreters and relational databases. These approaches were motivated in the declarative nature of logic programming languages and the data-handling capabilities of database systems. Several researchers have built intelligent database systems coupling Prolog and a relational DBMS or extending Prolog with database facilities [8]. These works were motivated by the fact that Prolog attracted attention in the 80's for its ability to support rapid development of complex applications. Besides, Prolog is based on Horn Clauses that are close relatives of database query languages and its language is more powerful than SQL [24].

Brodie and Jarke [4] envisioned several years ago that large scale data processing would require more efficient and more intelligent access to databases. He proposed the integration of logic programming and databases to meet future requirements. First, he identified two different approaches for coupling a Prolog interpreter and a Relational DBMS, which are usually called "tight coupling" and "loose coupling". In the tight coupling approach the Prolog interpreter and the Relational DBMS are strongly integrated. For example, the Prolog interpreter can directly use low level functionalities of the DBMS, like relation management in secondary memory, and relation access via indexes [12]. In contrast, in the loose coupling approach, the Relational DBMS is called by the Prolog interpreter at the top level, that acts like a standard user. It sends Relational queries to the DBMS, and the corresponding answers are treated as ground clauses by the interpreter.

Brodie and Jarke also identified four basic architectural frameworks for combining Prolog and a database system:

- Loose coupling of an existing Prolog implementation to an existing relational database system;
- Extending Prolog to include some facilities of the relational database system;
- Extending an existing relational database to include some features of Prolog;
- Tightly integrating logic programming techniques with those of relational database systems.



They recommend the fourth alternative (tight integration), based on the belief that a special purposed language for large scale knowledge base systems would best address issues regarding performance, knowledge representation and software engineering. They also put forward a number of issues concerning the best division of tasks between logic programming and a DBMS.

Zaniolo [25] proposed an approach to intelligent databases based on deductive databases. He advocates for elevating database programming to the more abstract level of rules and knowledge base programming to create an environment more supportive of the new wave of database applications. To achieve these goals the *LDL/DDL+* project was developed. During the project a new logic-based language was designed along with the definition of its formal semantics, new implementation techniques were developed for the efficient support of rule-based logic languages, and it was successfully used in a wide range of application domains. The system supported an open architecture and SQL schema from external databases could be incorporated into the *LDL* program seamlessly.

In the following section we present the system architecture for our proposal. We believe that argumentation can offer a new perspective into the problem of reasoning with large databases, giving more expressive power to the reasoning component, making it able to decide even in the presence of uncertain and/or contradictory information. This addresses a limitation that was present in each of the deductive database systems considered in this section.

### 3 System Architecture

In this section we present an architectural pattern for our multiagent system that can be applied to design information-based applications where a certain level of intelligence is required. Such applications will be engineered for contexts where: (1) information is uncertain and heterogeneous, (2) handling of great volume of data flows is needed, and (3) data may be incomplete, vague, or contradictory. These applications are also expected to integrate multiple information systems such as databases, knowledge bases, source systems, etc.

Our system architecture is presented in Figure 1. The architecture is modular and is independent of any particular domain or application. We have used a layered architectural style, where every layer provides a series of services to the one above. The first of our layers concerns data and knowledge acquisition. This layer will receive heterogeneous sources of data and will extract and transform this data into the formats required of the particular application. It can work with diverse sources, such as laboratory data, different types of sensors, knowledge bases, etc.

The received data will be formatted to comply with the relational models provided by a group of federated databases that share a common export schema. In our system, each one of the databases is represented by an agent. The common export schema will be the result of a negotiation process among these agents. The union of the views of these databases will generate a key element of our framework, the extensional database that contains the information needed for the

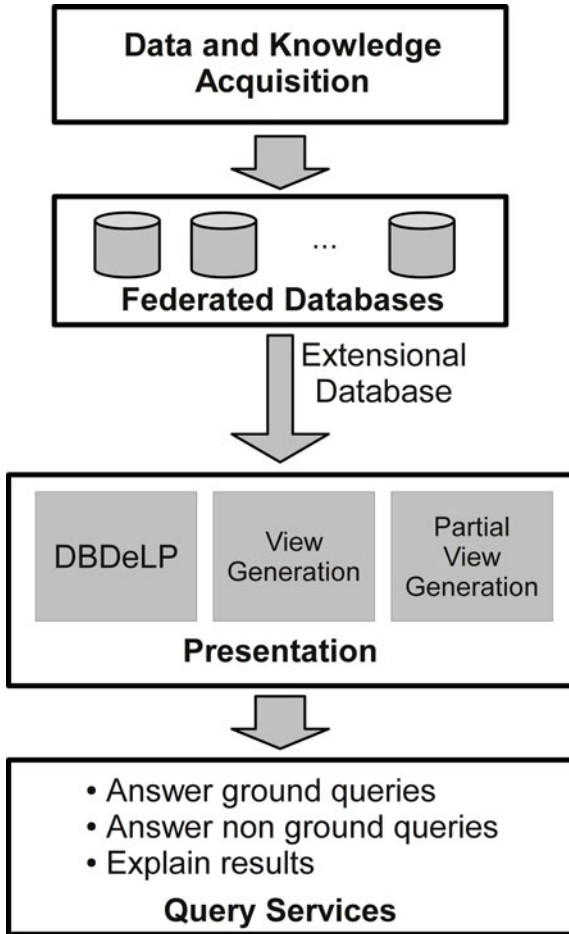


Fig. 1. Proposed architectural pattern

reasoning module. The extensional database also will provide the data elements with a certainty degree that depends of the credibility of the data source from where it was obtained.

We have chosen to use a multi-source perspective for the characterization of data quality [3]. In this case, the quality of data can be evaluated by comparison with the quality of other homologous data (i.e., data from different information sources which represent the same reality but may have contradictory values). The approaches usually adopted to reconcile heterogeneity between values of data are: (1) to prefer the values of the most reliable sources, (2) to mention the source ID for each value, or (3) to store quality meta-data with the data.

For our proposed architecture, we have used the second approach. In multi-source databases, each attribute of a multiple source element has multiple values with the ID of their source and their associated *Quality of Expertise*, which is

represented as meta-data associated with each value, such as a given certainty degree. This degree may be obtained weighting the plausibility of the data value, its accuracy, the credibility of its source, and the freshness of the data.

The federated database layer provides the extensional database to the presentation layer. The extensional database can be computed on demand and is not necessarily stored in a physical component. The presentation layer contains the services related with the reasoning process and its optimization. This is the core of our proposal and will be described later on in Sections 4 and 5. The reasoning agent that generated the consolidated view is part of this layer. It contains the set of rules that encode the specific knowledge of the application. These rules will be used by the argumentation-based inference engine. The presentation layer also commands the generation of the extensional database, and the selection if it is going to be done on demand (following a lazy approach) or if it has to be computed completely. It can also generate a partial view of the system according to these rules, resulting in an optimization mechanism. This partial view depends only on the set of rules and must be changed accordingly if changes on the rules are produced. Finally, the query services layer is composed by an interactive agent that receives user queries, provides answers, and can also explain the reasons backing these answers.

## 4 The DB\_DeLP Argumentation Framework

In this section we formally define the argumentation system that is used by the reasoning agent in the Query Services Layer of the proposed system architecture. Here we detail the semantics and proof theory of the framework and we also show some practical examples. A simplified view of our system would describe it as a deductive database whose inference engine is based on a specialization of the DeLP language [15]. This particular framework will be known as *Database Defeasible Logic Programming* (DB\_DeLP). Formally, DB\_DeLP is a language for knowledge representation and reasoning that uses *defeasible argumentation* to decide between contradictory conclusions through a *dialectical analysis*. DB\_DeLP also incorporates uncertainty management, taking elements from Possibilistic Defeasible Logic Programming (P\_DeLP) [21], an extension of DeLP in which the elements of the language have the form  $(\phi, \alpha)$ , where  $\phi$  is a DeLP clause or fact and  $\alpha$  expresses a lower bound for the certainty of  $\phi$  in terms of a necessity measure. Conceptually, our deductive database consists of an extensional database **EDB**, an intensional database **IDB**, and a set of integrity constraints **IC**. In what follows, we formally define these elements.

The language of DB\_DeLP follows a logic programming style. Standard logic programming concepts (such as signature, variables, functions, *etc.*) are defined in the usual way. Literals are atoms that may be preceded by the symbol “ $\sim$ ” denoting *strict* negation, as in extended logic programming.

**Definition 1.** [Literal–Weighted Literal] *Let  $\Sigma$  be a signature, then every atom  $A$  of  $\Sigma$  is a positive literal, while every negated atom  $\sim A$  is a negative literal. A literal of  $\Sigma$  is a positive literal or a negative literal. A certainty weighted literal,*

or simply a weighted literal, is a pair  $(L, \alpha)$  where  $L$  is a literal and  $\alpha \in [0, 1]$  expresses a lower bound for the certainty of  $L$  in terms of a necessity measure.

The extensional database **EDB** is composed by a set of certainty weighted literals, according to the export schema of the federated database that is part of our architecture. Conceptually, it accounts for the union of the views of every particular database that belongs to the federation [18]. When implementing the system, this set of ground literals may not be physically stored in any place, and may simply be obtained on demand when information about a particular literal is needed.

The certainty degree associated with every literal is assigned by the federated database layer that assigns a particular degree to every data source according to its plausibility. The resulting extensional database is not necessarily consistent, in the sense that a literal and its complement w.r.t. strong negation may both be present, with different or the same certainty degrees. In this case, the system decides according to a given criterion which fact will prevail and which one will be removed from the view.

The intensional part of a DB\_DeLP database is formed by a set of *defeasible rules* and *integrity constraints*. Defeasible rules provide a way of performing tentative reasoning as in other argumentation formalisms [9,19,14].

**Definition 2.** [Defeasible Rule] A defeasible rule expresses a tentative, weighted, relation between a literal  $L_0$  and a finite set of literals  $\{L_1, L_2, \dots, L_k\}$ . It has the form  $(L_0 \prec L_1, L_2, \dots, L_k, \alpha)$  where  $\alpha \in [0, 1]$  expresses a lower bound for the certainty of the rule in terms of a necessity measure.

In previously defined argumentation systems, the meaning of defeasible rules  $L_0 \prec L_1, L_2, \dots, L_k$  was understood as “ $L_1, L_2, \dots, L_k$  provide tentative reasons to believe in  $L_0$ ” [22], but these rules did not have an associated certainty degree. In contrast, DB\_DeLP adds the certainty degree, that expresses how strong is the connection between the premises and the conclusion. A defeasible rule with a certainty degree 1 will model a strong rule. Figures 2 and 3 show an extensional and an intensional database in our formalism.

Note that DB\_DeLP programs are *range-restricted*, a common condition for deductive databases: a program is said to be range-restricted if and only if every variable that appears in the head of the clause also appears in its body. This implies that all the facts in the program must be ground (cannot contain variables). These programs can be interpreted more efficiently since full unification

<u>species(X,Y)</u>	<u>age(X,Y)</u>
(species(simba,lion),0.6)	(age(simba,young),0.65)
(species(mufasa,lion),0.7)	(age(mufasa,old),0.7)
(species(grace,lion),0.6)	(age(grace,adult),0.8)
(species(grace,leopard),0.4)	(age(dumbo,baby),0.8)
...	...

**Fig. 2.** An Extensional Database in DB\_DeLP

```

(feline(X)  $\leftarrow$  species(X,lion),1)
(climbs_tree(X)  $\leftarrow$  feline(X),0.65)
( $\sim$ climbs_tree(X)  $\leftarrow$  species(X,lion),0.70)
(climbs_tree(X)  $\leftarrow$  species(X,lion), age(X,young),.0.75)
( $\sim$ climbs_tree(X)  $\leftarrow$  sick(X),0.45)

```

**Fig. 3.** An Intensional Database in DB\_DeLP

is not needed, only matching that is significantly more efficient. Nevertheless, the reason for this decision comes from a semantic standpoint, given that a defeasible reason in which there is no connection between the head and the body has no clear meaning; the range restriction ensures this connection.

*Integrity constraints* are rules of the form  $L \leftarrow L_0, L_1, \dots, L_n$  where  $L$  is a literal, and  $L_0, L_1, \dots, L_n$  is a non-empty finite set of literals. These rules are used to compute the *derived negations* as follows. For the extensional and intensional databases regarding felines, consider that the set of integrity constraints is composed by  $\{\sim\text{leopard}(X) \leftarrow \text{lion}(X), \sim\text{lion}(X) \leftarrow \text{leopard}(X)\}$  and the negations  $\{(\sim\text{species}(\text{grace}, \text{lion}), 0.4), (\sim\text{species}(\text{grace}, \text{leopard}), 0.6)\}$  are then added to the extensional database. The certainty degree of the added rule is calculated as the minimum of the certainty degree of the literals that are present in the body of the integrity constraint rule used to obtain it. Note that a conflict may arise with information received from other knowledge bases, since we may want to add a literal and its complement may be already present in the extensional database. Then the system will decide according to a given criterion which fact will prevail and which one will be removed from the view. A standard criterion in this case would be using the plausibility of the source, the certainty degree of the literals, or a combination of both. Databases in DB\_DeLP, for short called *defeasible databases*, can also include built-in predicates as needed along with their corresponding axioms.

The P\_DeLP language [11], which presented the novel idea of mixing argumentation and possibilistic logic, is based on Possibilistic Gödel Logic or PGL [21], which is able to model both uncertainty and fuzziness and allows for a partial matching mechanism between fuzzy propositional variables. In DB\_DeLP, for simplicity reasons, we will avoid fuzzy propositions, and hence it will be based on the necessity-valued classical Possibilistic logic [13]. As a consequence, possibilistic models are defined by possibility distributions on the set of classical interpretations, and the proof theory for our formulas, written  $\vdash$ , is defined by derivation based on the following instance of the Generalized Modus Ponens rule (GMP):  $(L_0 \leftarrow L_1 \wedge \dots \wedge L_k, \gamma), (L_1, \beta_1), \dots, (L_k, \beta_k) \vdash (L_0, \min(\gamma, \beta_1, \dots, \beta_k))$ , which is a particular instance of the well-known possibilistic resolution rule, and which provides the *non-fuzzy* fragment of DB\_DeLP with a complete calculus for determining the maximum degree of possibilistic entailment for weighted literals. Literals in the extensional database are the base case of the derivation sequence; for every literal  $Q$  in **EDB** with a certainty degree  $\alpha$  it holds that  $(Q, \alpha)$  can be derived from the corresponding program.

A query presented to a DB\_DeLP database is a ground literal  $Q$  which must be supported by an *argument*. Deduction in DB\_DeLP is argumentation-based, thus a derivation is not enough to endorse a particular fact, and queries must be supported by arguments. In the following definition  $instances(IGB)$  accounts for any set of ground instances of the rules in **IGB**, replacing free variables for ground literals in the usual way.

**Definition 3.** [Argument]–[Subargument] *Let  $DB = (EDB, IGB, IC)$  be a de-feasible database,  $\mathcal{A} \subseteq instances(IGB)$  is an argument for a goal  $Q$  with necessity degree  $\alpha > 0$ , denoted as  $\langle \mathcal{A}, Q, \alpha \rangle$ , iff:*

1.  $\Psi \cup \mathcal{A} \vdash (Q, \alpha)$ ,
2.  $\Psi \cup \mathcal{A}$  is non contradictory, and
3. there is no  $\mathcal{A}_1 \subset \mathcal{A}$  such that  $\Psi \cup \mathcal{A}_1 \vdash (Q, \beta)$ ,  $\beta > 0$ .

An argument  $\langle \mathcal{A}, Q, \alpha \rangle$  is a subargument of  $\langle \mathcal{B}, R, \beta \rangle$  iff  $\mathcal{A} \subseteq \mathcal{B}$ .

Arguments in DB\_DeLP can attack each other; this situation is captured by the notion of *counterargument*. An argument  $\langle \mathcal{A}_1, Q_1, \alpha \rangle$  *counter-argues* an argument  $\langle \mathcal{A}_2, Q_2, \beta \rangle$  at a literal  $Q$  if and only if there is a sub-argument  $\langle \mathcal{A}, Q, \gamma \rangle$  of  $\langle \mathcal{A}_2, Q_2, \beta \rangle$ , (called *disagreement subargument*), such that  $Q_1$  and  $Q$  are complementary literals. Defeat among arguments is defined combining the counter-argument relation and a preference criterion “ $\succeq$ ”. This criterion is defined on the basis of the necessity measures associated with arguments.

**Definition 4.** [Preference criterion  $\succeq$ ] [11] *Let  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$  be a counterargument for  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$ . We will say that  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$  is preferred over  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$  (denoted  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle \succeq \langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$ ) iff  $\alpha_1 \geq \alpha_2$ . If it is the case that  $\alpha_1 > \alpha_2$ , then we will say that  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$  is strictly preferred over  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$ , denoted  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle \succ \langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$ . Otherwise, if  $\alpha_1 = \alpha_2$  we will say that both arguments are equi-preferred, denoted  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle \approx \langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$ .*

**Definition 5.** [Defeat] [11] *Let  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$  and  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$  be two arguments built from a program  $\mathcal{P}$ . Then  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$  defeats  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$  (or equivalently  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$  is a defeater for  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$ ) iff (1) Argument  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$  counter-argues argument  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$  with disagreement subargument  $\langle \mathcal{A}, Q, \alpha \rangle$ ; and (2) Either it is true that  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle \succ \langle \mathcal{A}, Q, \alpha \rangle$ , in which case  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$  will be called a proper defeater for  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$ , or  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle \approx \langle \mathcal{A}, Q, \alpha \rangle$ , in which case  $\langle \mathcal{A}_1, Q_1, \alpha_1 \rangle$  will be called a blocking defeater for  $\langle \mathcal{A}_2, Q_2, \alpha_2 \rangle$ .*

As in most argumentation systems [9,20], DB\_DeLP relies on an exhaustive dialectical analysis which allows to determine if a given argument is *ultimately* undefeated (or *warranted*) w.r.t. a program  $\mathcal{P}$ . An *argumentation line* starting with an argument  $\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle$  is a sequence  $[\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle, \langle \mathcal{A}_1, Q_1, \alpha_1 \rangle, \dots, \langle \mathcal{A}_n, Q_n, \alpha_n \rangle, \dots]$  that can be thought of as an exchange of arguments between two parties, a *proponent* (even-numbered arguments) and an *opponent* (odd-numbered arguments).

Given a program  $\mathcal{P}$  and an argument  $\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle$ , the set of all acceptable argumentation lines starting with  $\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle$  accounts for a whole dialectical

analysis for  $\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle$  (i.e. all possible dialogs rooted in  $\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle$ ), formalized as a *dialectical tree* and denoted  $\mathcal{T}_{\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle}$ . Nodes in a dialectical tree  $\mathcal{T}_{\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle}$  can be marked as *undefeated* or *defeated* nodes (U-nodes and D-nodes, resp.). A dialectical tree will be marked as an AND-OR tree: all leaves in  $\mathcal{T}_{\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle}$  will be marked as U-nodes (as they have no defeaters), and every inner node is to be marked as a *D-node* iff it has at least one U-node as a child, and as a *U-node* otherwise. An argument  $\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle$  is ultimately accepted as valid (or *warranted*) iff the root of  $\mathcal{T}_{\langle \mathcal{A}_0, Q_0, \alpha_0 \rangle}$  is labeled as a *U-node*.

**Definition 6.** [Warrant] [11] *Given a database  $\mathcal{DB}$ , and a literal  $Q$ ,  $Q$  is warranted w.r.t.  $\mathcal{DB}$  iff there exists a warranted argument  $\langle \mathcal{A}, Q, \alpha \rangle$  that can be built from  $\mathcal{P}$ .*

*Example 1.* Suppose the system has to solve the query  $\text{climbs}(\text{simba})$ . Then argument

$$\mathcal{A}_2 = \{(\text{climbs}(\text{simba}) \multimap \text{feline}(\text{simba}), 0.65), (\text{feline}(\text{simba}) \multimap \text{species}(\text{simba}, \text{lion}), 1)\}$$

must be built. This argument has a certainty degree of 0.6, taking into account the certainty degree of the literals on which the deduction is founded.

Next, the system looks for the defeaters. The only defeater is:

$$\langle \mathcal{A}_4, \sim \text{climbs}(\text{simba}), 0.6 \rangle, \mathcal{A}_4 = \{(\sim \text{climbs}(\text{simba}) \multimap \text{species}(\text{simba}, \text{lion}), 0.75)\}$$

But this argument is in turn defeated by  $\langle \mathcal{A}_3, \text{climbs}(\text{simba}), 0.6 \rangle$ ,

$$\mathcal{A}_3 = \{(\text{climbs}(\text{simba}) \multimap \text{species}(\text{simba}, \text{lion}), \text{age}(\text{simba}, \text{young}), 0.75)\}$$

Thus,  $\text{climbs}(\text{simba})$  is warranted.

## 5 Optimization of DB\_DeLP's Dialectical Process

To obtain faster query processing in the DB\_DeLP system we integrate pre-compiled knowledge to avoid the construction of arguments which were already computed. The approach follows the proposal presented in [6] where the pre-compiled knowledge component is required to: (1) minimize the number of stored arguments in the pre-compiled base of arguments (for instance, using one structure to represent the set of arguments that use the same defeasible rules); and (2) maintain independence from the observations that may change with new perception in order to avoid modifying also the pre-compiled knowledge when new observations are incorporated.

Considering these requirements, we define a database structure called *dialectical graph*, which will keep a record of all possible *arguments* in an DB\_DeLP database  $\mathcal{DB}$  (by means of a special structure named potential argument) as well as the counterargument relation among them. Potential arguments, originally defined in [6], contain non-grounded defeasible rules, thus depending only on the set of rules in the **IDB**, i.e., they are independent from the extensional database.

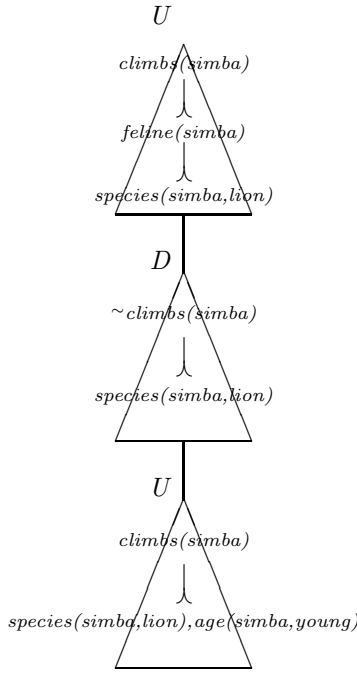


Fig. 4. Dialectical tree from Example 1

Potential arguments have been can be thought as schemata that sum up arguments that are obtained using *different* instances of the *same* defeasible rules. Recording every generated argument could result in storing many arguments which are structurally identical, only differing on the constants being used to build the corresponding derivations. Thus, a potential argument stands for several arguments which use the same defeasible rules. Attack relations among potential arguments can be also captured, and in some cases even defeat can be pre-compiled. In what follows we introduce the formal definitions, adapted from [6] to fit the DB\_DeLP system.

**Definition 7.** [Weighted Potential Argument] *Let  $IDB$  be an intensional database. A subset  $A$  of  $IDB$  is a potential argument for a literal  $Q$  with an upper bound  $\gamma$  for its certainty degree, noted as  $\langle\langle A, Q, \gamma \rangle\rangle$  if there exists a non-contradictory set of weighted literals  $\Phi$  and an instance  $\mathcal{A}$  that is obtained by finding an instance for every rule in  $A$ , such that  $\langle\mathcal{A}, Q, \alpha\rangle$  is an argument w.r.t. the database with  $\Phi$  as its extensional database and  $IDB$  as its intensional database ( $\alpha \leq \gamma$ ) and there is no instance  $\langle\mathcal{B}, Q, \beta\rangle$  of  $A$  such that  $\beta > \gamma$ .*

Definition 7 does not specify how to obtain the set of potential arguments from a given database. The interested reader may consult [6] for a constructive definition and its associated algorithm. The calculation of the upper bound  $\gamma$  deserves special mention, since the algorithm in [6] was devised for a different system,



without uncertainty management. This element will be used later on to speedup the extraction of the dialectical tree from the dialectical graph for a given query. To calculate  $\gamma$  for a potential argument  $A$  we simply choose the lower certainty degree of the defeasible rules present in  $A$ .

The nodes of the dialectical graph are the potential arguments. The arcs of our graph are obtained by calculating the counterargument relation among the nodes previously obtained. To do this, we extend the concept of counterargument for potential arguments. A potential argument  $\langle\langle A_1, Q_1, \alpha \rangle\rangle$  *counter-argues*  $\langle\langle A_2, Q_2, \beta \rangle\rangle$  at a literal  $Q$  if and only if there is a non-empty potential sub-argument  $\langle\langle A, Q, \gamma \rangle\rangle$  of  $\langle\langle A_2, Q_2, \beta \rangle\rangle$  such that  $Q_1$  and  $Q$  are contradictory literals.<sup>2</sup> Note that potential counter-arguments may or may not result in a real conflict between the instances (arguments) associated with the corresponding potential arguments. In some cases instances of these arguments cannot co-exist in any scenario (*e.g.*, consider two potential arguments based on contradictory observations). Now we can finally define the concept of dialectical graph:

**Definition 8.** [Dialectical Graph] *Let  $\mathcal{DB} = (EDB, IDB, IC)$  be a defeasible database. The dialectical graph of  $IDB$ , denoted as  $\mathcal{G}_{IDB}$ , is a pair*

$$(\text{PotArgs}(IDB), C)$$

such that:

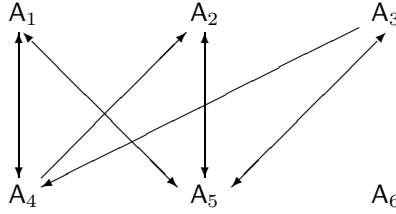
1.  $\text{PotArgs}(IDB)$  is the set  $\{\langle\langle A_1, Q_1, \alpha_1 \rangle\rangle, \dots, \langle\langle A_k, Q_k, \alpha_k \rangle\rangle\}$  of all the potential arguments that can be built from  $IDB$ ;
2.  $C$  is the counterargument relation over the elements of  $\text{PotArgs}(IDB)$ .

*Example 2.* Consider the feline database previously presented; its dialectical graph is composed by:

```
(feline(X)  $\neg$  species(X,lion),1)
(climbs_tree(X)  $\neg$  feline(X),0.65)
( $\sim$ climbs_tree(X)  $\neg$  species(X,lion),0.70)
(climbs_tree(X)  $\neg$  species(X,lion), age(X,young),0.75)
( $\sim$ climbs_tree(X)  $\neg$  sick(X),0.45)
```

- $\langle\langle A_1, \text{climbs}(X), 0.65 \rangle\rangle, A_1 = \{(\text{climbs}(X) \neg \text{feline}(X), 0.65)\}$ .
- $\langle\langle A_2, \text{climbs}(X), 0.65 \rangle\rangle, A_2 = \{(\text{climbs}(X) \neg \text{feline}(X), 0.65),$   
 $(\text{feline}(X) \neg \text{species}(X,\text{lion}), 1)\}$ .
- $\langle\langle A_3, \text{climbs}(X), 0.75 \rangle\rangle,$   
 $A_3 = \{(\text{climbs}(X) \neg \text{species}(X,\text{lion}), \text{age}(X,\text{young}), 0.75)\}$ .
- $\langle\langle A_4, \sim\text{climbs}(X), 0.75 \rangle\rangle, A_4 = \{(\sim\text{climbs}(X) \neg \text{species}(X,\text{lion}), 0.75)\}$ .
- $\langle\langle A_5, \sim\text{climbs}(X), 0.45 \rangle\rangle, A_5 = \{(\sim\text{climbs}(X) \neg \text{sick}(X), 0.45)\}$ .
- $\langle\langle A_6, \text{feline}(X), 1 \rangle\rangle, A_6 = \{(\text{feline}(X) \neg \text{species}(X,\text{lion}), 1)\}$ .
- $D_p = \{(A_2, A_4), (A_4, A_3)\}$
- $D_b = \{(A_1, A_4), (A_4, A_1), (A_1, A_5), (A_5, A_1), (A_2, A_5), (A_5, A_2), (A_3, A_5), (A_5, A_3)\}$ .

<sup>2</sup> Note that  $P(X)$  and  $\sim P(X)$  are contradictory literals even though they are non-grounded. The same idea is applied to identify contradiction in potential arguments.



**Fig. 5.** Dialectical graph corresponding to Example 2

The relations  $D_b$  and  $D_p$  can be depicted as shown in Figure 2, where blocking defeat is indicated with a double headed arrow.

Having defined the dialectical graph we can now use a specific graph traversing algorithm to extract a particular dialectical tree rooted in a given potential argument. The facts present in the **EDB** will be used as evidence to instantiate the potential arguments in the dialectical graph that depend on the intensional database **IDB**. This gives rise to the inference process of the system. This process starts when a new query is formulated. Consider the dialectical graph in Example 2 and suppose the set of facts in Figure 2 is present in the extensional database. Lets see how the system works when faced with the query  $climbs(simba)$ .

First, the set of potential arguments in the dialectical graph is searched to see if there exists an element whose conclusion can be instantiated to match the query. It finds  $\langle\langle A_2, climbs(X), 0.65 \rangle\rangle$ ,

$$A_2 = \{(\text{climbs}(X) \multimap \text{feline}(X), 0.65), (\text{feline}(X) \multimap \text{species}(X, \text{lion}), 1)\}$$

$A_2$  can be instantiated to

$$A_2 = \{(\text{climbs}(\text{simba}) \multimap \text{feline}(\text{simba}), 0.65), (\text{feline}(\text{simba}) \multimap \text{species}(\text{simba}, \text{lion}), 1)\}$$

that has a certainty degree of 0.6 taking into account the certainty degree of the literals on which the deduction is founded.

Now, to see if  $climbs(simba)$  is inferred by the system from the intensional and the extensional database, we must check whether  $A_2$  can sustain its conclusion when confronted with its counterarguments. Using the links in the dialectical graph we find one defeater for  $A_2$ , instantiating potential argument

$$A_4 = \{(\sim \text{climbs}(X) \multimap \text{species}(X, \text{lion}), 0.75)\}$$

to

$$A_4 = \{(\sim \text{climbs}(\text{simba}) \multimap \text{species}(\text{simba}, \text{lion}), 0.75)\}$$

The argument  $\langle A_4, \sim climbs(simba), 0.6 \rangle$  is defeated by  $\langle A_3, climbs(simba), 0.6 \rangle$  (an instantiation of  $\langle\langle A_3, climbs(X), 0.75 \rangle\rangle$ ). Thus,  $climbs(simba)$  is warranted and we found the same dialectical tree that was found in example 1 with an optimized inference mechanism. Note that the links for the defeaters present in the

dialectical graph are used to find the conflicts. This makes it easier to recover the tree from the dialectical graph of the framework.

The deductive database can be subject to constant changes as is the case with every real world database. The only restriction is that it must not be changed while a query is being solved. The dialectical graph is not affected by changes in the extensional database.

We present now a classic example in traditional deductive database systems based on logic programming, that usually causes problems with the semantics. In our case the system follows a cautious semantics, not deriving either  $p(a)$  or  $q(a)$ .

*Example 3.* Consider a deductive database composed by:

- **EDB** =  $\{(r, 0.6), (s, 0.6)\}$ ,
- **IDB** =  $\{(p(X) \multimap \sim q(X), 0.8), (q(X) \multimap \sim p(X), 0.8)\}$

The dialectical graph  $\mathbf{G}_{IDB}$  is composed by the two potential arguments:

- $\langle\langle A_1, p(X), \rangle\rangle A_1 = \{(p(X) \multimap \sim q(X), 0.8)\}$ .
- $\langle\langle A_2, q(X), \rangle\rangle A_1 = \{(q(X) \multimap \sim p(X), 0.8)\}$ .

and the defeat relation  $D_b = \{(A_1, A_2), (A_2, A_1)\}$ .

Suppose the system is faced with the query  $p(a)$ . The dialectical tree for this query is formed by argument  $\langle A_1, \sim q(a), 0.6 \rangle$ ,  $\mathcal{A}_1 = \{(p(a) \multimap \sim q(a), 0.6)\}$  that is in turn defeated by  $\langle A_2, q(a), 0.6 \rangle$ ,  $\mathcal{A}_1 = \{(q(a) \multimap \sim p(a), 0.6)\}$ .

The situation with query  $q(a)$  is analogous and therefore the system cannot derive  $p(a)$  nor  $q(a)$ .

Note that the DB\_DeLP system can seamlessly treat this example without semantic or operational problems. Furthermore, there is no need for imposing additional restrictions, such as requiring predicate stratification. Traditional systems would enter a loop jumping from one rule to the other. This is prevented in DB\_DeLP by the circularity condition imposed on argumentation lines of dialectical trees<sup>3</sup>. This condition does not allow the re-introduction of  $\mathcal{A}_1$  as a defeater of  $\mathcal{A}_2$  in the dialectical tree of the previous example.

## 6 A Worked Example

In this section we present an example to illustrate the practical uses of defeasible databases. The example is based on a classical benchmark in deductive databases concerning data on prescriptions, physicians and patients [21]. The system is a DSS to help employees decide whether a given medication should be sold to a patient. The relation *prescription* (pres) means that there is a prescription for a given drug to be administered to a given patient. *Allergic* shows known allergic reactions in patients, *physician* lists where physicians work, *patient* lists insurance company and clinics to which a patient usually goes, and *psychiatrist* (psy) establishes that a physician is also a psychiatrist (see Figure 6).

<sup>3</sup> This condition was inherited from the original DeLP system, the interested reader may consult [15].

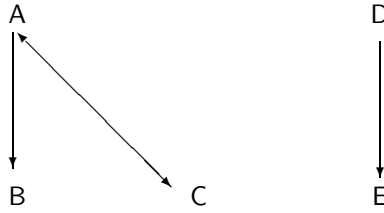
<u>patients(patient_id,clinic,insurance)</u>	<u>physicians(phy_id,clinic)</u>
(patients(456,new_line,hope),0.6)	(physicians(432,star),0.7)
(patients(587,delta,hope),0.6)	(physicians(54,delta),0.7)
(patients(234,new_line,trust),0.6)	(physicians(672,new_line),0.7)
(patients(1211,delta,trust),0.6)	(physicians(432,delta),0.7)
(patients(254,star,trust),0.6)	...
...	
<u>pres(note_id,patient_id,phy_id,drug,text)</u>	<u>allergic(patient_id,drug)</u>
(pres(23445,587,432,pen,text1),0.6)	(allergic(587,pen),0.7)
(pres(23446,587,54,amoxicillin,text2),0.6)	(allergic(1211,pen),0.6)
(pres(23447,587,54,vicodin,text3),0.6)	(allergic(1211,morphine),0.6)
(pres(23448,1211,54,morphine,text4),0.6)	...
(pres(23449,234,672,diazepam,text5),0.6)	
...	<u>psy(phy_id)</u>
	(psy(672),0.8)
	(psy(54),0.8)
	...

The intensional database is formed by the rules in Figure 6. The first rule says that a medication should be sold if there is prescription for it. The second rule says that it should not be sold if the physician is suspended and the third says that it should not be sold if the patient is allergic. The fourth rule concerns special drugs that have to be authorized before being sold and for that should have been prescribed by a psychiatrist. The fifth rule establishes the drug is not authorized when it is prescribed by a psychiatrist that has been suspended.

The dialectical graph contains arguments A, B, C, D, and E:

- $\langle\langle A, \text{sell}(\text{Patient}, \text{Drug}), 0.65 \rangle\rangle,$   
 $A = \{(\text{sell}(\text{Patient}, \text{Drug}) \multimap \text{pres}(X, \text{Patient}, Y, \text{Drug}), 0.65)\}.$
- $\langle\langle B, \sim\text{sell}(\text{Patient}, \text{Drug}), 0.75 \rangle\rangle,$   
 $B = \{(\sim\text{sell}(\text{Patient}, \text{Drug}) \multimap \text{pres}(X, \text{Patient}, Y, \text{Drug}, \text{Text}), \text{susp}(Y), 0.75)\}.$
- $\langle\langle C, \sim\text{sell}(\text{Patient}, \text{Drug}), 0.95 \rangle\rangle,$   
 $C = \{(\sim\text{sell}(\text{Patient}, \text{Drug}) \multimap \text{allergic}(\text{Patient}, \text{Drug}), 0.95)\}.$
- $\langle\langle D, \text{authorize\_pres}(\text{Patient}, \text{Drug}), 0.6 \rangle\rangle,$   
 $D = \{(\text{auth\_pres}(\text{Patient}, \text{Drug}) \multimap \text{pres}(X, \text{Patient}, Y, \text{Drug}, \text{Text}), \text{psy}(Y), 0.6)\}.$
- $\langle\langle E, \sim\text{authorize\_pres}(\text{Patient}, \text{Drug}), 0.7 \rangle\rangle,$   
 $E = \{(\sim\text{auth\_pres}(\text{Patient}, \text{Drug}) \multimap \text{pres}(X, \text{Patient}, Y, \text{Drug}, \text{Text}), \text{psy}(Y), \text{susp}(Y), 0.7)\}.$

(sell(Patient,Drug)  $\multimap$  pres(X,Patient,Y,Drug),0.65)  
 ( $\sim$ sell(Patient,Drug)  $\multimap$  pres(X,Patient,Y,Drug),susp(Y),0.75)  
 ( $\sim$ sell(Patient,Drug)  $\multimap$  allergic(Patient,Drug),0.95)  
 (auth\_pres(Patient,Drug)  $\multimap$  pres(X,Patient,Y,Drug),psychiatrist(Y),0.6)  
 ( $\sim$ auth\_pres(Patient,Drug)  $\multimap$  pres(X,Patient,Y,Drug),psy(Y),susp(Y),0.7)



**Fig. 6.** Dialectical graph for clinical database

Suppose the system is faced with a query for the fact  $\text{sell}(587, \text{vicodin})$ . It first finds a potential argument that can be instantiated to support this fact, so it selects A and instantiates it to:

$\mathcal{A} = \{(\text{sell}(787, \text{vicodin}) \rightarrow \text{pres}(23447, 587, 54, \text{vicodin}, \text{text3}), 0.6)\}$ . Using the dialectical graph we can see that there are two links that connect A with its defeaters, so we can explore to see if an instance of B or C can be built to attack argument  $\mathcal{A}$ . Since this is not the case argument  $\mathcal{A}$  is the only argument in the dialectical tree and the answer is yes.

Next, the system is faced with query  $\text{sell}(587, \text{pen})$ . The structure is similar to the previous case, but in this situation potential argument A is instantiated to

$$\mathcal{A} = \{(\text{sell}(587, \text{pen}) \rightarrow \text{pres}(23445, 587, 432, \text{pen}, \text{text1}), 0.6)\}$$

and following the links in the dialectical graph we find defeater B that can be instantiated to:

$\mathcal{B} = \{(\sim \text{sell}(587, \text{pen}) \rightarrow \text{allergic}(587, \text{pen}), 0.7)\}$ . No more defeaters can be added to this dialectical tree so the answer to  $\text{sell}(587, \text{pen})$  is no.

Now the query  $\text{auth\_pres}(234, \text{diazepam})$  is performed. In this case potential argument D is instantiated to:

$\mathcal{D} = \{\text{auth\_pres}(234, \text{diazepam}) \rightarrow \text{pres}(23449, 234, 672, \text{diazepam}, \text{text5}), \text{psy}(672), 0.6\}$  and no defeater can be found for  $\mathcal{D}$  thus the answer is yes.

Facts can be added to the database and also new tables can be created. Suppose we add a new table that contains a list of doctors that have been suspended due to legal issues. This table contains the fact  $(\text{suspended}(672), 0.8)$ . If query  $\text{authorize\_pres}(234, \text{diazepam})$  is re-processed by the system the answer would now be no, given that a new argument:

$\mathcal{E} = \{(\sim \text{auth\_pres}(234, \text{diazepam}) \rightarrow \text{pres}(23449, 234, 672, \text{diazepam}, \text{text5}), \text{psychiatrist}(672), \text{suspended}(672), 0.7)\}$ . can be built by instantiating E, resulting in a defeater for  $\mathcal{D}$ . Thus,  $\mathcal{D}$  is no longer warranted. Note how new tables and new facts can be added to the system without rebuilding the dialectical graph.

## 7 Conclusions and Future Work

In this work, we have defined a multi-agent system which virtually integrates different databases into a common view. We have also presented a layered architectural model that we have designed to develop practical applications for

reasoning with data from multiple sources. This model provides a novel system architecture for decision-support systems (DSS) that combines database technologies with an argumentation based framework.

We have also defined an argumentation-based formalism that integrates uncertainty management and is specifically tailored for database integration. This formalism was also provided with an optimization mechanism based on pre-compiled knowledge. Using this mechanism, the argumentation system can comply with real time requirements needed to manage data and model reasoning over this data in dynamic environments.

Future work may be done in different directions. First, many important and interesting issues could be considered in the general framework of database theory or information integration theory, such as how integrity constraints affect this set-up, how our proposal relates to local/global views, or which connections could be established with database repairs. Second, we will integrate DB\_DeLP with a massive data component to obtain experimental results regarding the system's behavior in such trying environment.

A related research line could also be obtained by extending the language of DB\_DeLP to use it in practical systems, particularly to implement argumentation-based active databases and decision support systems backed by large repositories of data.

**Acknowledgements.** This research was partially supported by CONICET (Argentina), and by the *Universidad Nacional del Sur*. The authors would also like to thank the anonymous reviewers for their valuable comments and suggestions that helped improve the quality of this work.

## References

1. Alsinet, T., Chesñevar, C.I., Godo, L., Simari, G.R.: A logic programming framework for possibilistic argumentation: Formalization and logical properties. *Fuzzy Sets and Systems* 159(10), 1208–1228 (2008)
2. Alsinet, T., Godo, L.: A complete calculus for possibilistic logic programming with fuzzy propositional variables. In: *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence (UAI 2000)*, pp. 1–10. ACM Press, New York (2000)
3. Berti, L.: Quality and recommendation of multi-source data for assisting technological intelligence applications. In: *Proc. of 10th International Conference on Database and Expert Systems Applications*, pp. 282–291. AAAI, Italy (1999)
4. Brodie, M.L., Jarke, M.: On integrating logic programming and databases. In: *Expert Database Workshop 1984*, pp. 191–207 (1984)
5. Bryant, D., Krause, P.: An implementation of a lightweight argumentation engine for agent applications. *Logics in Artificial Intelligence* 4160(1), 469–472 (2006)
6. Capobianco, M., Chesñevar, C.I., Simari, G.R.: Argumentation and the dynamics of warranted beliefs in changing environments. *Journal of Autonomous Agents and Multiagent Systems* 11, 127–151 (2005)
7. Carbogim, D., Robertson, D., Lee, J.: Argument-based applications to knowledge engineering. *The Knowledge Engineering Review* 15(2), 119–149 (2000)

8. Ceri, S., Gottlob, G., Tanca, L.: What you always wanted to know about datalog (and never dared to ask). *IEEE Trans. on Knowledge and Data Eng.* 1(1) (1989)
9. Chesñevar, C.I., Maguitman, A.G., Loui, R.P.: Logical Models of Argument. *ACM Computing Surveys* 32(4), 337–383 (2000)
10. Chesñevar, C.I., Maguitman, A.G., Simari, G.R.: Argument-based critics and recommenders: A qualitative perspective on user support systems. *Data & Knowledge Engineering* 59(2), 293–319 (2006)
11. Chesñevar, C.I., Simari, G.R., Alsinet, T., Godo, L.: A logic programming framework for possibilistic argumentation with vague knowledge. In: *Proc. of Uncertainty in Artificial Intelligence Conference (UAI 2004)*, Banff, Canada (2004)
12. Cuppens, F., Demolombe, R.: Cooperative answering: a method to provide intelligent access to databases. In: *Proc. 2nd Conf. on Expert Database Systems*, pp. 621–643 (1988)
13. Dubois, D., Lang, J., Prade, H.: Possibilistic logic. In: Gabbay, D., Hogger, C., Robinson, J. (eds.) *Handbook of Logic in Art. Int. and Logic Prog (Nonmonotonic Reasoning and Uncertain Reasoning)*, pp. 439–513. Oxford University Press, Oxford (1994)
14. Dung, P.M.: On the Acceptability of Arguments and its Fundamental Role in Nonmonotonic Reasoning and Logic Programming and n-Person Games. *Artificial Intelligence* 77(2), 321–357 (1995)
15. García, A.J., Simari, G.R.: Defeasible Logic Programming: An Argumentative Approach. *Theory and Practice of Logic Programming* 4(1), 95–138 (2004)
16. Lakshmanan, L.V.S., Sadri, F.: Probabilistic deductive databases. In: *Proc. of the Int. Logic Programming Symposium*, pp. 254–268 (1994)
17. Lakshmanan, L.V., Shiri, N.: A parametric approach to deductive databases with uncertainty. *Journal of Intelligent Information Systems* 13(4), 554–570 (2001)
18. McLeod, D., Heimbigner, D.: A federated architecture for information management. *ACM Transactions on Information Systems* 3(3), 253–278 (1985)
19. Prakken, H., Vreeswijk, G.: Logical systems for defeasible argumentation. In: *Handbook of Philosophical Logic*, vol. 4, pp. 219–318 (2002)
20. Prakken, H., Vreeswijk, G.: Logical systems for defeasible argumentation. In: *Handbook of Philosophical Logic*, vol. 4, pp. 219–318 (2002)
21. Quian, X.: Query folding. In: *Proc. 12th Intl. Conf on Data Engineering*, pp. 48–55 (1996)
22. Simari, G.R., Loui, R.P.: A Mathematical Treatment of Defeasible Reasoning and its Implementation. *Artificial Intelligence* 53(1-2), 125–157 (1992)
23. Subrahmanian, V.S.: Paraconsistent disjunctive deductive databases. *Theoretical Computer Science* 93(1), 115–141 (1992)
24. Zaniolo, C.: Prolog: A database query language for all seasons. In: *Expert Database Workshop 1984*, pp. 219–232 (1984)
25. Zaniolo, C.: Intelligent databases: Old challenges and new opportunities 3/4(1), 271–292 (1992)

# Argumentation in the View of Modal Logic

Davide Grossi

Institute for Logic, Language and Computation  
University of Amsterdam  
Science Park 904, 1098 XH Amsterdam, The Netherlands  
d.grossi@uva.nl

**Abstract.** The paper presents a study of abstract argumentation theory from the point of view of modal logic. The key thesis upon which the paper builds is that argumentation frameworks can be studied as Kripke frames. This simple observation allows us to import a number of techniques and results from modal logic to argumentation theory, and opens up new interesting avenues for further research. The paper gives a glimpse of the sort of techniques that can be imported, discussing complete calculi for argumentation, adequate model-checking and bisimulation games, and sketches an agenda for future research at the interface of modal logic and argumentation theory.

**Keywords:** Argumentation theory, modal logic.

## 1 Introduction

The paper advocates a perspective on abstract argumentation theory based on techniques and results borrowed from the field of formal logic and, in particular, of modal logic [3]. First steps in this line of research have been moved in [12] and [13]. The present paper recapitulates some of the results presented in those works and sketches a number of theoretical problems arising at the interface of logic and argumentation which constitute, in our view, an interesting and challenging agenda for future research in both disciplines.

The key point of the paper is that standard results in argumentation theory obtain elegant reformulations within well-investigated modal logics. Once this link is established a number of techniques (e.g., calculi, logical games), as well as results related to those techniques (e.g. completeness, adequacy), can be easily imported from modal logic to argumentation theory.

The paper presupposes some familiarity with both modal logic and abstract argumentation theory. Proofs are omitted for space reasons. The interested reader is referred to [12][13].

**Outline of the paper.** Section 2 starts off by applying a well-known modal logic to study a first set of notions of argumentation theory. This enables the possibility of using calculi to derive argumentation-theoretic results such as the Fundamental Lemma [7]. Along the same line, Section 3 tackles the formalization of the notion of grounded extension within the modal  $\mu$ -calculus. In Section 4 semantic games are studied for the logic introduced in Section 2 which provide a version of



**Table 1.** Basic notions of argumentation theory ( $X$  denotes a set of arguments)

---

$c_{\mathcal{A}}$ is the characteristic function of $\mathcal{A}$ iff $c_{\mathcal{A}} : 2^A \rightarrow 2^A$ s.t.	$c_{\mathcal{A}}(X) = \{a \mid \forall b : [b \rightarrow a \Rightarrow \exists c \in X : c \rightarrow b]\}$
$X$ is acceptable w.r.t. $Y$ in $\mathcal{A}$	iff $X \subseteq c_{\mathcal{A}}(Y)$
$X$ is conflict-free in $\mathcal{A}$	iff $\nexists a, b \in X$ s.t. $a \rightarrow b$
$X$ is admissible in $\mathcal{A}$	iff $X$ is conflict-free and $X \subseteq c_{\mathcal{A}}(X)$ iff $X$ is a post-fixpoint of $c_{\mathcal{A}}$
$X$ is a complete extension of $\mathcal{A}$	iff $X$ is conflict-free and $X = c_{\mathcal{A}}(X)$ ( $X$ is a conflict-free fixpoint of $c_{\mathcal{A}}$ )
$X$ is a stable extension of $\mathcal{A}$	iff $X$ is a complete extension of $\mathcal{A}$ and $\forall b \notin X, \exists a \in X : a \rightarrow b$ iff $X = \{a \in A \mid \nexists b \in X : b \rightarrow a\}$
$X$ is the grounded extension of $\mathcal{A}$	iff $X$ is the minimal complete extension of $\mathcal{A}$ iff $X$ is the least fixpoint of $c_{\mathcal{A}}$
$X$ is a preferred extension of $\mathcal{A}$	iff $X$ is a maximal complete extension of $\mathcal{A}$

---

games for argumentation by means of model-checking games. Section 5 tackles the question—not yet addressed in the literature—of when two arguments, or two argumentation frameworks, are equivalent from the point of view of argumentation theory. For this purpose the model-theoretic notion of bisimulation is introduced and bisimulation games are presented as a procedural method to check the ‘behavioral equivalence’ of two argumentation frameworks. Section 6 sketches some of the possible lines of research that we consider worth pursuing by applying logic-based methods to abstract argumentation. Section 7 briefly concludes.

## 2 Arguments in Modal Disguise

The section moves the first steps towards looking at argumentation frameworks as structures upon which to interpret modal languages.

### 2.1 Argumentation Frameworks

Let us start with the basic structures of argumentation theory [7].

**Definition 1 (Argumentation frameworks).** *An argumentation framework is a relational structure  $\mathcal{A} = (A, \rightarrow)$  where  $A$  is a non-empty set of arguments,*

and  $\rightarrow_{\subseteq} A^2$  is a so-called ‘attack’ relation on  $A$ . A pointed argumentation framework is a pair  $(\mathcal{A}, a)$  with  $a \in A$ . The set of all argumentation frameworks is called  $\mathfrak{A}$ .

The intuitive reading of “ $a \rightarrow b$ ” is that argument  $a$  attacks argument  $b$ . Doing abstract argumentation theory means, essentially, to study specific properties of subsets of the set of arguments  $A$  in a given  $\mathcal{A}$ . For space reasons the paper cannot introduce argumentation theory in an extensive way but, to make it as most self-contained as possible, the main argumentation-theoretic notions from [7] have been recapitulated in Table 1. As such notions are formalized along the paper, their intuitive reading will also be provided.

The paper is based on the simple idea of viewing argumentation frameworks as the structures known in modal logic as Kripke frames, that is, structures  $(S, R)$  where  $S$  is taken to be a non-empty set of states, and  $R$  a binary relation on elements of  $S$  [3]. In essence, the paper studies what modal logic can say about argumentation frameworks when  $S$  is set to be  $A$ , i.e., the modal states are taken to be arguments, and  $R$  is set to be the inverse of the attack relation, that is, relation  $\rightarrow^{-1}$ . The entire paper and all its results hinge on this simple assumptions.

The reader might ask himself why  $R$  is taken to be the inverse  $\rightarrow^{-1}$  of the attack relation instead of the attack relation  $\rightarrow$  itself. This will become clear as the paper develops. However, a simple inspection of Table 1 should already show that all the key argumentation-theoretic notions can be defined in terms of the characteristic function, and that the characteristic function is defined by taking, for any argument  $a$  in the given input  $X$ , the set of attackers  $b$  of  $a$ —that is, the set of arguments by which  $a$  is attacked—for which there always exists another attacker  $c$ —that is, an argument by which the attacker of  $a$  is attacked. So, the characteristic function looks, for any argument, at whether its attackers are attacked. To put it in modal logic terms, the characteristic function is defined by looking at the tree-unraveling of  $\rightarrow^{-1}$  at each point  $a$ , and not at the tree-unraveling of  $\rightarrow$ . We will come back to this issue in Section 3.1, now we proceed to the use of argumentation frameworks in a modal logic setting.

## 2.2 Argumentation Models

If an argumentation framework can be viewed as a Kripke frame, then an argumentation framework plus a function assigning names from a set  $\mathbf{P}$  to sets of arguments can be viewed as a Kripke model [3].

**Definition 2 (Argumentation models).** *Let  $\mathbf{P}$  be a set of propositional atoms. An argumentation model  $\mathcal{M} = (\mathcal{A}, \mathcal{I})$  is a structure such that:  $\mathcal{A} = (A, \rightarrow)$  is an argumentation framework;  $\mathcal{I} : \mathbf{P} \rightarrow 2^A$  is an assignment from  $\mathbf{P}$  to subsets of  $A$ . The set of all argumentation models is called  $\mathfrak{M}$ . A pointed argumentation model is a pair  $(\mathcal{M}, a)$  where  $\mathcal{M}$  is an argumentation model and  $a$  an argument from  $A$ .*

Argumentation models are nothing but argumentation frames together with a way of ‘naming’ sets of arguments or, to put it otherwise, of ‘labeling’ arguments.

The fact that an argument  $a$  belongs to  $\mathcal{I}(p)$  in a given model  $\mathcal{M}$ , which in logical notation reads  $(\mathcal{A}, \mathcal{I}), a \models p$ , can be interpreted as stating that “argument  $a$  has property  $p$ ”, or that “ $p$  is true of  $a$ ”. By substituting  $p$  with a Boolean compound  $\varphi$  (e.g.,  $\varphi := p \wedge q$ ) we can say that “ $a$  belongs to both the sets called  $p$  and  $q$ ”, and the same can be done for all other Boolean connectives.

This much as to Boolean properties of arguments. But what about statements of the sort: “argument  $a$  is attacked by an argument in a set  $\varphi$ ”; “argument  $a$  is defended by the set  $\varphi$ ”, or, “ $\varphi$  attacks an attacker of argument  $a$ ”? These are modal statements, and in order to express them, it suffices to introduce a dedicated modal operator  $\langle \leftarrow \rangle$  whose intuitive reading is “there exists an attacking argument such that”. To this we turn in the next section.

### 2.3 Argumentation and Logic $\mathbf{K}^\forall$

This section introduces logic  $\mathbf{K}^\forall$ , an extension of the minimal modal logic  $\mathbf{K}$  with universal modality. The section shows how such a simple and standard modal logic is already able of capturing quite a few argumentation-theoretic notions.

**Language.** The language of  $\mathbf{K}^\forall$  is a standard modal language with two modalities:  $\langle \leftarrow \rangle$  and  $\langle \forall \rangle$ , i.e., the universal modality. It is built on the set of atoms  $\mathbf{P}$  by the following BNF:

$$\mathcal{L}^{\mathbf{K}^\forall} : \varphi ::= p \mid \perp \mid \neg\varphi \mid \varphi \wedge \varphi \mid \langle \leftarrow \rangle\varphi \mid \langle \forall \rangle\varphi$$

where  $p$  ranges over  $\mathbf{P}$ . The other standard boolean  $\{\top, \vee, \rightarrow\}$  and modal  $\{\langle \leftarrow \rangle, \langle \forall \rangle\}$  connectives are defined as usual.

#### Semantics

**Definition 3 (Satisfaction).** *Let  $\varphi \in \mathcal{L}^{\mathbf{K}^\forall}$ . The satisfaction of  $\varphi$  by a pointed argumentation model  $(\mathcal{M}, a)$  is inductively defined as follows (Boolean clauses are omitted):*

$$\begin{aligned} \mathcal{M}, a \models \langle \leftarrow \rangle\varphi & \text{ iff } \exists b \in A : (a, b) \in \rightarrow^{-1} \text{ AND } \mathcal{M}, b \models \varphi \\ \mathcal{M}, a \models \langle \forall \rangle\varphi & \text{ iff } \exists b \in A : \mathcal{M}, b \models \varphi \end{aligned}$$

As usual,  $\varphi$  is valid in an argumentation model  $\mathcal{M}$  iff it is satisfied in all pointed models of  $\mathcal{M}$ , i.e.,  $\mathcal{M} \models \varphi$ . The truth-set of a formula  $\varphi$  is denoted  $|\varphi|_{\mathcal{M}}$ .

Logic  $\mathbf{K}^\forall$  is therefore endowed with modal operators of the type “there exists an argument attacking the current one such that”, i.e.,  $\langle \leftarrow \rangle$ , and “there exists an argument such that”, i.e.,  $\langle \forall \rangle$ , together with their duals. Given an argumentation model  $\mathcal{M}$  we can thereby express statements such as the ones adverted to above: “ $a$  is attacked by an argument in a set called  $\varphi$ ” corresponds to  $\langle \leftarrow \rangle\varphi$  being true in the pointed model  $(\mathcal{M}, a)$  and “ $a$  is defended by the set  $\varphi$ ” corresponds to  $\langle \leftarrow \rangle\langle \leftarrow \rangle\varphi$  being true in the pointed model  $(\mathcal{M}, a)$ .

On the ground of this semantics, it becomes already clear that logic  $\mathbf{K}^\forall$  is expressive enough to capture several basic notions of argumentation theory such

as: conflict freeness, acceptability, admissibility, complete extensions, stable extensions.

$$Acc(\varphi, \psi) := [\forall](\varphi \rightarrow [\leftarrow]\langle \leftarrow \rangle \psi) \quad (1)$$

$$CFree(\varphi) := [\forall](\varphi \rightarrow [\leftarrow]\neg\varphi) \quad (2)$$

$$Adm(\varphi) := [\forall](\varphi \rightarrow ([\leftarrow]\neg\varphi \wedge [\leftarrow]\langle \leftarrow \rangle \varphi)) \quad (3)$$

$$Compl(\varphi) := [\forall](\varphi \rightarrow [\leftarrow]\neg\varphi \wedge (\varphi \leftrightarrow [\leftarrow]\langle \leftarrow \rangle \varphi)) \quad (4)$$

$$Stable(\varphi) := [\forall](\varphi \leftrightarrow [\leftarrow]\neg\varphi) \quad (5)$$

Intuitively, a set of arguments  $\varphi$  is acceptable with respect to the set of arguments  $\psi$  if and only all  $\varphi$ -arguments are such that for all their attackers there exists a defender in  $\psi$  (Formula [1](#)). A set of arguments  $\varphi$  is conflict free if and only if all  $\varphi$ -arguments are such that none of their attackers is in  $\varphi$  (Formula [2](#)). A set of arguments  $\varphi$  is admissible if and only if it is conflict free and acceptable with respect to itself (Formula [3](#)). A set  $\varphi$  is a complete extension if and only if it is conflict free and it is equivalent to the set of arguments all the attackers of which are attacked by some  $\varphi$ -argument (Formula [4](#)). Finally, a set  $\varphi$  is a stable extension if and only if it is equivalent to the set of arguments whose attackers are not in  $\varphi$  (Formula [5](#)). The adequacy of these definitions with respect to the ones in Table [1](#) is easily checked.

**Axiomatics.** Logic  $K^\forall$  is axiomatized as follows, where  $i \in \{\leftarrow, \forall\}$ :

- (**Prop**) propositional tautologies
- (**K**)  $[i](\varphi_1 \rightarrow \varphi_2) \rightarrow ([i]\varphi_1 \rightarrow [i]\varphi_2)$
- (**T**)  $[\forall]\varphi \rightarrow \varphi$
- (**4**)  $[\forall]\varphi \rightarrow [\forall][\forall]\varphi$
- (**5**)  $\neg[\forall]\varphi \rightarrow [\forall]\neg[\forall]\varphi$
- (**Incl**)  $[\forall]\varphi \rightarrow [i]\varphi$
- (**Dual**)  $\langle i \rangle \varphi \leftrightarrow \neg[i]\neg\varphi$

The axiom system combines the axioms of logic **K** for the  $[\leftarrow]$  operator, the axioms of logic **S5** for the universal operator  $[\forall]$ , and the interaction axiom **Incl**. It can be proven that this axiomatics is sound and strongly complete for the class  $\mathfrak{A}$  of argumentation frames [\[3, Ch. 7\]](#).

The fact that  $K^\forall$  is axiomatized by the axioms and rules above gives us a calculus by means of which we can prove theorems of abstract argumentation theory in a purely formal manner. A notable example is the following generalized version of the fundamental lemma from [\[7\]](#), which states that if  $\varphi$  is admissible and both  $\psi$  and  $\xi$  are acceptable with respect to it, then also  $\psi \vee \xi$  is admissible and  $\xi$  is acceptable with respect to  $\varphi \vee \psi$ .

**Theorem 1 (Fundamental Lemma [\[7\]](#)).** *The following formula is a theorem of  $K^\forall$ :*

$$Adm(\varphi) \wedge Acc(\psi \vee \xi, \varphi) \rightarrow Adm(\varphi \vee \psi) \wedge Acc(\xi, \varphi \vee \psi) \quad (6)$$

The theorem could be proven semantically by then calling in completeness. However, to give a detailed example of an application of the above axiomatics, a formal derivation of the theorem is provided in the appendix.

Other examples of theorems of [7] that could be casted in this logic are, for instance:  $Stable(\varphi) \rightarrow Adm(\varphi)$  and  $Stable(\varphi) \rightarrow Compl(\varphi)$ .

### 3 Modal Fixpoints

The present section shows what kind of modal machinery is needed to capture the notion of grounded extension left aside in Section 2. In [7], the grounded extension is defined as the smallest fixpoint of the characteristic function of an argumentation framework (see Table 1).

#### 3.1 Characteristic Functions in Modal Logic

Each argumentation framework  $\mathcal{A} = (A, \rightarrow)$  determines a *characteristic function*  $c_{\mathcal{A}} : 2^A \rightarrow 2^A$  such that for any set of arguments  $X$ ,  $c_{\mathcal{A}}(X)$  yields the set of arguments in  $A$  which are acceptable with respect to  $X$ , i.e.,  $\{a \in A \mid \forall b \in A : [b \rightarrow a \Rightarrow \exists c \in X : c \rightarrow b]\}$ . Does logic  $\mathbf{K}^{\vee}$  have a syntactic counterpart of the characteristic function? The answer turns out to be yes.

Let  $\mathcal{L}^{[\leftarrow]\langle\leftarrow\rangle}$  be the language defined by the following BNF:

$$\mathcal{L}^{[\leftarrow]\langle\leftarrow\rangle} : \varphi ::= p \mid \perp \mid \neg\varphi \mid \varphi \wedge \varphi \mid [\leftarrow]\langle\leftarrow\rangle\varphi$$

where  $p$  belongs to the set of atoms  $\mathbf{P}$ . Language  $\mathcal{L}^{[\leftarrow]\langle\leftarrow\rangle}$  is the fragment of  $\mathcal{L}^{\mathbf{K}^{\vee}}$  containing only the compounded modal operator  $[\leftarrow]\langle\leftarrow\rangle$  or, also, simply the fragment of  $\mathcal{L}^{\mathbf{K}}$  (i.e.,  $\mathcal{L}^{\mathbf{K}^{\vee}}$  without universal modality) containing only the  $[\leftarrow]\langle\leftarrow\rangle$ -operator. Let  $\mathcal{A}^+ = (2^A, \cap, -, \emptyset, c_{\mathcal{A}})$  be the power set algebra on  $2^A$  extended with operator  $c_{\mathcal{A}}$ , and consider the term algebra  $\mathbf{ter}_{\mathcal{L}^{[\leftarrow]\langle\leftarrow\rangle}} = (\mathcal{L}^{[\leftarrow]\langle\leftarrow\rangle}, \wedge, \neg, \perp, [\leftarrow]\langle\leftarrow\rangle)$ . Finally, let  $\mathcal{I}^* : \mathcal{L}^{[\leftarrow]\langle\leftarrow\rangle} \rightarrow 2^A$  be the inductive extension of a valuation function  $\mathcal{I} : \mathbf{P} \rightarrow 2^A$  according to the semantics given in Definition 3. We can prove the following result.

**Theorem 2** ( $c_{\mathcal{A}}$  vs.  $[\leftarrow]\langle\leftarrow\rangle$ ). *Let  $\mathcal{M} = (\mathcal{A}, \mathcal{I})$  be an argumentation model. Function  $\mathcal{I}^*$  is a homomorphism from  $\mathbf{ter}_{\mathcal{L}^{[\leftarrow]\langle\leftarrow\rangle}}$  to  $\mathcal{A}^+$ .*

In other words, Theorem 2 shows that the complex modal operator  $[\leftarrow]\langle\leftarrow\rangle$ , under the semantics provided in Definition 3, behaves exactly like the characteristic function of the argumentation frameworks on which the argumentation models are built. To put it yet otherwise, formulae of the form  $[\leftarrow]\langle\leftarrow\rangle\varphi$  denote the value of the characteristic function applied to the set  $\varphi$  of arguments. Notice also that from Theorem 2 the adequacy of Formulae 1-5 with respect to the definitions in Table 1 follows straightforwardly.

Characteristic functions are known to be monotonic [7] hence, by Theorem 2, we get that  $[\leftarrow]\langle\leftarrow\rangle$  denotes a monotonic function and therefore, by the Knaster-Tarski theorem [1] we have that there always exist a greatest and a least  $[\leftarrow]\langle\leftarrow\rangle$ -fixpoint. From a logical point of view this means that, in order to be able to

<sup>1</sup> We refer the interested reader to [5].

express the grounded extension, it suffices to add to the  $K$  fragment of  $K^\forall$  a least fixpoint operator. This takes us to the realm of  $\mu$ -calculus.

### 3.2 Argumentation and the $\mu$ -Calculus

**Language.** To add the least fixpoint operator  $\mu$  to logic  $K$  we first define language  $\mathcal{L}^{K^\mu}$  via the following BNF:

$$\mathcal{L}^{K^\mu} : \varphi ::= p \mid \perp \mid \neg\varphi \mid \varphi \wedge \varphi \mid \langle \leftarrow \rangle \varphi \mid \mu p. \varphi(p)$$

where  $p$  ranges over  $\mathbf{P}$  and  $\varphi(p)$  indicates that  $p$  occurs free in  $\varphi$  (i.e., it is not bounded by fixpoint operators) and under an even number of negations.<sup>2</sup> In general, the notation  $\varphi(\psi)$  stands for  $\psi$  occurs in  $\varphi$ . The usual definitions for Boolean and modal operators can be applied. Intuitively,  $\mu p. \varphi(p)$  denotes the smallest formula  $p$  such that  $p \leftrightarrow \varphi(p)$ . This intuition is made precise in the semantics of  $\mathcal{L}^{K^\mu}$ .

#### Semantics

**Definition 4 (Satisfaction).** Let  $\varphi \in \mathcal{L}^{K^\mu}$ . The satisfaction of  $\varphi$  by a pointed model  $(\mathcal{M}, a)$ , with  $\mathcal{M} = (\mathcal{A}, \mathcal{I})$ , is inductively defined as follows (Boolean clauses, as well as the clause for  $\langle \leftarrow \rangle$ , are as in Definition 3):

$$\mathcal{M}, a \models \mu p. \varphi(p) \text{ iff } a \in \bigcap \{X \in 2^A \mid |\varphi|_{\mathcal{M}[p:=X]} \subseteq X\}$$

where  $|\varphi|_{\mathcal{M}[p:=X]}$  denotes the truth-set of  $\varphi$  once  $\mathcal{I}(p)$  is set to be  $X$ . As usual, we say that:  $\varphi$  is valid in an argumentation model  $\mathcal{M}$  iff it is satisfied in all pointed models of  $\mathcal{M}$ , i.e.,  $\mathcal{M} \models \varphi$ ;  $\varphi$  is valid in a class  $\mathfrak{M}$  of argumentation models iff it is valid in all its models, i.e.,  $\mathfrak{M} \models \varphi$ .

We have now all the logical machinery in place to express the notion of grounded extension. Set  $\varphi(p) := [\leftarrow] \langle \leftarrow \rangle p$ , that is, take  $\varphi(p)$  to be the modal version  $[\leftarrow] \langle \leftarrow \rangle$  of the characteristic function, and apply it to formula  $p$ . What we obtain is a modal formula expressing the least fixpoint of a characteristic function, that is, the grounded extension:

$$\text{Grounded} := \mu p. [\leftarrow] \langle \leftarrow \rangle p \tag{7}$$

Notice that, unlike the notions formalized in Formulae 1.5, the grounded extension of a framework is always unique and does not depend on the particular labeling of a given model.

We refrain here from providing a sound and complete axiomatization of  $\mu$ -calculus. The interested reader is referred to [19]. However, just like we did for logic  $K^\forall$  we give now a couple of examples of the kind of argumentation-theoretic results formalizable in  $K^\mu$ . Well-known theorems of argumentation theory are provable formulae of the  $\mu$ -calculus.

<sup>2</sup> This syntactic restriction guarantees that every formula  $\varphi(p)$  defines a monotonic set transformation.

**Theorem 3 (The grounded extension is conflict-free).** *The following formula is a theorem of  $K^\mu$ :*

$$\text{Grounded} \rightarrow [\leftarrow] \neg \text{Grounded} \tag{8}$$

We can also study complexity results from this modal perspective and, unsurprisingly, the results are in accordance with complexity studies in argumentation theory [8], although the proofs take different routes.

**Theorem 4 (Model-checking grounded).** *Let  $\mathcal{A}$  be an argumentation framework. It can be decided in polynomial time whether an argument  $a$  belongs to the grounded extension of  $\mathcal{A}$ , that is, whether  $\mathcal{A}, a \models \text{Grounded}$ .*

## 4 Dialogic Games and Logic Games

The proof-theory of abstract argumentation is commonly given in terms of dialogue games [16]. The present section introduces a new game-theoretic proof procedure for argumentation theory based on model-checking games. In model-checking games, a proponent or verifier ( $\exists$ ve) tries to prove that a given formula  $\varphi$  holds in a point  $a$  of a model  $\mathcal{M}$ , while an opponent or falsifier ( $\forall$ dam) tries to disprove it. The present section deals with the model-checking game for  $K^\forall$ . For the  $K^\mu$ -variant of this game we refer the reader to [18].

### 4.1 Model-Checking Game for $K^\forall$

A model-checking game is a *graph game*, that is, a game played by two agents on a directed graph, where each node—called position—is labelled by the player that is supposed to move next. The structure of the graph determines which are the *admissible moves* at any given position. If a player has to move in a certain position but there are no available moves, then it loses and its opponent wins. In general, graph games might have infinite paths, but this is not the case in the game we are going to introduce. A match of a graph game is then just the set of positions visited during play, that is, a complete path through the graph.

**Definition 5 ( $K^\forall$ -model-checking game).** *Let  $\varphi \in \mathcal{L}^{K^\forall}$ , and  $\mathcal{M}$  be an argumentation model. The model-checking game  $\mathcal{C}(\varphi, \mathcal{M})$  is defined by the following items. **Players:** The set of players is  $\{\exists, \forall\}$ . An element from  $\{\exists, \forall\}$  will be denoted  $P$  and its opponent  $\bar{P}$ . **Game form:** The game form of  $\mathcal{C}(\varphi, \mathcal{M})$  is defined by the board game in Table 2. **Winning conditions:** Player  $P$  wins if and only if  $\bar{P}$  has to play in a position with no available moves. **Instantiation:** The instance of  $\mathcal{C}(\varphi, \mathcal{M})$  with starting point  $(\varphi, a)$  is denoted  $\mathcal{C}(\varphi, \mathcal{M})@(\varphi, a)$ .*

The important thing to notice is that positions of the game are pairs of a formula and an argument, and that the type of formula in the position determines which player has to play:  $\exists$  if the formula is a disjunction, a diamond, a false atom or  $\perp$ , and  $\forall$  in the remaining cases.<sup>3</sup>

We now define what it means to have a winning strategy and to be in a winning position in this type of games.

<sup>3</sup> Notice that positions use formulae in positive normal form.

**Table 2.** Rules of the model-checking game for  $K^\forall$

Position	Turn	Available moves
$(\varphi_1 \vee \varphi_2, a)$	$\exists$	$\{(\varphi_1, a), (\varphi_2, a)\}$
$(\varphi_1 \wedge \varphi_2, a)$	$\forall$	$\{(\varphi_1, a), (\varphi_2, a)\}$
$(\langle \leftarrow \rangle \varphi, a)$	$\exists$	$\{(\varphi, b) \mid (a, b) \in \rightarrow^{-1}\}$
$([\leftarrow] \varphi, a)$	$\forall$	$\{(\varphi, b) \mid (a, b) \in \rightarrow^{-1}\}$
$(\langle \forall \rangle \varphi, a)$	$\exists$	$\{(\varphi, b) \mid b \in A\}$
$([\forall] \varphi, a)$	$\forall$	$\{(\varphi, b) \mid b \in A\}$
$(\perp, a)$	$\exists$	$\emptyset$
$(\top, a)$	$\forall$	$\emptyset$
$(p, a) \ \& \ a \notin \mathcal{I}(p)$	$\exists$	$\emptyset$
$(p, a) \ \& \ a \in \mathcal{I}(p)$	$\forall$	$\emptyset$
$(\neg p, a) \ \& \ a \in \mathcal{I}(p)$	$\exists$	$\emptyset$
$(\neg p, a) \ \& \ a \notin \mathcal{I}(p)$	$\forall$	$\emptyset$

**Definition 6 (Winning strategies and positions).** A strategy for player  $P$  in  $\mathcal{C}(\varphi, \mathcal{M})@(\varphi, a)$  is a function telling  $P$  what to do in any match played from position  $(\varphi, a)$ . Such a strategy is winning for  $P$  if and only if, in any match played according to the strategy,  $P$  wins. A position  $(\varphi, a)$  in  $\mathcal{C}(\varphi, \mathcal{M})$  is winning for  $P$  if and only if  $P$  has a winning strategy in  $\mathcal{C}(\varphi, \mathcal{M})@(\varphi, a)$ . The set of winning positions of  $\mathcal{C}(\varphi, \mathcal{M})$  is denoted  $Win_P(\mathcal{C}(\varphi, \mathcal{M}))$ .

By Definitions 4.2 and 6 it follows that the model-checking game is a two-players zero-sum game with perfect information. It is known that such games are determined, that is, each match has a winner [20].

These games can be proven adequate. This means that if Eve has a winning strategy then the formula defining the game is true at the point of instantiation and, vice versa, that if a formula is true at a point in a model, then Eve has a winning strategy in the corresponding game instantiated at that point.

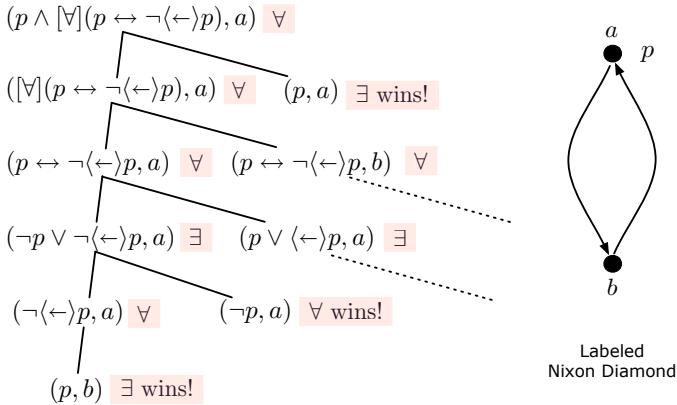
**Theorem 5 (Adequacy).** Let  $\varphi \in \mathcal{L}^{K^\forall}$ , and let  $\mathcal{M} = (\mathcal{A}, \mathcal{I})$  be an argumentation model. Then, for all  $a \in A$ :

$$(\varphi, a) \in Win_{\exists}(\mathcal{C}(\varphi, \mathcal{M})) \iff \mathcal{M}, a \models \varphi.$$

### 4.2 Games for Model-Checking Extensions

The next example illustrates a model-checking game for stable extensions run on the so-called Nixon diamond [16].





**Fig. 1.** Game for stable extensions in the 2-cycle with labeling (valuation) function

*Example 1 (Model-checking the Nixon diamond).* Let  $\mathcal{A} = (\{a, b\}, \{(a, b), (b, a)\})$  be an argumentation framework consisting of two arguments  $a$  and  $b$  attacking each other (i.e., the Nixon diamond), and consider the labeling  $\mathcal{I}$  assigning  $p$  to  $a$  and  $\neg p$  to  $b$  (top right corner of Figure 1). We now want to run an evaluation game for checking whether  $a$  belongs to a stable extension corresponding to the truth-set of  $p$ . Such game is the game  $\mathcal{C}(p \wedge Stable(p), (\mathcal{A}, \mathcal{I}))$  initialized at position  $(p \wedge Stable(p), a)$ . That is, spelling out the definition of  $Stable(p)$ :  $\mathcal{C}(p \wedge [\forall](p \leftrightarrow \neg\langle \leftarrow \rangle p)) @ (p \wedge [\forall](p \leftrightarrow \neg\langle \leftarrow \rangle p), a)$ . Such a game, played according to the rules in Definitions 4.2 and 6, gives rise to the tree in Figure 1.

In general, model-checking games provide a proof procedure for checking whether an argument belongs to a certain extension given an argumentation model. What must be noted is that the structure of such proof procedure is invariant, and the different games are obtained simply by choosing the right formula to be checked (Table 3).<sup>4</sup> This feature confers a high systematic flavor to this sort of games.

Now the natural question arises of what the precise relationship is between model-checking games and the sort of games studied in argumentation, sometimes called dialogue games [16, 14]. The difference is as follows.

In model-checking games you are given a model  $\mathcal{M} = (\mathcal{A}, \mathcal{I})$ , a formula  $\varphi$  and an argument  $a$ , and Eve is asked to prove that  $\mathcal{M}, a \models \varphi$ . In dialogue games, the check appointed to Eve is inherently more complex since the input consists only of an argumentation framework  $\mathcal{A}$ , a formula  $\varphi$  and an argument  $a$ . Eve is then asked to prove one of the two following things:

- that there exists a labeling function  $\mathcal{I}$  such that  $(\mathcal{A}, \mathcal{I}), a \models \varphi$  (the so-called *credulous* semantics for  $\varphi$ );
- that for all the labeling functions  $\mathcal{I}$ ,  $(\mathcal{A}, \mathcal{I}), a \models \varphi$  (the so-called *skeptical* semantics for  $\varphi$ ).

<sup>4</sup> Note that the game for checking grounded extensions is, obviously, the model-checking game for  $K^\mu$  [18].

**Table 3.** Games for model-checking extensions in argumentation models

---

$Adm : \mathcal{E}(\varphi \wedge Adm(\varphi), \mathcal{M})@(\varphi \wedge Adm(\varphi), a)$
$Complete : \mathcal{E}(\varphi \wedge Complete(\varphi), \mathcal{M})@(\varphi \wedge Complete(\varphi), a)$
$Stable : \mathcal{E}(\varphi \wedge Stable(\varphi), \mathcal{M})@(\varphi \wedge Stable(\varphi), a)$
$Grounded : \mathcal{E}(Grounded, \mathcal{M})@(Grounded, a)$

---

These are not a model-checking problems, but satisfiability problems in a pointed frame [3] which, in turn, are essentially model-checking problems in some fragment of monadic second-order logic. That is the problem of checking, given a frame  $\mathcal{A}$  and an argument  $a$ , whether the following is the case:

$$\begin{aligned} \mathcal{A} & \models \exists p_1, \dots, p_n ST_x(\varphi)[a] \\ \mathcal{A} & \models \forall p_1, \dots, p_n ST_x(\varphi)[a] \end{aligned}$$

where  $p_1, \dots, p_n$  are the atoms occurring in  $\varphi$  and  $ST_x(\varphi)[a]$  is the standard translation of  $\varphi$  realized in state  $a$  [5].

To conclude, we might say that the games defined in Section 4.1 provide a proof procedure for a reasoning task which is computationally simpler than the one tackled by standard dialogue games. It should be noted, however, that this is no intrinsic limitation to the logic-based approach advocated in the present paper. Model-checking games for monadic second-order logic (or rather for appropriate fragments of it) would be able to perform the sort of tasks demanded in dialogue games and do that in the same systematic manner of modal model-checking games. We will come back to this issue in Section 6.

## 5 Equivalent Arguments

Since abstract argumentation neglects the internal structure of arguments, the natural question arises of when two arguments can be said to be equivalent. Such a notion of equivalence will necessarily be of a structural nature. The study of a notion of equivalence for argumentation has not received attention yet by the argumentation theory community, except for one recent notable exception [15], which defines a notion of strong equivalence for argumentation frameworks, borrowed from the analogous notion developed in logic programming.

Modal logic offers a readily available notion of structural equivalence, the notion of bisimulation (with all its variants) [3, 11]. This section sketches the use of bisimulation for argumentation theoretic purposes. To illustrate the issue we use a simple motivating example depicted in Figure 2. We have two labelled argumentation frameworks which both contain an argument labeled  $p$  which is attacked by some arguments labelled  $q$ . Now the question would be: are the two  $p$ -arguments equivalent as far as abstract argumentation theory is concerned? The answer is yes, and the next sections explain why.

<sup>5</sup> For the definition of the standard translation we refer the reader to [3].

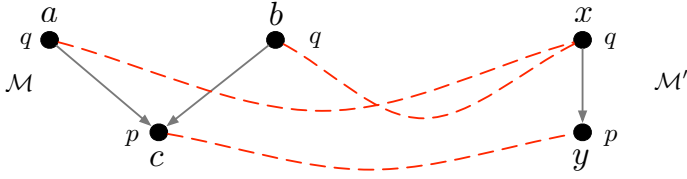


Fig. 2. Two (totally) bisimilar arguments ( $c$  and  $y$ ) in two argumentation models

### 5.1 Bisimilar Arguments

It is well-known that logic  $K^\mu$  is invariant under bisimulation [18]. In the present section we will focus on the specific notion of bisimulation which is tailored to  $K^\forall$ , also called *total bisimulation*.

**Definition 7 (Bisimulation).** Let  $\mathcal{M} = (A, \rightarrow, \mathcal{I})$  and  $\mathcal{M}' = (A', \rightarrow', \mathcal{I}')$  be two argumentation models. A bisimulation between  $\mathcal{M}$  and  $\mathcal{M}'$  is a non-empty relation  $Z \subseteq A \times A'$  such that for any  $a, a'$  s.t.  $aZa'$ : **Atom:**  $a$  and  $a'$  are propositionally equivalent; **Zig:** if  $a \rightarrow^{-1} b$  for some  $b \in A$ , then  $a' \rightarrow'^{-1} lb'$  for some  $b' \in A'$  and  $bZb'$ ; **Zag:** if  $a' \rightarrow'^{-1} b'$  for some  $b' \in A'$  then  $a \rightarrow^{-1} b$  for some  $b \in A$  and  $aZa'$ . A total bisimulation is a bisimulation  $Z \subseteq A \times A'$  such that its left projection covers  $A$  and its right projection covers  $A'$ . When a total bisimulation exists between  $\mathcal{M}$  and  $\mathcal{M}'$  we write  $(\mathcal{M}, a) \simeq (\mathcal{M}', a')$ .

Now, since logic  $K^\forall$  is invariant under total bisimulation [3] and logic  $K^\mu$  under bisimulation [11], we obtain a natural notion of equivalence of arguments, which is weaker than the notion of isomorphism of argumentation frameworks. If two arguments are equivalent in this perspective, then they are equivalent from the point of view of argumentation theory, as far as the notions expressible in those logics are concerned. In particular, we obtain the following simple theorem.

**Theorem 6 (Bisimilar arguments).** Let  $(\mathcal{M}, a)$  and  $(\mathcal{M}', a')$  be two pointed models, and let  $Z$  be a total bisimulation between  $\mathcal{M}$  and  $\mathcal{M}'$ . It holds that:

$$\begin{aligned} \mathcal{M}, a \models \text{Adm}(\varphi) \wedge \varphi &\iff \mathcal{M}', a' \models \text{Adm}(\varphi) \wedge \varphi \\ \mathcal{M}, a \models \text{CFree}(\varphi) \wedge \varphi &\iff \mathcal{M}', a' \models \text{CFree}(\varphi) \wedge \varphi \\ \mathcal{M}, a \models \text{Compl}(\varphi) \wedge \varphi &\iff \mathcal{M}', a' \models \text{Compl}(\varphi) \wedge \varphi \\ \mathcal{M}, a \models \text{Stable}(\varphi) \wedge \varphi &\iff \mathcal{M}', a' \models \text{Stable}(\varphi) \wedge \varphi \\ \mathcal{M}, a \models \text{Grounded} &\iff \mathcal{M}', a' \models \text{Grounded} \end{aligned}$$

In other words, Theorem 6 states that if two arguments are totally bisimilar, then they are indistinguishable from the point of view of abstract argumentation in the sense that the first belongs to a given conflict-free, or admissible set  $\varphi$  if and only if also the second does, and the first belongs to a given stable, complete extension  $\varphi$ , or to the grounded extension, if and only if also the second does. Arguments  $c$  and  $y$  in Figure 2 are totally bisimilar arguments.

**Table 4.** Rules of the bisimulation game

Position	Available moves
$((\mathcal{M}, a)(\mathcal{M}', a'))$	$\{((\mathcal{M}, a)(\mathcal{M}', b')) \mid \exists b' \in A' : a' \leftarrow b'\}$ $\cup\{((\mathcal{M}, b)(\mathcal{M}', a')) \mid \exists b \in A : a \leftarrow b\}$ $\cup\{((\mathcal{M}, a)(\mathcal{M}', b')) \mid \exists b' \in A'\}$ $\cup\{((\mathcal{M}, b)(\mathcal{M}', a')) \mid \exists b \in A\}$

## 5.2 Total Bisimulation Games

We can associate a game to Definition 7. Such game checks whether two given pointed models  $(\mathcal{M}, a)$  and  $(\mathcal{M}', a')$  are bisimilar or not. The game is played by two players: **Spoiler**, which tries to show that the two given pointed models are not bisimilar, and **Duplicator** which pursues the opposite goal. A match is started by **S**, then **D** responds, and so on. If and only if **D** moves to a position where the two pointed models are not propositionally equivalent, or if it cannot move any more, **S** wins.

**Definition 8 (Total bisimulation game).** Take two pointed models  $\mathcal{M}$  and  $\mathcal{M}'$ . The total bisimulation game  $\mathcal{B}(\mathcal{M}, \mathcal{M}')$  is defined by the following items. **Players:** The set of players is  $\{\mathbf{D}, \mathbf{S}\}$ . An element from  $\{\mathbf{D}, \mathbf{S}\}$  will be denoted  $P$  and its opponent  $\bar{P}$ . **Game form:** The game form of  $\mathcal{B}(\mathcal{M}, \mathcal{M}')$  is defined by Table 4. **Turn function:** If the round is even **S** plays, if it is odd **D** plays. **Winning conditions:** **S** wins if and only if either **D** has moved to a position  $((\mathcal{M}, a)(\mathcal{M}', a'))$  where  $a$  and  $a'$  do not satisfy the same labels, or **D** has no available moves. Otherwise **D** wins. **Instantiation:** The instance of  $\mathcal{B}(\mathcal{M}, \mathcal{M}')$  with starting position  $((\mathcal{M}, a)(\mathcal{M}', a'))$  is denoted  $\mathcal{B}(\mathcal{M}, \mathcal{M}')@ (a, a')$ .

So, as we might expect, positions in a (total) bisimulation game are pairs of pointed models, that is, the pointed models that **D** tries to show are bisimilar. It might also be instructive to notice that such a game can have infinite matches, which, according to Definition 8, are won by **D**.

From Definition 8 we obtain the following notions of winning strategies and winning positions.

**Definition 9 (Winning strategies and positions).** A strategy for player  $P$  in  $\mathcal{B}(\mathcal{M}, \mathcal{M}')@ (a, a')$  is a function telling  $P$  what to do in any match played from position  $(a, a')$ . Such a strategy is winning for  $P$  if and only if, in any match played according to the strategy,  $P$  wins. A position  $((\mathcal{M}, a)(\mathcal{M}', a'))$  in  $\mathcal{B}(\mathcal{M}, \mathcal{M}')$  is winning for  $P$  if and only if  $P$  has a winning strategy in  $\mathcal{B}(\mathcal{M}, \mathcal{M}')@ (a, a')$ . The set of all winning positions of game  $\mathcal{B}(\mathcal{M}, \mathcal{M}')$  for  $P$  is denoted by  $Win_P(\mathcal{B}(\mathcal{M}, \mathcal{M}'))$ .

We have the following adequacy theorem. The proof is standard and the reader is referred to [11].

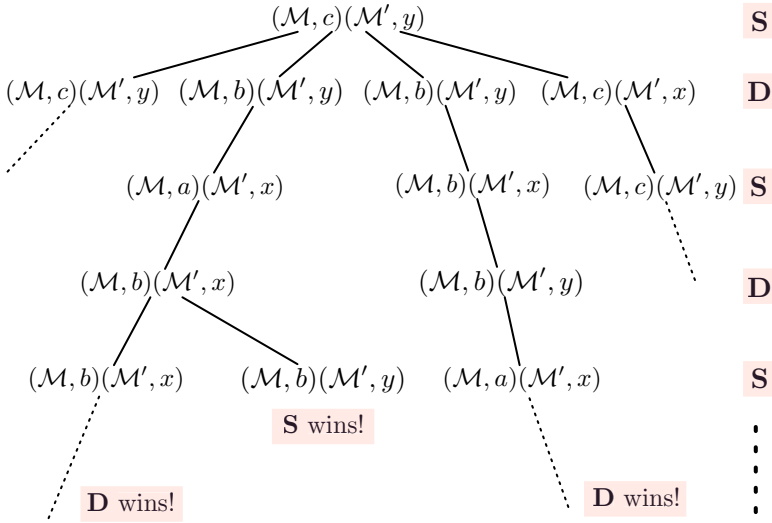


Fig. 3. Part of the total bisimulation game played on the models in Figure 2

**Theorem 7 (Adequacy).** Take  $(\mathcal{M}, a)$  and  $(\mathcal{M}', a')$  to be two argumentation models. It holds that:

$$((\mathcal{M}, a)(\mathcal{M}', a')) \in Win_{\mathbf{D}}(\mathcal{B}(\mathcal{M}, \mathcal{M}')) \iff (\mathcal{M}, a) \rightleftharpoons (\mathcal{M}', a').$$

In words, **D** has a winning strategy in the game  $\mathcal{B}(\mathcal{M}, \mathcal{M}')@(\mathcal{M}, a, \mathcal{M}', a')$  if and only if  $\mathcal{M}, a$  and  $\mathcal{M}', a'$  are totally bisimilar. An example of such a game follows.

*Example 2 (A total bisimulation game).* Let us play a total bisimulation game on the two models  $\mathcal{M}$  and  $\mathcal{M}'$  given in Figure 2. A total bisimulation connects  $c$  with  $y$ , and  $a$  and  $b$  with  $x$ . Part of the extensive bisimulation game  $\mathcal{B}(\mathcal{M}, \mathcal{M}')@(\mathcal{M}, c, \mathcal{M}', y)$  is depicted in Figure 3. Notice that **D** wins on those infinite paths where it can always duplicate **S**'s moves. On the other hand, it loses for instance when it replies to one of **S**'s moves  $((\mathcal{M}, b)(\mathcal{M}', x))$  by moving in the second model to argument  $y$ , which is labelled  $p$  while  $b$  is not.

## 6 A Research Agenda between Logic and Argumentation

By recapitulating results presented in [12,13], the paper has given a glimpse of the sort of results that can be obtained about abstract argumentation theory by resorting to quite standard methods and techniques of modal logic. The present section proposes an agenda for this line of research which, in the author's view, is of definite interest for a deeper mathematical understanding of abstract argumentation theory.

### 6.1 Other Extension-Based Semantics in Modal Logic

The paper has left aside one key notion of argumentation: preferred extensions. In [7], preferred extensions are defined as maximal, with respect to set-inclusion, complete extensions. The natural question is whether the logics we have introduced are expressive enough to capture also this notion.

Technically, this means looking for a formula  $\varphi(p)$  such that for any pointed model  $\mathcal{M} = ((\mathcal{A}, \mathcal{I}), a)$   $\mathcal{M}, a \models \varphi(p)$  iff  $a \in |p|_{\mathcal{M}}$  and  $|p|_{\mathcal{M}}$  is a preferred extension of  $\mathcal{A}$ , where  $p \in \mathbf{P}$ . It is easy to see that such  $\varphi(p)$  can be expressed in monadic second-order logic with a  $\Pi_1^1$  quantification:

$$p \wedge ST_x(Compl(p)) \wedge \forall q(ST_x(Compl(q)) \rightarrow \neg(p \sqsubseteq q)) \tag{9}$$

where  $ST_x(Compl(p))$  denotes the standard translation [3] of the  $K^\forall$  formula for complete extensions (Formula 4) and  $q \sqsubseteq p$  means just that  $|q|_{\mathcal{M}} \subseteq |p|_{\mathcal{M}}$ , i.e., the truth set of  $q$  is included in the truth-set of  $p$ . The same question of representability within (possibly extended) modal languages can be posed for other types of extensions, such as the semi-stable one [4].

### 6.2 A Unified Game-Theoretic Proof-Theory for Argumentation

Section 4.2 has shown how model-checking games can be used to provide a form of game-theoretic proof theory to check the membership of a given argument to a given extension. Although Section 4.2 has then pointed out how these games differ from the standard dialogue games studied in argumentation theory, it is our thesis that a suitable extension of the expressivity of the modal languages used in this paper can offer a unified game-theoretic proof-theory for argumentation.

For instance, the question whether there exists, given a pointed frame  $(\mathcal{A}, a)$ , a stable extension of  $\mathcal{A}$  containing  $a$  could be phrased as the model checking of a formula of the extension of  $K^\forall$  with second order quantification limited to alternation depth 1:6

$$\mathcal{A}, a \models \exists p.(Stable(p) \wedge p).$$

The prospect of an extension of this type is to provide each argumentation-theoretic notion with a game-theoretic proof-theory (both for its skeptical and credulous versions) which would directly follow from the model-checking game of the underlying logic. We would thereby obtain also games for extensions which have not yet found a game-theoretic proof-theory in the literature on abstract argumentation theory such as, for instance, skeptical and credulous stable extensions, or skeptical preferred extensions.

### 6.3 Equivalence in Argumentation

Another original application of modal logic which could open up new venues for research is the study of invariance, or equivalence, in argumentation theory. Section 5 has shown how to tackle the question of when two labelled argumentation

---

<sup>6</sup> A well-known logical language for this purpose could be second order propositional modal logic [9].

frameworks can be considered equivalent, by looking at the existence of a (total) bisimulation relation between them.

The key observation in this case is that, depending on the features we consider relevant for the comparison of two argumentation frameworks, different modal languages can be chosen, which come with their characteristic notion of bisimulation, i.e., structural equivalence. For instance, if we were to compare two argumentation frameworks by considering, as a relevant property for the comparison, also the number of attackers, then the two arguments considered equivalent in Figure 2 would cease to be such, as the first one has two attackers, while the second has only one.

The modal language with the sort of expressivity necessary to ‘count’ the number of attackers of a given argument is called *graded modal logic* [10]. In such a language it becomes possible to say that:

$$\mathcal{A}, a \models \diamond_2 \top$$

that is,  $a$  has at least two attackers. Going back to the example given in Figure 2, while argument  $c$  in the first framework satisfies  $\diamond_2 \top$ , argument  $y$  in the second does not. In modal logic terms this implies that  $c$  and  $y$  are not bisimilar with respect to the language of graded modal logic or, put it otherwise, they are not *graded bisimilar* [17]. So, mapping all the relevant modal languages for argumentation theory would automatically provide a whole landscape of different equivalence notions which can be used to compare argumentation frameworks.

## 6.4 Argumentation Dynamics

The whole of abstract argumentation theory is built on structures—the argumentation frameworks—which are essentially static. To date, no theory has yet been systematically developed about how to modify argumentation frameworks by operations of addition and deletion of arguments and links.

The link with modal logic could offer again a wealth of techniques, stemming from dynamic logic [6,2], which might prove themselves useful for the development of such a theory of argumentation dynamics. The possibly simplest example in this line is sabotage modal logic [1], where formulae of the type:

$$(\mathcal{A}, \mathcal{I}), a \models \blacksquare \varphi$$

express, in an argumentation-theoretic reading, that after any possible removal of an attack relation, argument  $a$  still belongs to the truth-set of  $\varphi$ .

## 7 Conclusions

The paper has shown how rather standard modal logics—the extensions of  $K$  with universal modality and least fixpoint operator—can be applied to argumentation theory in an almost direct way. Both these logics come equipped with complete calculi, in which, therefore, theorems from argumentation theory can

be formally derived, with model-checking games, which can be used to provide a game-theoretic proof-theory on argumentation models, and with characteristic notions of structural equivalence (bisimulation) which can be used to provide a formalization of notions of equivalence for argumentation frameworks.

We have concluded by pointing at several directions for future work, ranging from the problem of the formalization of preferred extensions, to second-order model checking games, to the study of argumentation equivalence via bisimulation, and to argumentation dynamics.

**Acknowledgments.** This work is sponsored by the *Nederlandse Organisatie voor Wetenschappelijk Onderzoek* (NWO) under the VENI grant 639.021.816. The author wishes to thank Sanjay Modgil for the inspiring conversation that sparked this study.

## References

1. van Benthem, J.: An essay on sabotage and obstruction. In: Hutter, D., Stephan, W. (eds.) *Mechanizing Mathematical Reasoning*. LNCS (LNAI), vol. 2605, pp. 268–276. Springer, Heidelberg (2005)
2. van Benthem, J.: *Logical Dynamics of Information and Interaction*. Cambridge University Press, Cambridge (forthcoming)
3. Blackburn, P., de Rijke, M., Venema, Y.: *Modal Logic*. Cambridge University Press, Cambridge (2001)
4. Caminada, M.: Semi-stable semantics. In: Dunne, P.E., Bench-Capon, T. (eds.) *Proceedings of Computational Models of Argument, COMMA 2006*, pp. 121–130 (2006)
5. Davey, B.A., Priestley, H.A.: *Introduction to Lattices and Order*. Cambridge University Press, Cambridge (1990)
6. van Ditmarsch, H., Kooi, B., van der Hoek, W.: *Dynamic Epistemic Logic*. Synthese Library Series, vol. 337. Springer, Heidelberg (2007)
7. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artificial Intelligence* 77(2), 321–358 (1995)
8. Dunne, P., Bench-Capon, T.: *Complexity and combinatorial properties of argument systems*. Technical report, University of Liverpool (2001)
9. Fine, K.: Propositional quantifiers in modal logic. *Theoria* 36, 336–346 (1936)
10. Fine, K.: In so many possible worlds. *Notre Dame Journal of Formal Logic* 13, 516–520 (1972)
11. Goranko, V., Otto, M.: Model theory of modal logic. In: Blackburn, P., van Benthem, J., Wolter, F. (eds.) *Handbook of Modal Logic*. Studies in Logic and Practical Reasoning, vol. 3, pp. 249–329. Elsevier, Amsterdam (2007)
12. Grossi, D.: *Doing argumentation theory in modal logic*. ILLC Prepublication Series PP-2009-24, Institute for Logic, Language and Computation (2009)
13. Grossi, D.: On the logic of argumentation theory. In: van der Hoek, W., Kaminka, G., Lespérance, Y., Sen, S. (eds.) *Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2010)*, IFAAMAS, pp. 409–416 (2010)



14. Modgil, S., Caminada, M.: Proof theories and algorithms for abstract argumentation frameworks. In: Rahwan, Y., Simari, G. (eds.) *Argumentation in AI*. Springer, Heidelberg (2009)
15. Oikarinen, E., Woltran, S.: Characterizing strong equivalence for argumentation frameworks. In: Lin, F., Sattler, U., Truszczynski, M. (eds.) *Principles of Knowledge Representation and Reasoning: Proceedings of the Twelfth International Conference (KR 2010)*, May 9–13. AAAI Press, Toronto (2010)
16. Prakken, H., Vreeswijk, G.: Logics for defeasible argumentation. In: *Handbook of Philosophical Logic*, 2nd edn., vol. IV, pp. 218–319 (2002)
17. de Rijke, M.: A note on graded modal logic. *Studia Logica* 64(2), 271–283 (2000)
18. Venema, Y.: *Lectures on the modal  $\mu$ -calculus*. Renmin University in Beijing, China (2008)
19. Walukiewicz, I.: Completeness of Kozen’s axiomatization of the propositional  $\mu$ -calculus. *Information and Computation* 157, 142–182 (2000)
20. Zermelo, E.: Über eine anwendung der mengenlehre auf die theorie des schachspiels. In: *Proceedings of the 5th Congress Mathematicians*, pp. 501–504. Cambridge University Press, Cambridge (1913)

## Appendix: A Formal Proof of the *Fundamental Lemma*

- |     |  |                        |
|-----|--|------------------------|
| 1.  | $\varphi \rightarrow \varphi \vee \psi$  | <b>Prop</b>            |
| 2.  | $\langle \leftarrow \rangle \varphi \rightarrow \langle \leftarrow \rangle (\varphi \vee \psi)$  | 1, K – derived rule    |
| 3.  | $[\leftarrow] \langle \leftarrow \rangle \varphi \rightarrow [\leftarrow] \langle \leftarrow \rangle (\varphi \vee \psi)$  | 2, K – derived rule    |
| 4.  | $(\alpha \vee \beta \rightarrow \gamma) \rightarrow (\beta \rightarrow \gamma)$  | <b>Prop</b>            |
| 5.  | $(\psi \vee \xi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi) \rightarrow (\xi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi)$  | 4, instance            |
| 6.  | $(\psi \vee \xi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi) \rightarrow (\xi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi \vee \psi)$  | 5, 3, <b>Prop, MP</b>  |
| 7.  | $[\forall](\psi \vee \xi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi) \rightarrow [\forall](\xi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi \vee \psi)$  | 6, K – derived rule    |
| 8.  | $Acc(\psi \vee \xi, \varphi) \rightarrow Acc(\xi, \varphi \vee \psi)$  | 7, definition          |
| 9.  | $(\psi \vee \xi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi) \rightarrow (\psi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi)$   | 4, instance            |
| 10. | $[\forall](\psi \vee \xi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi) \rightarrow [\forall](\psi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi)$   | 9, K – derived rule    |
| 11. | $Acc(\psi \vee \xi, \varphi) \rightarrow Acc(\psi, \varphi)$   | 10, definition         |
| 12. | $((\alpha \rightarrow \gamma) \wedge (\beta \rightarrow \gamma)) \rightarrow (\alpha \vee \beta \rightarrow \gamma)$   | <b>Prop</b>            |
| 13. | $([\forall](\alpha \rightarrow \gamma) \wedge [\forall](\beta \rightarrow \gamma)) \rightarrow [\forall](\alpha \vee \beta \rightarrow \gamma)$  | 12, <b>N, K, MP</b>    |
| 14. | $([\forall](\varphi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi) \wedge [\forall](\psi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi)) \rightarrow [\forall](\varphi \vee \psi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi)$           | 13, Instance           |
| 15. | $[\leftarrow] \langle \leftarrow \rangle \varphi \rightarrow [\leftarrow] \langle \leftarrow \rangle (\varphi \vee \psi)$  | 14, <b>Prop, K, N</b>  |
| 16. | $([\forall](\varphi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi) \wedge [\forall](\psi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi)) \rightarrow [\forall](\varphi \vee \psi \rightarrow [\leftarrow] \langle \leftarrow \rangle \varphi \vee \psi)$ | 15, <b>Prop, K, N</b>  |
| 17. | $Acc(\varphi, \varphi) \wedge Acc(\psi, \varphi) \rightarrow Acc(\varphi \vee \psi, \varphi \vee \psi)$  | 16, definition         |
| 18. | $Acc(\varphi, \varphi) \wedge Acc(\psi \vee \xi, \varphi) \rightarrow Acc(\varphi \vee \psi, \varphi \vee \psi)$   | 17, 9, <b>Prop, MP</b> |
| 19. | $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi) \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi)$   | <b>Incl</b>            |

20.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi) \rightarrow ([\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow [\leftarrow]\neg \varphi)$  19, **Prop**, **MP**
21.  $[\forall][\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi) \rightarrow [\forall]([\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow [\leftarrow]\neg \varphi)$  20, **K** – derived rule
22.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi) \rightarrow [\forall]([\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow [\leftarrow]\neg \varphi)$  21, **S5** – derived rule
23.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi) \wedge [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi)$   
 $\rightarrow [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi) \wedge [\forall]([\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow [\leftarrow]\neg \varphi)$  22, **Prop**, **MP**
24.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi) \wedge [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi) \rightarrow [\forall](\varphi \vee \psi \rightarrow [\leftarrow]\neg \varphi)$  23, **Prop**, **MP**
25.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi \wedge \neg \psi) \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi \wedge \neg \psi)$  **Incl**
26.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi \wedge \neg \psi) \rightarrow ([\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow [\leftarrow]\neg \varphi \wedge \neg \psi)$  25, **K**, **Prop**, **MP**
27.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi \wedge \neg \psi) \rightarrow [\forall]([\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow [\leftarrow]\neg \varphi \wedge \neg \psi)$  26, **S5** – derived rule
28.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi) \wedge [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi)$   
 $\rightarrow [\forall]([\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow [\leftarrow]\neg \varphi \wedge \neg \psi)$  24, 27, **Prop**, **MP**
29.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi) \wedge [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi)$   
 $\rightarrow [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi) \wedge [\forall]([\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow [\leftarrow]\neg \varphi \wedge \neg \psi)$  28, **Prop**, **MP**
30.  $[\forall](\alpha \rightarrow \beta) \wedge [\forall](\beta \rightarrow \gamma) \rightarrow [\forall](\alpha \rightarrow \gamma)$  **S5** – theorem
31.  $[\forall]([\leftarrow](\langle \leftarrow \rangle \varphi \rightarrow [\leftarrow](\neg \varphi \wedge \neg \psi)) \wedge [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi)$   
 $\rightarrow [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\neg \varphi \wedge \neg \psi))$  30, instance
32.  $[\forall](\langle \leftarrow \rangle \varphi \rightarrow \neg \varphi) \wedge [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\langle \leftarrow \rangle \varphi)$   
 $\rightarrow [\forall](\varphi \vee \psi \rightarrow [\leftarrow](\neg \varphi \wedge \neg \psi))$  29, 31, **Prop**, **MP**
- 33.**  $CFree(\varphi) \wedge Acc(\varphi \vee \psi, \varphi) \rightarrow CFree(\varphi \vee \psi)$  32, definition
34.  $Acc(\varphi, \varphi) \wedge Acc(\psi, \varphi) \rightarrow Acc(\varphi \vee \psi, \varphi)$  14, definition
35.  $CFree(\varphi) \wedge Acc(\varphi, \varphi) \wedge Acc(\psi, \varphi) \rightarrow CFree(\varphi \vee \psi)$  33, 34, **Prop**, **MP**
36.  $CFree(\varphi) \wedge Acc(\varphi, \varphi) \wedge Acc(\psi \vee \xi, \varphi) \rightarrow CFree(\varphi \vee \psi)$  35, 9, **Prop**, **MP**
37.  $CFree(\varphi) \wedge Acc(\varphi, \varphi) \wedge Acc(\psi \vee \xi, \varphi)$   
 $\rightarrow CFree(\varphi \vee \psi) \wedge Acc(\varphi \vee \psi, \varphi \vee \psi)$  36, 18, **Prop**, **MP**
38.  $CFree(\varphi) \wedge Acc(\varphi, \varphi) \wedge Acc(\psi \vee \xi, \varphi)$   
 $\rightarrow CFree(\varphi \vee \psi) \wedge Acc(\varphi \vee \psi, \varphi \vee \psi) \wedge Acc(\xi, \varphi \vee \psi)$  37, 8, **Prop**, **MP**
- 39.**  $Adm(\varphi) \wedge Acc(\psi \vee \xi, \varphi) \rightarrow Adm(\varphi \vee \psi)Acc(\xi, \varphi \vee \psi)$  38, definition

# Towards Pragmatic Argumentative Agents within a Fuzzy Description Logic Framework

Ioan Alfred Letia and Adrian Groza

Technical University of Cluj-Napoca  
Department of Computer Science  
Baritiu 28, RO-400391 Cluj-Napoca, Romania  
{letia,adrian}@cs-gw.utcluj.ro

**Abstract.** To bring the level of current argumentation to the expressive and flexible status expected by human agents, we introduce fuzzy reasoning on top of the classical Dung abstract argumentation framework. The system is built around Fuzzy Description Logic and exploits the integration of ontologies with argumentation theory, attaining the advantage of facilitating communication of domain knowledge between agents. The formal properties of fuzzy relations are used to provide semantics to the different types of conflicts and supporting roles in the argumentation. The usefulness of the framework is illustrated in a supply chain scenario.

## 1 Introduction

Abstract argumentation frameworks lack high-level conveniences such as ease of understanding, an aspect required by human agents. Many challenges still exist in order to build intelligent systems based on abstract argumentation frameworks [2].

On the one hand, humans manifest a lot of flexibility when they convey arguments from supporting and attacking a claim. One can disagree, can provide a counter example, can rebut or undercut a claim. Currently, these common patterns of attacking relations are encapsulated as argumentation schemes [18]. This informal reasoning does not exploit the formal properties of the attacking relations. The semantics of the support relation *agree* contains the transitivity property: *agree*(*a*, *b*) and *agree*(*b*, *c*) implies that *a* agrees with *c*. Similarly, the rebutting relation is symmetrical. That is, *rebut*(*a*,  $\neg a$ ) implies *rebut*( $\neg a$ , *a*). In this paper, we advocate to use such properties when deciding on the status of an argument. We provide software agents with description logic based reasoning capabilities to exploit the formal properties of the attacking relations.

On the other hand, people do not express their arguments precisely in their daily life. Such vague notions as: strongly, moderately, don't fully agree, tend to disagree are used during an argumentative dialog. Real arguments are also a mixture of fuzzy linguistic variables and ontological knowledge. Arguments conveyed by people are incomplete, normally *enthymemes* [11], where the opponent of the arguments assumes that his partner understands the missing part. Thus, a common knowledge on the debate domain is assumed by the agents.

**Table 1.** Operators in Fuzzy Logics

Operation	Lukasiewicz Logic	Gödel Logic
intersection $\alpha \otimes_S \beta$	$\max\{\alpha + \beta - 1, 0\}$	$\min\{\alpha, \beta\}$
union $\alpha \oplus_S \beta$	$\min\{\alpha + \beta, 1\}$	$\max\{\alpha, \beta\}$
negation $\ominus_S \alpha$	$1 - \alpha$	1, if $\alpha = 0$ , 0, otherwise
implication $\alpha \Rightarrow_S \beta$	$\min\{1, 1 - \alpha + \beta\}$	1, if $\alpha \leq \beta$ , $\beta$ , otherwise

We introduce Fuzzy Description Logic on top of the argumentation theory, as the adequate technical instrumentation needed to model real-life debates.

## 2 Preliminaries

### 2.1 Fuzzy Sets and Relations

A fuzzy relation  $R$  between two set  $A$  and  $B$  has degree of membership whose value lies in  $[0, 1]$ :  $\mu_R : A \times B \rightarrow [0, 1]$ .  $\mu_R(x, y)$  is interpreted as strength of relation  $R$  between  $x$  and  $y$ . When  $\mu_R(x, y) \geq \mu_R(x', y')$ ,  $(x, y)$  is more strongly related than  $(x', y')$ . A fuzzy relation  $R$  over  $X \times X$  is called:

- *transitive*:  $\forall a, b, c \in X, R(a, c) \geq \sup_{b \in X} \{\otimes_S(R(a, b), R(b, c))\}$
- *reflexive*:  $\forall a \in X, R(a, a) = 1$
- *irreflexive*:  $\exists a \in X, R(a, a) \neq 1$
- *antireflexive*:  $\forall a \in X, R(a, a) \neq 1$
- *symmetric*:  $\forall a, b \in X, R(a, b) \rightarrow R(b, a)$
- *antisymmetric*:  $\forall a, b \in X, R(a, b) \rightarrow \neg R(b, a)$
- *disjoint*:  $\forall a, b \in X, \otimes_S(R(a, b), S(a, b)) = 0$

The *inverse* of a fuzzy relation  $R \subseteq X \times Y$  is a fuzzy relation  $R^- \subseteq Y \times X$  defined as  $R^-(b, a) = R(a, b)$ . Given two fuzzy relations  $R_1 \subseteq X \times Y$  and  $R_2 \subseteq Y \times Z$  we define the *composition* as  $[R_1 \circ R_2](a, c) = \sup_{b \in Y} \{\otimes_S(R(a, b), R(b, c))\}$  (table 1). The composition satisfies the following properties:  $(R_1 \circ R_2) \circ R_3 = R_1 \circ (R_2 \circ R_3)$ , and  $(R_1 \circ R_2)^- = (R_2^- \circ R_1^-)$ . Due to the associativity property we can extend the composition operation to any number of fuzzy relations:  $[R_1 \circ^t R_2 \circ^t \dots \circ^t R_n](a, b)$ . If a relation is reflexive, antisymmetric, and transitive it is called *order relation*.

### 2.2 Fuzzy Description Logic

In the following paragraphs the differences introduced by fuzzy reasoning on top of classical description logic are presented. The complete formalization of the fuzzy description logic can be found in [4]. The *syntax* of fuzzy SHIF concepts [4] is as follows:

$$\begin{aligned}
 C, D &= \top \mid \perp \mid A \mid C \sqcap_S D \mid C \sqcup_S D \mid C \sqsubseteq_S D \mid \neg_L C \mid \\
 &\quad \forall R.C \mid \exists R.C \mid \forall T.d \mid \exists T.d \mid \leq nR \mid \geq nR \mid m(C) \mid \{a_1, \dots, a_n\} \\
 d &= \text{crisp}(a, b) \mid \text{triangular}(a, b, c) \mid \text{trapezoidal}(a, b, c, d)
 \end{aligned}$$

$\perp^I(x) = 0$	$(\forall R.C)^I(x) = \inf_{y \in \Delta^I} R^I(x, y) \Rightarrow_S C^I(y)$
$\top^I(x) = 1$	$(\exists R.C)^I(x) = \sup_{y \in \Delta^I} R^I(x, y) \otimes_S C^I(y)$
$(\neg C)^I = \ominus C^I(x)$	$(\forall T.d)^I(x) = \inf_{y \in \Delta^I} R^I(x, y) \Rightarrow_S d^I(y)$
$(C \sqcap_S D)^I(x) = C^I(x) \otimes_S D^I(x)$	$(\exists R.d)^I(x) = \sup_{y \in \Delta^I} R^I(x, y) \otimes_S d^I(y)$
$(C \sqcup_S D)^I(x) = C^I(x) \oplus_S D^I(x)$	$(x : C)^I = C^I(x^I)$
$(C \rightarrow_S D)^I(x) = C^I(x) \Rightarrow_S D^I(x)$	$((x, y) : R)^I = R^I(x^I, y^I)$
$(m(C))^I(x) = f_m(C^I(x))$	$(C \sqsubseteq D)^I(x) = \inf_{x \in \Delta^I} C^I(x) \Rightarrow_S D^I(x)$

**Fig. 1.** Semantics of fuzzy concepts

where  $S = \{L, G, C\}$ ,  $L$  comes from Lukasiewicz semantics,  $G$  from Gödel semantics, and  $C$  stands for classical logic (see table [II](#)). The modifier  $m(C) = \text{linear}(a) \mid \text{triangular}(a, b, c)$  can be used to alter the membership functions of the fuzzy concepts. Fuzzy modifiers such as *very*, *more-or-less*, *slightly* can be applied to fuzzy sets to change their membership functions. They are defined in terms of linear hedges. For instance, one can define  $\text{very} = \text{linear}(0.8)$ . A functional role  $S$  can always be obtained by means of the axiom  $\top \sqsubseteq (\leq 1S)$ .

*Example 1.* The definition of junk food is applied to some food which has little nutritional value, or to products with nutritional value but which also have ingredients considered unhealthy:  $\text{JunkFood} = \text{Food} \sqcap (\exists \text{hasNutritionalValue.Little} \sqcup \exists \text{hasIngredients.Unhealthy})$ . In this definition, there are two roles which point to the fuzzy concepts *Little* and *Unhealthy*, which could be represented as  $\text{Little} = \text{triangular}(10, 20, 30)$ , or  $\text{Unhealthy} = \exists \text{hasSalt.} \geq 2\text{mg} \sqcup \text{hasAdditive.} > 0.5\text{mg}$ .

The main idea of *semantics* of FDL is that concepts and roles are interpreted as fuzzy subsets of an interpretation's domain [4](#). A fuzzy interpretation  $I = (\Delta^I, \bullet^I)$  consists of a non empty set  $\Delta^I$  (the domain) and a fuzzy interpretation function  $\bullet^I$ . The mapping  $\bullet^I$  is extended to roles and complex concepts as specified in figure [II](#).

### 3 Fuzzy Argumentation Systems

#### 3.1 Fuzzy Resolution Argumentation Base

An argumentation framework [7](#) consists of a set of arguments, some of which attack each other. In our approach, the arguments represent instances of concepts, while different types of attack relations are instantiations of roles defined on these concepts. Both, the concepts and the roles can be fuzzy.

**Definition 1.** A fuzzy resolution argumentation base is a tuple  $FRA = \langle A, \mathcal{J}, \mathcal{R} \rangle$ , consisting of a fuzzy Abox  $A$ , representing argument instances and their attacking relations, a fuzzy Tbox  $\mathcal{J}$  representing concepts, and a fuzzy Rbox  $\mathcal{R}$  encapsulating attack-like and support-like relations.

**Definition 2.** A fuzzy Abox  $A$  is a tuple  $\prec \text{Arg}, \text{Attacks} \succ$ , where  $\text{Arg}$  of a finite set of assertion axioms for fuzzy arguments  $\{a_1 : C_1 \bowtie \alpha_1, a_2 : C_2 \bowtie$

$\alpha_2, \dots, a_n : C_n, \bowtie \alpha_n\}$ , and *Attacks* is a set of fuzzy roles  $\subseteq \text{Arg} \times \text{Arg}$  of the form  $\{(a_i, a_j) : R_k \bowtie \alpha_l\}$ , where  $\alpha_l \in [0, 1]$ ,  $C_i$  are concepts,  $R_k$  are attack and support-like relations, and  $\bowtie = \{<, \leq, >, \geq\}$ .

*Example 2.* Let  $\mathcal{A} = \prec \{\text{funghi} : \text{CheapPizza} \geq 0.8\}, \{(\text{funghi}, \text{vegetarian}) : \text{Attack} \geq 0.7\} \succ$  states that *funghi* is a *CheapPizza* with degree at least 0.8, and it attacks the *vegetarian* argument with degree at least 0.7.

If  $\alpha$  is omitted, the maximum degree of 1 is assumed. We use  $\bowtie^-$  as the reflection of inequalities  $\leq^- = \geq$  and  $<^- = >$ .

**Definition 3.** A fuzzy Tbox  $\mathcal{T}$  is a finite set of inclusion axioms  $\{C_i \sqsubseteq_S D_i, \geq \alpha_i\}$ , where  $\alpha_i \in [0, 1]$ ,  $C_i, D_i$  are concepts, and  $S$  specifies the implication function (Lukasiewicz, Gödel) to be used. The axioms state that the subsumption degree between  $C$  and  $D$  is at least  $\alpha$ .

*Example 3.* Let’s take the common example of pizza. Can it be categorized as junk food or nutrition food? Associated with some food outlets, it is labeled as ”junk”, while in others it is seen as being acceptable and trendy. Rather, one can consider that it belongs to both concepts with different degree of truth, let’s say 0.7 for *JunkFood* and 0.5 to *NutritionFood*, encapsulated as  $\mathcal{T} = \{Pizza \sqsubseteq_L JunkFood \geq 0.7, Pizza \sqsubseteq_L NutritionalFood \geq 0.5, FreshFruits \sqsubseteq_L NutritionalFood, CandyBar \sqsubseteq_L JunkFood\}$ . Note the subconcept *CandyBar* is subsumed by the concept *JunkFood* with a degree of 1.

**Definition 4.** The argumentation core  $\mathcal{R}^k$  of the fuzzy Rbox  $\mathcal{R}$  consists of two relations *Attack* and *Support* (noted by  $\bar{A}$ , respectively  $\bar{S}$ ), having the property  $dis(\bar{A}, \bar{S})$ , meaning that  $\forall a, b \in \text{Arg}, \otimes_S((a, b) : \bar{A}, (a, b) : \bar{S}) = 0$ . Formally,  $\mathcal{R}^k = \{\bar{A}, \bar{S}, dis(\bar{A}, \bar{S})\}$ .

Under the Gödel semantics, the disjoint property of the *Attack* and *Support* roles states that  $\otimes_G((a, b) : \bar{A}, (a, b) : \bar{S}) = \min((a, b) : \bar{A}, (a, b) : \bar{S}) = 0 \Leftrightarrow$  if  $(a, b) : \bar{S} \geq 0$  then  $(a, b) : \bar{A} \leq 0$  and if  $(a, b) : \bar{A} \geq 0$  then  $(a, b) : \bar{S} \leq 0$ . In other words, if  $a$  attacks  $b$  there is no support relation from  $a$  to  $b$ , and similarly if  $a$  supports  $b$  there is no attack relation from  $a$  to  $b$ . The Lukasiewicz semantics leads to a more flexible interpretation, given by  $\otimes_L((a, b) : \bar{A}, (a, b) : \bar{S}) = \max((a, b) : \bar{A} + (a, b) : \bar{S} - 1, 0) = 0 \Leftrightarrow (a, b) : \bar{A} + (a, b) : \bar{S} \leq 1$ . Thus, if  $a$  attacks  $b$  to a certain degree  $\alpha$ , there exists the possibility that also  $a$  supports  $b$  with a maximum degree of  $1 - \alpha$ . While the Gödel interpretation fits perfectly to the general case of argumentative debates, some special examples lay under the Lukasiewicz semantics.

*Example 4.* Several government strategies focus on decreasing expenses with personal in order to re-allocate funds to investments, hoping to support the growth of the economy. Two relations of type support are used to express this:  $(\text{decreaseExpenses}, \text{largerInvestments}) : \bar{S}$ , respectively  $(\text{largerInvestments}, \text{growth}) : \bar{S}$ . A different chain of reasoning follows by the fact the decreasing

salaries leads to lowering the consumption which threatens the growth of the economy, formalised by a support role: (*decreaseExpenses, smallerConsumption*) : $\bar{S}$  and an attack-like relation (*smallerConsumption, growth*) : *threat*, where is a special type of attacking role (*threat*  $\sqsubseteq$   $\bar{A}$ ). Consequently, the strategy of decreasing salaries supports with a degree  $\alpha$ , but also attacks with a degree  $\beta$ , the objective of growing the economy. This is permitted in a Lukasiewicz setting, to an inconsistency budget of  $\alpha + \beta < 1$ .

**Definition 5.** *The fuzzy Rbox  $\mathcal{R}$  consists of i) the argumentation core  $\mathcal{R}^k$ ; ii) an hierarchy of disjoint attack and support-like relations  $R$ , defined by role inclusion axioms:  $R \sqsubseteq \text{Attack}$ , or  $R \sqsubseteq \text{Support}$ ; and iii) a set of role assertions of the form: (*fun*  $R$ ), (*trans*  $R$ ), *sym*( $R$ ), (*inv*  $R$   $R^-$ ), stating that the role  $R$  is functional, transitive, symmetric, respectively its inverse relation is  $R^-$ .*

There are two types of relations in the set  $\mathcal{R}$ : *supporting roles* (denoted by  $\mathcal{R}^S$ ), opposite to *attacking roles* (denoted by  $\mathcal{R}^A$ ), where  $\mathcal{R}^S \cap \mathcal{R}^A = \emptyset$ . We note that  $a_1$  supports  $a_2$  by  $a_1 \rightarrow a_2$  and  $a_1$  attacks  $a_2$  by  $a_1 \rightarrow a_2$ .

*Example 5.* Let  $\mathcal{R} = \mathcal{R}^k \cup \{\text{Defeat, Disagree, Agree, Defeat} \sqsubseteq \bar{A}, \text{Agree} \sqsubseteq \bar{S}, \text{Disagree} \sqsubseteq \bar{A}, \text{tra}(\text{Agree}), \text{sym}(\text{Agree}), \text{ref}(\text{Agree})\}$ . The two hierarchies are  $\mathcal{R}^A = \{\bar{A}, \text{Defeat, Disagree}\}$ , respectively  $\mathcal{R}^S = \{\bar{S}, \text{Agree}\}$ . Note that *Support* relation is transitive, while *Attack* role is not a transitive one; *Agree* is a particular instance of *Support* relation, while *Disagree* and *Defeat* relations are *Attack*-type relations. The following properties hold:

**Proposition 1.** *Attack and Support-like relations are not functional, i.e. the same argument  $a$  can attack two different arguments  $b_1 \neq b_2$ :  $(a, b_1) : \text{Attack}$  and  $(a, b_2) : \text{Attack}$ .*

**Proposition 2.** *The inverse of the Attack relation is an attack-like relation (*inv* *Attack*  $\sqsubseteq$  *Attack*).*

**Proposition 3.** *An argument  $a$  agrees to itself  $(a, a) : \text{Agree}$ , given by the the reflexivity property of the Agree relation.*

*Example 6.* Consider  $FRA = \langle \prec \{a : A, b : B, c : C\}, \{(a, b) : \text{Agree} \geq 0.9, (b, c) : \text{Agree} \geq 0.8\}, \{A, B, C\}, \mathcal{R}^k \cup \{\text{Agree} \sqsubseteq \bar{S}, \text{tra}(\text{Agree})\}$ . *Agree* being a transitive relation, the argument  $a$  also agrees to  $c$  with a degree of  $\alpha \geq \sup_{b \in \text{Arg}} \{\otimes_S((a, b) : \text{Agree}, (b, c) : \text{Agree})\}$ , which gives  $\max(0.9 + 0.8 - 1, 0) = 0.7$  under Lukasiewicz semantics and  $\min(0.9, 0.8) = 0.8$  in Gödel interpretation.

**Proposition 4.** *The relation  $S$  is the complement of the relation  $R$  if (*inv*  $R$   $S$ ) and  $(x, y) : R \bowtie \alpha \rightarrow (y, x) : S \bowtie^- (1 - \alpha)$ . Here, *Agree* and *Disagree* are complement relations,  $(a, b) : \text{Agree} \geq \alpha$  implies  $(b, a) : \text{Disagree} \leq 1 - \alpha$ . Informally, if  $a$  and  $b$  agree each other with at least  $\alpha$ , the disagreement degree between them should be less than  $1 - \alpha$ .*

**Definition 6.** (*Argumentation Chain*) *An argument  $b$  is supported by the argument  $a$  if there is a finite path  $p = (a, x_1) : R_1, (x_1, x_2) : R_2, \dots, (x_{n-1}, b) :$*

$R_n, \forall R_i, R_i \sqsubseteq \text{Support}$ . An argument  $b$  is attacked by the argument  $a$  if their is a finite path  $p = (a, x_1) : R_1, (x_1, x_2) : R_2, \dots, (x_{n-1}, b) : R_n, \forall R_i, R_i \sqsubseteq \text{Support} \sqcup \text{Attack}$ , and the number of attack relations  $|\mathcal{R}^A|$  is odd.

**Proposition 5.** (Indirect Support) *By composing an even number of attack relations, one gets an indirect support relation. Formally,  $R_1 \sqsubseteq \text{Attack}, R_1 \sqsubseteq \text{Attack}$ , implies  $R_1 \circ^t R_2 \sqsubseteq \text{Support}$ . The norm used to compute the strength of the attack is  $\circ^t_{\bar{A}} = \otimes^2_{\bar{S}}$ , where the power 2 models the fact that an indirect attack should be smaller than a direct one.*

*Example 7.* Let  $FRA = \langle \prec \{a : A, b : B, c : C\}, \{(a, b) : \text{Undercut} \geq 0.9 (b, c) : \text{Attack} \geq 0.7\} \succ, \{A, B, C\}, \mathcal{R}^k \cup \{\text{Disagree}, \text{Undercut}, \text{Disagree} \sqsubseteq \text{Attack}, \text{Undercut} \sqsubseteq \text{Attack}\}$ . By applying complex role inclusion we obtain  $\text{Attack} \circ^t \text{Disagree} \sqsubseteq \text{Support}$ . In other words, if the argument  $a$  attacks  $b$  and  $b$  disagrees with  $c$  we say that there is a support-like relation between  $a$  and  $c$ :  $(a, c) : R > 0, R \in \mathcal{R}^S$ . The degree of support is given under the Lukasiewicz semantics as

$$\text{Undercut} \otimes_L^2 \text{Disagree} = (\sup_{b:B} \{\max(0, 0.9 + 0.7 - 1)\})^2 = (0.6)^2 = 0.36$$

and under Gödel semantics:

$$\text{Undercut} \otimes_G^2 \text{Disagree} = (\sup_{b:B} \{\min(0.9 + 0.7)\})^2 = 0.49$$

### 3.2 Aggregation of Arguments

Several issues are raised by merging description logic and fuzzy argumentation: What happens when there is more than one attack-like relation between two concepts? What happens when one argument belongs with different membership functions to several concepts, which are linked by different attack-like relations with the opposite argument? What happens when two independent arguments attack the same argument? Should one take into consideration the strongest argument, or both of them may contribute to the degree of truth of that concept? Given an argumentation system, a semantic attaches a status to an argument. Different semantics may lead to different outputs [13].

One advantage of fuzzy logic is that it provides technical instrumentation (Lukasiewicz semantics, Gödel semantics) to handle all the above cases in an argumentative debate. The interpretation of Gödel operators maps the *weakest link principle* [16] in argumentation, which states that an argument supported by a conjunction of antecedents  $\alpha$  and  $\beta$ , is as good as the weakest premise  $\otimes_G = \min(\alpha, \beta)$ . The reason behind this principle is the fact that the opponent of the argument will attack the weakest premise in order to defeat the entire argumentation chain. When two reasons supporting the same consequent are available, having the strength  $\alpha$  and  $\beta$ , the strongest justification is chosen to be conveyed in a debate, which can be modeled by the Gödel union operator  $\oplus_G \max\{\alpha, \beta\}$ .



The Lukasiewicz semantics fits better to the concept of *accrual of arguments*. In some cases, *independent* reasons supporting the same consequent provide stronger arguments in favor of that conclusion. Under the Lukasiewicz logic, the strength of the premises  $(\alpha, \beta)$  contributes to the confidence of the conclusion, given by  $\oplus_L = \min\{\alpha + \beta, 0\}$ . For instance, the testimony of two witnesses is required in judicial cases. Similarly, several reasons against a statement act as a form of collaborative defeat [16]. One issue related to applying Lukasiewicz operators to argumentation regards the difficulty to identify independent reasons. Thus, an argument presented in different forms contributes with all its avatars to the alteration of the current degree of truth.

Thus, the description logic provides the technical instrumentation needed to identify independent justifications, whilst the Lukasiewicz semantics offers a formula to compute the accrual of arguments. The accrual of dependent arguments is not necessarily useless. Changing the perspective, this case can be valuable in persuasion dialogs, where an agent, by repeatedly posting the same argument in different forms, will end in convincing his partner to accept it.

### 3.3 Resolution Schemes

The key limitation of conventional systems is that, even if they guarantee to compute a solution for consistent sets, admissible or preferred extensions, it is possible that the only answer to be the empty set.

**Definition 7.** *The preference relation  $Pref \subseteq Arg \times Arg$  is a fuzzy role having the following properties: (ref  $P$ ), (tran  $P$ ), and (antysim  $P$ ). The Rbox  $\mathcal{R}$  extended with preferences is given by  $\mathcal{R}^P = \mathcal{R} \cup \{Pref\}$ . An argument  $a$  is preferred to  $b$  ( $a \succ b$ ) based on the preference role  $Pref$  with a degree  $\alpha$  if  $(a, b) : Pref \bowtie \alpha$ .*

*Example 8.* Consider the task to classify a compound according to potential toxicity. In the guidelines of U.S. Environmental Protection Agency for the assessment the health impacts of potential carcinogens, an argument for carcinogenicity that is based on human epidemiological evidence is considered to outweigh arguments against carcinogenicity that are based only on animal studies. The corresponding FRA (figure 2) augmented with preferences is

$$\begin{aligned} FRA^P = \langle & \prec (h : HumanStudy, a : A, b : B, c : Carcinogenicity, \\ & d : AnimalStudy, (a, h) : BasedOn, (a, c) : For, (b, d) : BasedOn, \\ & (b, c) : Against, (a, b) : Outweigh \succ, \\ & \{HumanStudy, AnimalStudy, Carcinogenicity, A, B\}, \mathcal{R}^k \cup \\ & \{For \sqsubseteq Support, Against \sqsubseteq Attack, (inv BasedOn \sqsubseteq Support)\} \\ & \cup \{Outweigh \sqsubseteq Pref\} \end{aligned}$$

Observe that if  $a$  is based on  $h$  then there exists a support-like relation from  $h$  towards  $a$ . Formally  $(inv BasedOn \sqsubseteq Support)$ .

In the above example, the preference assertion *Outweigh* between the arguments  $a$  and  $b$  was explicitly given, stating that  $a$  clearly outweighs  $b$  ( $\alpha = 1$ ). Note that any preference role can be a fuzzy one. When this explicitness does not exist, FDL offers the possibility to infer preference relations among arguments based on various *conflict resolution strategies*, like the following ones.

- *Fuzzy membership value (M)*. The status of an argument is assessed by comparing the membership degrees of the arguments to their concepts. Prior information is usually provided by an expert or knowledge engineer.
- *Specificity (S)*. This heuristic can be applied both on concepts, in which case the most specific argument dominates, and roles, where the most specific relation in the hierarchy prevails (figure 3).
- *Value based argumentation (V)*. The argument which promotes the highest value according to some strict partial ordering on values will defeat its counter-argument. In a FRA, arguments can promote (or demote) values to a given degree, so that if the arguments  $a$  and  $b$  promote the same value  $v$ , we consider that  $a$  successfully attacks  $b$  if it promotes  $v$  to a greater degree than  $b$ . In the current framework, values can be used from an ontology of values, providing a reasoning mechanism over values.

*Example 9.* (Specificity on concepts) Consider the FRA base

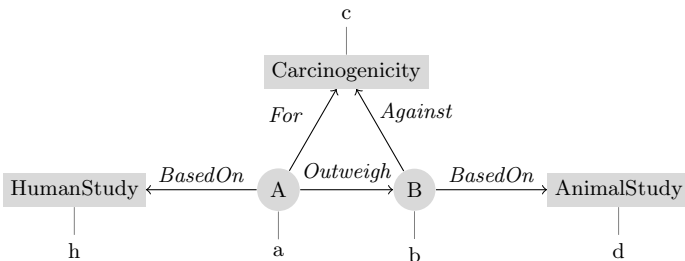
$$FRA^{\mathcal{P}} = \langle \prec (a : A, c : C, d : D, A \sqsubseteq B) \succ, \{A, B, C, D\}, \mathcal{R}^k\{(C, B) : Attack, (D, A) : Support\} \rangle$$

So,  $a$  as an element of  $A \sqsubseteq B$  is supported by  $d$ , and attacked by  $c$  (figure 3a). In this case the specificity principles says that the support relation will prevail.

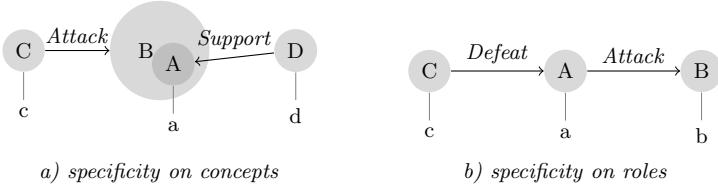
*Example 10.* Now consider the case in which  $(a, b) : Defeat, (b, c) : Attack$  and  $Defeat \sqsubseteq Attack$  (figure 3b). Based on the specificity heuristic on roles, the *Defeat* relation is stronger (more specific) than *Attack*. Consequently, the only admissible set is  $\{a, c\}$ .

The specificity preference is also illustrated by the following dialog in figure 4. Here, the agent  $B$  accepts the argument *Food*, defined as  $Food \sqsubseteq \exists canEat$  and supported by the argument *hungry*. Observe that the support is stronger from the agent  $A$ , given by the fuzzy modifier *very*, and not so convincing as denoted by the modifier *little*. However, the more specific argument *pizza*, which is a kind of food ( $Pizza \sqsubseteq JunkFood \sqsubseteq Food$ ) is rejected by the agent  $B$ . Similarly, the argument *fish* conflicts with the argument *expensive*.

In many real life discussions, people have agreements at certain level of generality, while they manifest divergent opinions starting with a given level of



**Fig. 2.** Explicit preference among arguments



**Fig. 3.** Conflict Resolution Strategies

A: *I am very hungry. Let's go eat something.*  
 B: *I am a little hungry too. I agree.*  
 A: *I don't have too much time. Let's have a pizza.*  
 B: *It's not healthy. I prefer something else. What about fish and wine.*  
 A: *It's too expensive.*

**Fig. 4.** Illustrating the specificity principle

specificity. Description logic is particularly useful to define the edge between agreement and disagreement. In this particular case, the agreed concept would be  $NutritionalFood \sqcap \forall hasPrice. \neg Expensive$ . The instance that belongs to this concept with the highest degree will best satisfy the both agents. If such an assertion does not exist, the agreement is reached by a preference relation over the common constraints: *healthy, not expensive, quick*.

Note that preferences are fuzzy relations, meaning that linguistic scale can be defined on them. Preferences like: *just equal, weakly more important, fairly strongly preferred, absolutely outweighed* are accepted in a  $FRA^P$ .

**Proposition 6.** *By composing two preference relations we get a preference relation. In order not to breach the transitivity property, the composition function that we use is  $\circ^t_{Pref} = \otimes_S^{1/2}$ .*

*Example 11.* Let  $(a, b) : Outweigh \geq 0.9$  and  $(b, c) : Pref \geq 0.7$ . The preference degree between  $a$  and  $c$  is given by  $Outweigh(a, b) \circ^t Pref(b, c) = \min(0.9, 0.7)^{1/2} = 0.83$ .

### 3.4 Semantic Inconsistency

An important aspect is that inconsistency is naturally accommodated in fuzzy logic: the intersection between the fuzzy concept  $A$  and its negation is not 0 ( $A \sqcap \neg A \neq 0$ ). Similarly, the disjoint property of an attack  $A_1 \sqsubseteq Attack$  and support  $S_1 \sqsubseteq Support$  relation  $\otimes_S((a, b) : A_1 \geq \alpha, (a, b) : S_1 \geq \beta) = 0$ , under the Lukasiewicz interpretation leads to an inconsistent argumentation base if  $\alpha + \beta > 1$ .

Consider the concept *InternallyConsistentArguments* (ICA) defined in FDL as:  $ICA \equiv Argument \sqcap \neg \exists Attack. ICA$ . Based on the  $(\forall R.C)^I = (\neg \exists R. \neg C)^I$ , which holds under the Lukasiewicz semantics [20], follows that:  $ICA \equiv Argument \sqcap \forall Attack. \neg ICA$ . The semantics of  $\forall R.C$  being  $(\forall R.C)^I = \inf_{y \in \Delta^I} = R^I(x, y) \rightarrow$

$C^I(y)$  implies in FRA that if  $(x, y) : Attack \bowtie \alpha \rightarrow_L y : \neg ICA \bowtie \beta$  The implication holds if  $1 - \alpha + \beta \geq 1$  (recall table [III](#)), or  $\alpha \leq \beta$ .

In order to keep the argumentation base semantically consistent the following constraints exist, where  $\gamma = 1 - \beta$ , represents the degree of  $y$  to the  $ICA$  concept:

- $(x, y) : Attack \leq \alpha \Rightarrow \alpha \leq \beta \Leftrightarrow \gamma < 1 - \alpha$ : If the attack relation between  $x$  and  $y$  is maximum  $\alpha$ , the knowledge base remains consistent as long as  $y$  belongs to the concept  $ICA$  no more than  $1 - \alpha$ .
- $(x, y) : Attack \geq \alpha \Rightarrow \beta = 1 \Leftrightarrow \gamma = 0$ : If the attack relation between  $x$  and  $y$  is at least  $\alpha$ , the knowledge base is guaranteed to remain consistent if  $y$  does not belong to  $ICA$  at all.
- $y : \neg ICA \leq \beta \Rightarrow \alpha = 0$ : If  $y$  belongs to the concept  $\neg ICA$  with maximum  $\beta$ , it means that it should belong to the opposite concept  $IAC$  at least  $1 - \beta$ . Consequently, no attack relation should exist between  $x$  and  $y$ .
- $y : \neg ICA \geq \beta \Rightarrow \alpha \leq \beta \Leftrightarrow \gamma \leq 1 - \alpha$ .

The notion of indirect support in combination with the disjoint property of the attack and support relation may help to signal semantic inconsistencies in an argument bases. If  $A$  attacks  $B$  attacks  $C$  and  $A$  attacks  $C$ , then  $A$  indirectly both supports  $C$  and attacks  $C$ .

*Example 12.* Consider the situation in which  $(A, B) : attack_{0.6}$ ,  $(B, C) : attack_{0.9}$ , and also  $(A, C) : attack_{0.7}$ . Under the Lukasiewicz semantics, the indirect support from  $A$  to  $C$  equals  $\max(0.6 + 0.9 - 1, 0)^2 = 0.25$ . The disjoint property of attack and support holds because  $\max(0.25 + 0.7 - 1, 0) = 0$ . If the strength of the attack from  $A$  to  $C$  increases to  $0.7$  the disjoint property is violated. In this case, the framework signals to the human agent that the argument base is semantically inconsistent. In other words, the alert means that on the these particular argumentation chains the strengths of the attacks or supports might be incorrect stated and the initial facts should be reconsidered.

There is no need to explicitly define simple negation on roles, as it exists in FDL systems by mean of assertions that use the inequalities,  $\leq$  and  $<$  [\[19\]](#). For instance, the assertion  $x$  does not attack  $y$  can be defined as  $(x, y) : Attack \leq 0$

## 4 Argumentative Agents in FRA

**Definition 8.** A preference scheme specifies the order in which the conflict resolution strategies are applied. An example of preference scheme is MSV (fuzzy membership value, specificity, value-based).

Preference schemes come from the cognitive system of the agents. A  $FRA^P$  still allows a lack of complete transitive preference.

**Definition 9.** An agent is a tuple  $Ag = [PreferenceScheme, (\oplus, \otimes, \ominus, \rightarrow), (\circ_{Pref}, \circ_{\bar{A}}, \circ_S)]$ , where  $\oplus, \otimes, \ominus, \rightarrow$  represent the union, intersection, negation and implication operators, and  $\circ_{Pref}, \circ_{\bar{A}}, \circ_S$  the norm used for composing preferences, attacks, and support relations.  $KB$  encapsulates the private domain knowledge of the agent.

We assume that agents acting within the same  $FRA^P$  share a common vocabulary of attacking, supporting, and outranking fuzzy relations and common understanding of their formal properties. However, they have their own order of preferences and own functions for aggregating arguments. The inconsistency budget of each agent emerges from the combination of these functions. The personality of the agents can be encapsulated also by the above combination

*Example 13.* An agent *Judge* =  $[MSV, (\oplus_L, \otimes_L, \ominus_L, \rightarrow_L), (\otimes_L^{1/2}, \otimes_L^2, \otimes_L^2)]$ , by aggregating arguments under the Lukasiewicz semantics takes into consideration all the existing facts. It acts based on a hierarchy of values derived from the hierarchy of laws. In case of conflict, the most specific norm, which in general refers to exceptions, will be applied, then the argument from the most recent case (in case based law) or most recent norm (in norm based law). Afterward, the fuzziness of some linguistic terms from the law, will be considered in the decision.

For modeling practical scenarios we follow the steps:

1. **Identify** the relevant concepts and their attacking and supporting relations.
2. **Define** the membership functions.
3. **Formalise** the FRA:
  - (a) define the class hierarchy;
  - (b) define the role hierarchy;
  - (c) define the role properties;
  - (d) define the membership of arguments to their concepts;
  - (e) define the known strengths of the attacking and supporting relations.
4. **Build** the argumentation network.
5. **Reason** on the knowledge within the FRA based on their own preference schemes and aggregating functions:
  - (a) identify semantic inconsistencies;
  - (b) identify indirect attacks and supports;
  - (c) reduce the argumentation network by composing the fuzzy relations;
  - (d) aggregate arguments;
  - (e) compute the defeat status based on the active preference scheme;

## 5 A Case for Food Supply Chains

ISO 22000 is a recent standard designed to ensure safe food supply chains worldwide. Its main component is the HACCP (Hazard Analysis at Critical Control Points) system, which is a preventive method of securing food safety. Its objectives are the identification of consumer safety hazards that can occur in the supply chain and the establishment of a control process to guarantee a safer product.

## 5.1 Technology Drivers for HACCP

HACCP is based on the following steps and principles [9]. In the first step, the business entities within the supply chain determine the *food safety hazards* and identify the preventive measures for controlling them. Then, the *critical control points* (CCP) are identified. They represent point steps in a food process at which a specific action can be applied in order to prevent a safety hazard, or to reduce it to an acceptable level. Afterward, *critical limits* are established, representing criteria which separate acceptability from unacceptability. Criteria often used include measurements of time, temperature, moisture level, pH, Aw, available chlorine, and sensory parameters such as visual appearance and texture. Critical limits must be specified for each CCP and they should be justified. A *monitor process* is followed by the *establishment of critical actions* in order to deal with deviations when they occur. Then procedures for verifications are needed to confirm that the HACCP system is working effectively. Finally, the documentation is needed to encapsulate justifications for all the decisions which have been taken. The main goal of the standard is to build confidence between suppliers and customers. It demands that business entities follow specific well-documented procedures, in which the quality of the items should be demonstrated by different types of justifications, and not only by attaching a quality label to the product. The technical requirements for building an HACCP system lay around the need to integrate support for argumentative debates.

**Structured Argumentation.** The technical support for argumentation is needed during the HACCP development for various tasks.

*Justifying hazards.* The HACCP standard explicitly requests that arguments pro and against should be provided in order to justify all decisions to classify hazards as critical or not critical, formalised as  $Hazard = \exists hasJustification.Argument$ . Based on the above definition in DL, the reasoner can check that a justification is attached to both significant or not significant hazard.

*Justifying control options.* For each hazard which is considered significant, a control measure should be defined ( $SignificantHazard = Hazard \sqcap hasControl Measure.\top$ ). The absence of the control measure is signaled as an inconsistency by the reasoner. The advantages and disadvantages of each available option should be backed by supporting arguments, respectively counter-arguments.

*Justifying associated critical limits.* The recommended sources of information for justifying the chosen critical limits are: norms, experts, scientific publications, or experimental studies. The rationale and the reference material should become part of the HACCP plan [9].

**Domain knowledge.** When implementing the HACCP standard, the human experts need ontological knowledge in the following activities.

*Hazard identification.* The user can query hazard ontologies and their possible connections with ingredients and processing steps. Also, food and pathogen ontologies may be used to compare different risks which may stem from the production system.

*Automatic verification of the safety conditions.* Having formal descriptions about what a safety device, process, or service represent (encapsulated as TBoxes) and by having the current situation (encapsulated as ABoxes) the system can automatically point out possible contradictions with the norms in use.

**Fuzzy Reasoning.** It is used as a tool for qualitatively assessing during the following activities:

*Assessment of critical control points.* For each step of the production process, one should decide whether that stage will be a CCP. The decision depends on the hazard possibility of occurrence (terms such as rarely, often, sometimes, always are used in practice by the experts) and on its severity (usually assessed as low, medium, high).

*Supply chain integration.* An important source of hazards appears when receiving the input items. Depending on the potential risk, the company should decide to rely on the information from the product label or to conduct its own measurements of the product characteristics. This qualitative decision is based on fuzzy assessments. Also, the feedback received from the buyers, representing their preferences and perceptions is fuzzy. The costumers subjective evaluation can refer to attributes such as: color, smell, taste. Moreover, the company decides if it is able to deal with all the identified hazards or to outsource this task. For instance, the presence of rodents, insects, birds or other pests is unacceptable. The hazards are related both to the direct effects of these pests, and to the risks coming from the substances used to eliminate them. A good option is to contract a specialized company to handle these hazards.

*Process adjustment.* Actions need to be taken to bring the CCP under control before the critical limit is exceeded. The point where the operators take such action is called the operating limit. The process should be adjusted when the operating limit is reached to avoid violating the critical limit.

*Modeling fuzzy critical limits.* Consider some microbiological data. One rule can say: "The product is safe if it is kept no longer than 48 hours at a temperature below  $10^0 C$ ". What happens if the product is kept 47 hours at the temperature of  $9^0 C$ ? Is it safer comparing with the situation in which it is kept for 1 hour at a temperature of  $12^0 C$ ? According to the above rule, the second item is not considered safe. As the alteration of product features is gradual, fuzzy membership functions being able to model these cases.

## 5.2 An HACCP Scenario

The framework is exemplified by for a cooked shrimp company. The control measure has emerged after an argumentation process.

Agent HACCP Plan: Justifying control measures	
$E_1$	The first option for the control measure is to set a microbiological limit, under which the product is considered safe. This direct method minimizes error measurements, but I admit the monitoring process is expensive.
$E_2$	Several tests are necessary to determine critical limits derivations and samples may need to be large to be meaningful.
$E_3$	Moreover, the results are obtained in several days.
$E_2$	The second option is to set a minimum internal temperature at which the pathogens are destroyed. The method is practical and more sensitive.
$E_1$	But justification is needed to validate the chosen temperature value.
$E_3$	The third option is to control the factors that affect the internal temperature of the product (oil, thickness of the pane, or cooking time).
$E_1$	The method requires justifications between these limits and the internal temperature of food.
$E_3$	I agree. Nevertheless, it is very practical and it increases confidence in the measurements.
$B_1$	The business policy encourages practical and reliable solutions.

**Fig. 5.** Arguing for the adequate control measures in an HACCP plan

The first step is two identify the main concepts and the relations among them. In this case, three options exist, each supported by its advantages and attacked by the disadvantages that it brings (figure 6, step 1). In the second step the agents should agree on the fuzzy membership functions for each concept in the domain. In figure 6 three such membership functions are shown. The *ExpensiveToMonitor* concept is defined as a trapezoidal number  $ExpensiveToMonitor \equiv \exists hasCost.trapezoidal(10, 20, 30, 30)$ . Consider a particular pathogens limit  $d_1$ , which has the estimated cost at least 18 in order to be validated. The degree of membership to the concept *ExpensiveToMonitor* will be 0.8 ( $d_1 : ExpensiveToMonitor \geq 0.8$ ), reflected in the strength of the attack between  $d_1$  and the first option  $o_1$ .

We consider the agents  $E_1, E_2$  (decision makers) needing to be involved with other agents  $E_3, B_1$  (consultants) in the process.

$$E_1 = [M, (\oplus_G, \otimes_G, \ominus_L, \rightarrow_G), (\otimes_G^{1/2}, \otimes_G^2, \otimes_G^2)]$$

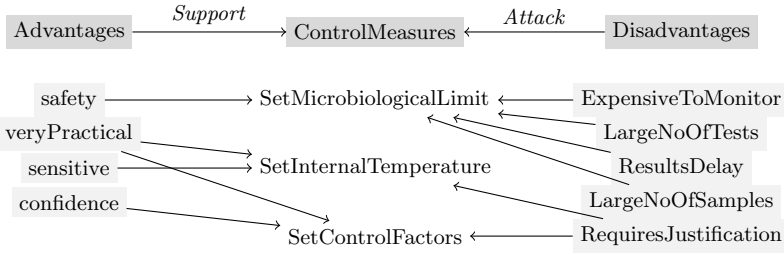
$$E_2 = [MS, (\oplus_L, \otimes_L, \ominus_L, \rightarrow_L), (\otimes_L^{1/2}, \otimes_L^2, \otimes_L^2)]$$

They start by aggregating the direct attack and support roles. The aggregation of  $d_1, d_2, d_3, d_4$  gives an attack strength of  $max(0.8, 0.5, 0.4, 0.7) = 0.8$  for the first option  $o_1$ , under the Gödel semantics (the agent  $E_1$ ) and  $min(min(min(0.8 + 0.5), 1) + 0.4, 1) + 0.7, 1) = 1$  (figure 7). Because both arguments  $a_2$  and  $a_3$  support the second option  $o_2$  there strengths are aggregated, giving  $max(0.55, 0.49) = 0.55$  for the agent  $E_1$ , respectively  $min(0.55 + 0.49, 1) = 1$  for the agent  $E_2$ .

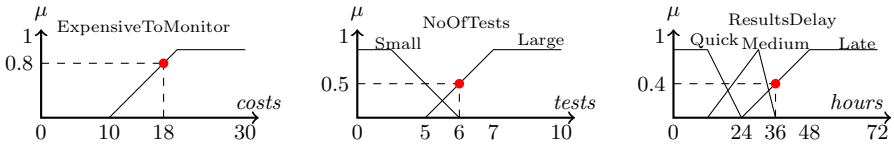
Next, the indirect relations are taken into consideration. By composing the supporting relation *Enc* with *Support* we get an indirect attack between  $b_1$  and



1. Identifying concepts and attacking and supporting relations.



2. Define the membership functions



4. Formalizing the FRA

(A=Advantages, D=Disadvantages, O=ControlMeasures, B=BusinessPolicy, Enc=Encourage)  
 $\mathcal{A} = \langle \{a_1 : \text{Safety} \geq 0.6, a_2 : \text{Practical} \geq 0.7, a_2 : \text{very}(\text{Practical}) \geq 0.49,$   
 $a_3 : \text{Sensitive} \geq 0.3, a_4 : \text{Confidence} \geq 0.5, d_1 : \text{ExpensiveToMonitor} \geq 0.8,$   
 $d_2 : \text{LargeNoOfTests} \geq 0.5, d_3 : \text{ResultsDelay} \geq 0.4, d_4 : \text{LargeNoOfSamples} \leq 0.8,$   
 $d_5 : \text{RequiresJustification}, b_1 : B, o_1 : O, o_2 : O, o_3 : O \rangle$   
 $\mathcal{T} = \{A, D, O, B, \text{Safeness} \sqsubseteq A, \text{Practical} \sqsubseteq A, \text{Sensitive} \sqsubseteq A, \text{Confidence} \sqsubseteq A,$   
 $\text{ExpensiveToMonitor} \sqsubseteq D, \text{LargeNoOfTests} \sqsubseteq D, \text{ResultsDelay} \sqsubseteq D,$   
 $\text{LargeNoOfSamples} \sqsubseteq D, \text{RequiresJustification} \sqsubseteq D, (A, O) : A, (A, O) : \bar{S}\}$   
 $\mathcal{R} = \mathcal{R}^k \cup \{\text{Enc} \sqsubseteq \bar{S}\}$

3. Building the Argumentation Network

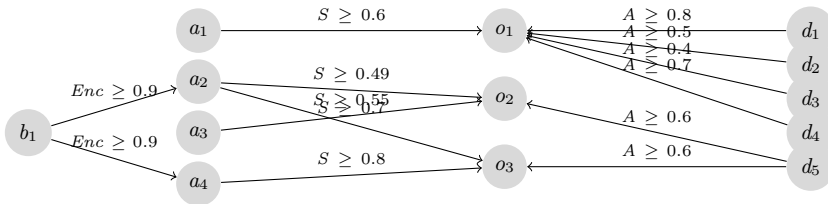
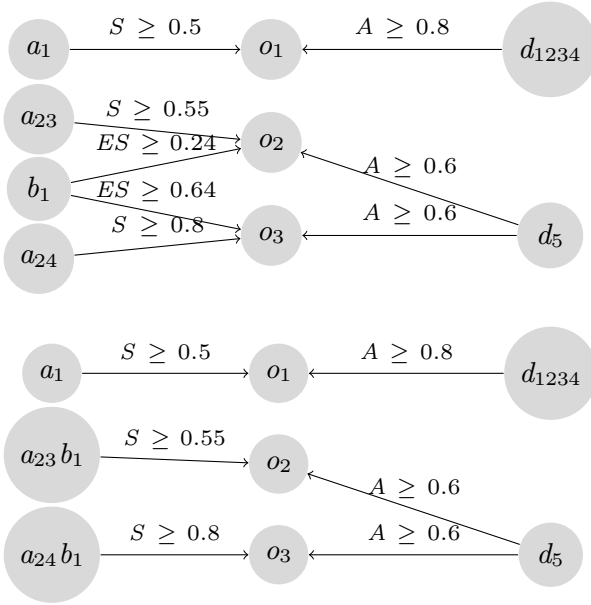


Fig. 6. A FRA for Justifying Control Measures

$o_2$  of strength  $\min(0.9, 0.49)^2 = 0.24$  under Logic semantics and  $\max(0.9 + 0.49 - 1, 0)^2 = 0.15$ . Hence,  $(b_1, o_1) : \text{Enc}(\otimes_G)^2 \geq 0.24$  for  $E_1$  and  $(b_1, o_1) : \text{Enc}(\otimes_L)^2 \geq 0.15$ . Note that  $b_1$  indirectly supports the option  $o_3$  through two intermediary nodes  $a_2$  and  $a_4$ . The amount of indirect support provided by  $b_1$  for  $o_3$  is  $\max(\min(0.9, 0.55)^2, \min(0.9, 0.8)^2) = 0.64$  from the  $E_1$ 's perspective. The agent  $E_2$  computes the degree of support from  $b_1$  towards  $o_3$  as  $\min(\max(0.9 + 0.55 - 1, 0)^2 + \max(0.9 + 0.7 - 1, 0)^2, 1) = 0.56$ .

$$E_1 = [M, (\oplus_G, \otimes_G, \ominus_L, \rightarrow_G), (\otimes_G^{1/2}, \otimes_G^2, \otimes_G^2)].$$



$$E_2 = [MS, (\oplus_L, \otimes_L, \ominus_L, \rightarrow_L), (\otimes_L^{1/2}, \otimes_L^2, \otimes_L^2)].$$

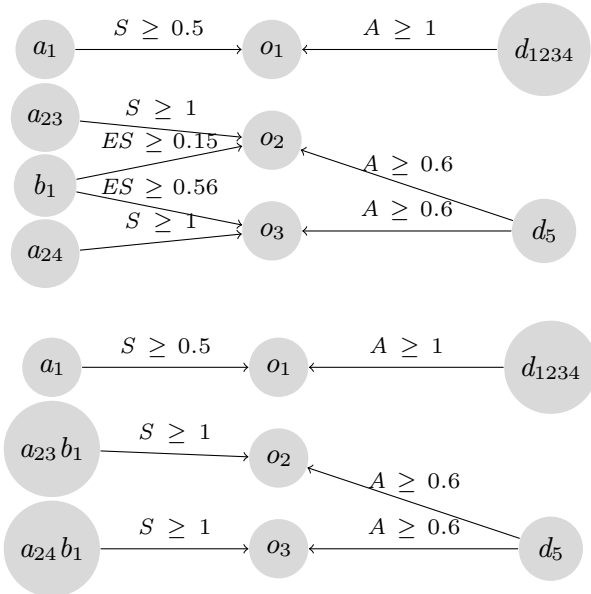


Fig. 7. Argumentative agents

The bottom part of figure 7 depicts the aggregation of direct and indirect relations. The compound concept  $a_{23}b_1$  supports the option  $o_2$  with at least  $\max(0.55, 0.24) = 0.55$  from the  $E_1$  viewpoint and with  $\min(1 + 0.15, 0) = 1$  from that of the expert  $E_2$ . The support given by  $a_{24}b_1$  to the third option equals  $\max(0.64, 0.8) = 0.8$ , respectively  $\min(1 + 0.49, 1) = 1$ .

The attack on the first option is stronger than its support, from both perspectives  $E_1$  and  $E_2$ . Hence, an agreement between the agent  $E_1$  and  $E_2$  exists to eliminate this option. Given the above information, the agent  $E_2$  equally accepts the options  $o_3$  and  $o_2$ . The degree of support is 1 and the attack 0.6 in both cases. On the other hand, the agent  $E_1$  rejects the option  $o_2$  but accepts option  $o_3$ . Consequently, the last control measure gets the support from both parties.

## 6 Discussion and Related Work

Arguments supporting both a consequent and its negation co-exist in the knowledge base. To overcome this drawback, *weighted argument systems* (WAS) have been introduced, with the notion of *inconsistency budget* [8] used to characterize how much inconsistency one is prepared to tolerate in an argumentation base. A FRA framework is a particular instance of a WAS, where the degree of inconsistency is accommodated by the semantics of fuzzy reasoning.

Other approaches have investigated imprecise argumentation [12,1,10]. The fuzzy argumentation framework [12] is an extension of the classical Dung model, while in our approach, the fuzzy component is meant to help software agents to exploit the real arguments conveyed by humans. Compared to the defeasible logic approach [1,10], where the ontological knowledge is embedded in the program, we have been interested in having a clear, separate representation of the ontology, allowing for transparency.

Rahwan and Banihashemi [17] address the idea of modeling argument structures in OWL, where arguments are represented in the Argument Interchange Format ontology (AIF), a current proposal for a standard to represent arguments. A mapping between the top level concepts of the AIF ontology and our research can be done as follows: support roles over the rule application nodes, attack roles over the conflict application nodes, conflict resolution strategies over preference nodes. Meta-argumentation [6,15] is supported by the AIF approach: one can apply a preference on preference, attack a support, or support a preference. Meta-argumentation can be handled in FRA indirectly by defining new concepts and applying roles on them. For instance, the preference relation  $P$  applied to the argument  $a$  over  $b$  can be encapsulated as the concept *a preferred to b*, which makes possible to apply further attack or support roles on it. If one wants to challenge the degree  $\alpha$  of membership of an element  $a$  to the concept  $A$ , the new concept  $A_{\geq\alpha}$  can be defined and the attack should be applied on it. A hierarchy of Dung frameworks is proposed in [14], in which level  $n$  arguments refer to level  $n-1$  arguments, and conflict based relations and preferences between level  $n-1$  arguments. Arguing hierarchically is handled by navigating through the concepts which are subject of dispute and which can be organized hierarchically based on description logic.

The role of ontologies to resolve conflicts among arguments based on the specificity principle appears also in [3]. The existing work has focused on concept properties only, and do not exploit the formal properties of the attack relations. Our formalism based on FDL contributes to the current vision of developing the infrastructure for the World Wide Argument Web.

## 7 Conclusions

The fuzzy-based approach presented in this paper makes a step towards practical applications, a fuzzy-based argumentative system being cognitively less demanding for the decision makers. Real argumentative debates implies other relations and concepts, not only attacking and support roles or arguments. This additional domain knowledge can be easily integrated into a FRA framework. The main contribution comes from the introduction of different types of attack and support roles with a specific semantics given by their formal properties, with no need to invent a new mechanism to compute the strength of the attack. We have just used the technical instrumentation provided by fuzzy logic for computing the status of argument.

We advocate the merging of argumentation theory with semantic technologies, which leads to the possibility to reuse the argumentation bases among multi-agent systems. The preference schemes have already proved their success as conflict resolution strategies in expert systems.

One drawback is that we assume a common ontology of attacking and supporting roles. In the presented scenario, this limitation is partial overcome by the fact that HACCP represents a standard, which implies a common description of concepts and procedures. Also, the implication of fuzzy composition to indirect attacks needs a deeper investigation to validate the proposed semantics. Third, a methodology for modeling practical applications with argumentation theory is needed. Reasoning with the finitely many-valued Lukasiewicz Fuzzy Description Logic SROIQ [5] should prove even more appropriate to deal with imprecise and vague knowledge, inherent to several real world domains.

## Acknowledgments

We are grateful to the anonymous reviewers for the very useful comments. Part of this work was supported by the grant ID\_170/672 from the National Research Council of the Romanian Ministry for Education and Research.

## References

1. Alsinet, T., Chesñevar, C.I., Godo, L., Simari, G.R.: A logic programming framework for possibilistic argumentation: Formalization and logical properties. *Fuzzy Sets and Systems* 159(10), 1208–1228 (2008)
2. Bench-Capon, T.J.M., Dunne, P.E.: Argumentation in artificial intelligence. *Artificial Intelligence* 171(10-15), 619–641 (2007)

3. Bench-Capon, T.J.M., Gordon, T.F.: Isomorphism and argumentation. In: Twelfth International Conference on Artificial Intelligence and Law, pp. 11–20 (2009)
4. Bobillo, F., Straccia, U.: fuzzyDL: An expressive fuzzy description logic reasoner. In: International Conference on Fuzzy Systems, Hong Kong, pp. 923–930. IEEE Computer Society Press, Los Alamitos (2008)
5. Bobillo, F., Straccia, U.: Reasoning with the finitely many-valued Lukasiewicz fuzzy description logic SROIQ. *Information Sciences* (2011) (to appear)
6. Boella, G., Gabbay, D.M., van der Torre, L., Villata, S.: Meta-argumentation modelling i: Methodology and techniques. *Studia Logica* 93(2-3), 297–355 (2009)
7. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artif. Intell.* 77, 321–357 (1995)
8. Dunne, P., Hunter, A., McBurney, P., Parsons, S., Wooldridge, M.: Inconsistency tolerance in weighted argument systems. In: The Ninth International Conference on Autonomous Agents and Multiagent Systems, pp. 851–858 (2009)
9. Food, of the United Nations World Health Organization, A.O.: *Codex Alimentarius* (1997)
10. Gomez, S.A., Chesñevar, C.I., Ricardo, G.: Reasoning with inconsistent ontologies through argumentation. *Applied Artificial Intelligence* 24, 102–148 (2010)
11. Hunter, A.: Real arguments are approximate arguments. In: Proceedings of the 22nd National Conference on Artificial Intelligence, pp. 66–71 (2007)
12. Janssen, J., De Cock, M., Vermier, D.: Fuzzy argumentation frameworks. In: 12th International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems, pp. 513–520 (2008)
13. Madakkatel, M.I., Rahwan, I., Bonnefon, J.F., Awan, R.N., Abdallah, S.: Formal argumentation and human reasoning: The case of reinstatement. In: Proceedings of the AAAI Fall Symposium on The Uses of Computational Argumentation, pp. 46–51. AAAI Press, Washington DC, USA (2009)
14. Modgil, S.: Hierarchical argumentation. In: 10th European Conference on Logics in Artificial Intelligence, pp. 319–332 (2006)
15. Modgil, S.: Reasoning about preferences in argumentation frameworks. *Artificial Intelligence* 173(9-10), 901–934 (2009)
16. Pollock, J.: Defeasible reasoning with variable degrees of justification. *Artificial Intelligence* 133, 233–282 (2002)
17. Rahwan, I., Banihashemi, B.: Arguments in OWL: A progress report. In: Second International Conference on Computational Models of Argument, pp. 297–310 (2008)
18. Rahwan, I., Zablith, F., Reed, C.: Laying the foundations for a World Wide Argument Web. *Artificial Intelligence* 171(10-15), 897–921 (2007)
19. Straccia, U.: Reasoning within fuzzy description logics. *Journal of Artificial Intelligence Research* 14, 137–166 (2001)
20. Straccia, U.: A fuzzy description logic for the semantic web. In: Sanchez, E. (ed.) *Capturing Intelligence: Fuzzy Logic and the Semantic Web*. Elsevier, Amsterdam (2006)

# Dynamic Argumentation in Abstract Dialogue Frameworks\*

M. Julieta Marcos, Marcelo A. Falappa, and Guillermo R. Simari

National Council of Scientific and Technical Research (CONICET)  
Artificial Intelligence Research & Development Laboratory (LIDIA)  
Universidad Nacional del Sur (UNS), Bahía Blanca, Argentina  
Avenida Alem 1253 (B8000BCP)  
Tel.: (0291) 459-5135; Fax: (0291) 459-5136  
{mjm, mfalappa, grs}@cs.uns.edu.ar

**Abstract.** In this work we present a formal model for collaborative argumentation based dialogues by combining an abstract dialogue framework with a formalism for dynamic argumentation. The proposed model allows any number of agents to interchange and jointly build arguments in order to decide the justification status of a given claim. The model is customizable in several aspects: the argument attack relation and acceptability semantics, the notion of relevance of contributions, and also the degree of collaboration are selectable. Important properties are ensured such as dialogue progress step by step, completeness of the sequence of steps, and termination. Under the higher degree of collaboration, the dialogue constitutes a sound and complete distributed argumentation process.

**ACM Categories and Subject Descriptors:** I.2.11 [Distributed Artificial Intelligence]: Coherence and coordination.

**General Terms:** Theory, Design.

**Keywords:** Collective intelligence, Dialogue, Argumentation.

## 1 Introduction and Motivation

Multi-agent systems (MAS) provide solutions to problems in terms of autonomous interactive components (agents). A *dialogue* is a kind of interaction in which a sequence of messages, over the same topic, is exchanged among a group of agents, with the purpose of jointly drawing some sort of conclusion. There is a subset of dialogues, which we call *collaborative*, in which the agents are willing to share any relevant knowledge to the topic at issue, having no other ambition than achieving the right conclusion on the basis of all the information they have.

*Argumentation-based dialogues* usually consist of interchanging arguments for and against certain claim. Mostly in the literature, these dialogues are held between two agents, one of them putting the arguments ‘for’ and the other putting the arguments

---

\* This research is partially supported by Sec. Gral. de Ciencia y Tecnología (Univ. Nac. del Sur), CONICET and Agencia Nac. de Prom. Científica y Técnica (ANPCyT).

‘against’. In order to achieve collaborative (in the sense described above) behavior, all the participants should contribute with both kinds of arguments, and also they should be able to jointly build new arguments. Even as part of non-collaborative dialogues (*e.g.* persuasion) it may be useful to build arguments in conjunction.

Classical *abstract argumentation* [3] assumes a static set of already built arguments, resulting insufficient for modeling collaborative dialogues. The set of arguments involved in a dialogue is, in contrast, dynamic: new arguments jointly constructed by the agents may arise, and also arguments may be invalidated (note this is not the same as defeated) at the light of new information. The argument construction step cannot be performed separately from the dialogue.

Recently, a *dynamic abstract argumentation framework (DAF)* has been proposed by Rotstein *et al.* [13], which extends the work done on acceptability of arguments, by taking into consideration their construction and their validity with respect to a varying set of evidence. This approach results, hence, very suitable for the modeling of collaborative dialogues. The main elements of the DAF are summarized in Sect. 2.

In [6] we have defined an *abstract dialogue framework ( $\mathcal{D}\mathfrak{F}$ )* together with a set of *collaborative semantics* which characterize different levels of collaboration in dialogues, in terms of a given reasoning model and a given notion for the relevance of contributions. Under certain natural conditions, the proposed semantics ensure important properties of collaborative dialogues, such as termination and outcome-determinism.

The aim of this work is to show how the abstract dialogue framework and semantics [6] can be applied to dynamic argumentation [13]. As will be seen, the agents will interchange both arguments (in Rotstein’s sense) and evidence, achieving the joint construction of arguments in the usual sense. A particular framework for argumentation-based dialogues will be obtained, which inherits the semantics and properties defined for the abstract framework. In sections 4 through 6 we will reintroduce the abstract concepts that constitute the  $\mathcal{D}\mathfrak{F}$  showing how they can be instantiated in terms of the DAF.

## 2 Background

Next we summarize an abstract argumentation framework capable of dealing with dynamics through the consideration of a varying set of evidence [13]. Depending on a particular situation (given by the content of the set of evidence), an instance of the framework will be determined, in which some arguments hold and others do not.

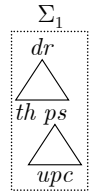
The formalization is coherent with classical abstractions [3], however arguments play a smaller role: they are aggregated in structures. These argumental structures can be thought as if they were arguments (in the usual sense), but they will not always guarantee their actual achievement of the claim.

A language  $\mathbf{L}$  will be assumed for the representation of evidence, premisses, and claims. An *argument*  $\mathcal{A}$  is a pair  $\langle \{s_1 \dots s_n\}, \delta \rangle$  consisting of a consistent set of *premisses*, noted  $\text{supp}(\mathcal{A})$ , and a *claim*, noted  $\text{cl}(\mathcal{A})$ . These basic premisses are considered the argument *support*. A *supporting argument* is one that claims for the premise of another argument. The language of all the possible arguments built from  $\mathbf{L}$  will be noted  $\mathbf{L}_A$ . Consider for instance the argument  $\mathcal{A}_1 = \langle \{th, ps\}, dr \rangle$  which assumes a route to be dangerous because there are known thieves in that area and the security there is poor.

Consider also the supporting argument  $\mathcal{A}_2 = \langle \{upc\}, ps \rangle$  saying that underpaid cops might provide poor security [13].

An argument is *coherent*, wrt. a set of evidence  $E$ , if its claim does not contradict, nor coincides with, any evidence in  $E$ . Then, a coherent argument is *active* if each of its premisses is either evidence or a claim of an active argument. Following the above example, the argument  $\mathcal{A}_1$  is active wrt. the set  $\{th, ps\}$  and also wrt. the set  $\{th, upc\}$ , but it is not wrt.  $\{th\}$  nor  $\{th, ps, \overline{dr}\}$ . Inactive arguments will be depicted in gray.

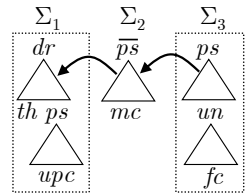
An *argumental structure* (*structure* for short), for a claim  $\delta$ , is a tree of arguments where the root argument claims for  $\delta$ , and every non-root argument supports the parent through a different premise (note there may be unsupported premisses). The arguments  $\mathcal{A}_1$  and  $\mathcal{A}_2$  from the previous example constitute an argumental structure  $\Sigma_1$ , shown on the right, for the claim ‘ $dr$ ’. Structures are depicted as dashed boxes. The box will be omitted when the structure consists of a unique argument.



Further constraints over structures (yielding *well-formed* structures) are imposed in order to ensure a sensible reasoning chain. These avoid arguments attacking each other within a structure, infinite structures, and heterogeneous support for a premise throughout a structure (see [13] for details). A well-formed argumental structure is *active*, wrt. a set of evidence  $E$ , if every argument in it is coherent wrt.  $E$ , and every unsupported premise is evidence in  $E$ . For instance, the previous structure  $\Sigma_1$  claiming a route as being dangerous, is active wrt. the set  $\{th, upc\}$ , but not wrt.  $\{th\}$  because the premise ‘ $upc$ ’ is not evidence nor the claim of another argument in the structure. Neither is  $\Sigma_1$  active wrt.  $\{th, ps\}$  since  $\mathcal{A}_2$  would not be coherent: its claim ‘ $ps$ ’ is redundant wrt. to the evidence. Inactive structures will be depicted in gray.

The *dynamic argumentation framework (DAF)* we will use is a pair  $\langle \mathbf{E}, (\mathbf{W}, \mathbf{R}) \rangle$  composed by a consistent set  $\mathbf{E}$  of *evidence*, a *working set of arguments*  $\mathbf{W}$ , and an *attack relation*  $\mathbf{R} \subseteq \mathbf{L}_A \times \mathbf{L}_A$  between arguments. We restrict the attacks to pairs of arguments with contradictory claims, and at least in one direction. That is, for every pair of arguments  $\mathcal{A}_1$  and  $\mathcal{A}_2$  whose claims are in contradiction, at least one of  $(\mathcal{A}_1, \mathcal{A}_2)$  or  $(\mathcal{A}_2, \mathcal{A}_1)$  will belong to  $\mathbf{R}$ . Contradictory sentences will be noted as  $a$  and  $\bar{a}$ .

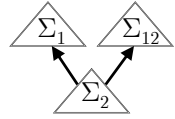
The notion of attack over arguments has a direct correlation to argumental structures: a structure *attacks* another if the root argument of the first attacks any argument of the second. Consider, for instance, the argument  $\mathcal{A}_3 = \langle \{mc\}, \overline{ps} \rangle$  which assumes the security to be good because there are many cops in the area. Consider also the arguments  $\mathcal{A}_4 = \langle \{un\}, ps \rangle$  and  $\mathcal{A}_5 = \langle \{fc\}, un \rangle$  saying that foreign cops might be unacquainted with the place, giving the idea of poor security [13]. Assume that  $\mathcal{A}_3$  attacks  $\mathcal{A}_2$  and  $\mathcal{A}_4$  attacks  $\mathcal{A}_3$ . Then, the structure  $\Sigma_1$  mentioned earlier is attacked by the structure  $\Sigma_2$  (composed by  $\mathcal{A}_3$ ), which is in turn attacked by the structure  $\Sigma_3$  (composed by  $\mathcal{A}_4$  and  $\mathcal{A}_5$ ). This sequence of attacks is depicted on the right.



<sup>1</sup> In [13], the attack relation is defined over the working set of arguments  $\mathbf{W}$ . Since a unified policy for comparing arguments from different agents is needed in a collaborative dialogue setting, here we introduce a slight variation generalizing the attack relation over the universal set of arguments  $\mathbf{L}_A$ .



At any moment, the *active instance* of a DAF is a pair  $(\mathbb{S}, \mathbb{R})$ , where  $\mathbb{S}$  is the set of active argumental structures, and  $\mathbb{R}$  is the resulting attack relation between them. This is equivalent to Dung's definition of abstract argumentation framework [3]. Therefore, classic argumentation semantics can be applied to the active instance. From the previous examples, if we consider the set of evidence  $\{th, upc, mc\}$  then the active instance, depicted on the right, consists of: the structures  $\Sigma_1$  and  $\Sigma_2$  mentioned earlier, and also the structure  $\Sigma_{12}$  composed only by argument  $\mathcal{A}_2$ . Note that  $\Sigma_3$  is not active, and hence does not belong to the active instance. Picking grounded semantics, for instance, the only accepted structure would be  $\Sigma_2$ .



In this work we will assume unique-extension semantics. The arguments belonging to the extension, along with their claims, will be considered *justified* from the DAF, as well as the whole evidence set. From the previous example, the claim ‘ $\overline{ps}$ ’ would be justified. Multiple-extension semantics could also be used, expanding the set of possible ‘justification statuses’ of a claim (e.g. ‘justified’, ‘not-justified’, or ‘undecided’).

### 3 Informal Requirements for Collaborative Dialogue Models

We believe that an ideal collaborative behavior of dialogues should satisfy the following, informally specified, requirements. Note that these requirements address desirable properties of *dialogues instances*, not to be confused with properties aimed to be achieved by *dialogue protocols* (like the ones proposed, for instance, in [8]).

- R<sub>1</sub>**: All the **relevant** information is exposed in the dialogue.
- R<sub>2</sub>**: The exchange of **irrelevant** information is avoided.
- R<sub>3</sub>**: The final conclusion **follows** from all what has been said.

On that basis, we will conduct our analysis of collaborative dialogue behavior in terms of two abstract elements: a *reasoning model* and a *relevance notion*<sup>2</sup>, assuming that the former gives a formal meaning to the word *follows*, and the latter to the word *relevant*. Both elements are domain-dependent and, as we shall see, they are not unattached concepts. It is important to mention that the relevance notion is assumed to work in a context of *complete information* (this will be clarified later).

We believe that the achievement of R<sub>1</sub>-R<sub>3</sub> should lead to achieving other important requirements, listed below. Later in this work we will state the conditions under which this hypothesis actually holds.

- R<sub>4</sub>**: The dialogue should always end.
- R<sub>5</sub>**: Once the dialogue ends, if the agents added all their still private information, and reasoned from there, the previously drawn conclusions should not change.

In the task of simultaneously achieving requirements R<sub>1</sub> and R<sub>2</sub>, in the context of a distributed MAS, a non-trivial problem arises: relevant information distributed in such

<sup>2</sup> The term *relevance* appears in many research areas: *epistemology*, *belief revision*, *economics*, *information retrieval*, etc. In this work we intend to use it in its most general sense, which may be closer to the epistemic one: *pertinence in relation to a given question*, but it should not be tied to any particular interpretation, except for concrete examples given in this work.

a way that none of the parts is relevant by itself. For instance, considering the DAF of Sect. 2 an agent may have an argument for a certain claim but the activating evidence resides in a different agent. This, in principle threatens  $R_1$  since the whole contribution would be left unseen. Besides, any attempt to detect these ‘non-self-relevant’ parts threatens  $R_2$  due to the risk of being mistaken. This could happen for instance, following with the previous example, if the argument is exposed but the activating evidence does not actually exist. There is a tradeoff between requirements  $R_1$  and  $R_2$ .

Because of the nature of collaborative dialogues, we believe  $R_1$  may be mandatory in many application domains, and hence we will seek solutions which achieve it, even at the expense of relegating  $R_2$  a bit. As will be seen later in Sect. 6 the basic idea will be to develop a new relevance notion (which will be called a *potential relevance notion*) able to detect parts of distributed relevant contributions (under the original notion).

## 4 The Dialogue Framework

Three languages are assumed to be involved in a dialogue: the *Knowledge Representation Language*  $\mathcal{L}$  for expressing the information exchanged by the agents, the *Topic Language*  $\mathcal{L}_T$  for expressing the *topic* that gives rise to the dialogue, and the *Outcome Language*  $\mathcal{L}_O$  for expressing the final conclusion (or *outcome*). Also assumed is a language  $\mathcal{L}_I$  for agent identifiers. As usual, a *dialogue* consists of a topic, a sequence of *moves*, and an outcome. In each move an agent makes a *contribution* (exposes a set of knowledge). This is a *public view of dialogue*: agents’ private knowledge is not taken into account yet.

**Definition 1 (Move).** A move is a pair  $\langle id, X \rangle$  where  $id \in \mathcal{L}_I$  is the identifier of the speaker, and  $X \subseteq \mathcal{L}$  is her contribution.

**Definition 2 (Dialogue).** A dialogue is a tuple  $\langle t, \langle m_j \rangle, o \rangle$  where  $t \in \mathcal{L}_T$  is the dialogue topic,  $\langle m_j \rangle$  is a sequence of moves, and  $o \in \mathcal{L}_O$  is the dialogue outcome.

As will be seen in short, in the argumentative approach based on the DAF, the agents will expose arguments and evidence, topics will correspond to claims, and dialogue outcomes might be Yes (justified) or No (not justified).

Note that the dialogue protocol here is very simple, since it consists of only one type of move: to assert a set of knowledge. For this reason, it is not modeled as a typical *dialogue-game protocol*<sup>3</sup>. Instead we are interested in specifying which subsets, from all the dialogues conforming with this basic protocol, represent desirable behaviors.

As anticipated in Sec. 3 we will study the dialogue behavior in terms of two abstract concepts: relevance and reasoning. To that end, an *Abstract Dialogue Framework* ( $\mathcal{D}\mathcal{F}$ ) is introduced, whose aim is to provide an environment for dialogues to take place, and which includes: the languages involved in the dialogue, a set of participating agents, a relevance notion and a reasoning model. An *agent* is represented by an *identifier* and a *private knowledge base* (*kb*), providing in this way a *complete view* of dialogues.

<sup>3</sup> *Dialogue-game protocols* usually specify several rules regarding different types of allowed moves (*locutions*), and the legal ways in which these may be combined during the dialogue. For more detail about formal dialogue-game protocols see, for instance, [7].

**Definition 3 (Agent).** An agent is a pair  $\langle id, K \rangle$ , noted  $K_{id}$ , where  $K \subseteq \mathcal{L}$  is a private finite knowledge base, and  $id \in \mathcal{L}_1$  is an agent identifier.

A *relevance notion* is a criterion for determining, given certain already known information and a topic, whether it would be relevant to add certain other information (*i.e.*, to make a contribution). We emphasize that this criterion works under an assumption of *complete information*, to be contrasted with the situation of a dialogue where each agent is unaware of the private knowledge of the others. This issue will be revisited in Sec. 5. Finally, a *reasoning model* will be understood as a mechanism for drawing a conclusion about a topic, on the basis of an individual knowledge base. The argumentation-based reasoning model, for instance, will determine the justification status of a claim from a given set of evidence and arguments.

**Definition 4 (Abstract Dialogue Framework).** An abstract dialogue framework  $(\mathfrak{D}\mathfrak{F})$  is a tuple  $\langle \mathcal{L}, \mathcal{L}_T, \mathcal{L}_O, \mathcal{L}_1, \mathcal{R}, \Phi, Ag \rangle$  where:

- $\mathcal{L}, \mathcal{L}_T, \mathcal{L}_O$  and  $\mathcal{L}_1$  are the languages involved in the dialogue,
- $Ag$  is a finite set of agents,
- $\mathcal{R} \subseteq 2^{\mathcal{L}} \times 2^{\mathcal{L}} \times \mathcal{L}_T$  is a relevance notion, and
- $\Phi : 2^{\mathcal{L}} \times \mathcal{L}_T \Rightarrow \mathcal{L}_O$  is a reasoning model.

The brief notation  $\langle \mathcal{R}, \Phi, Ag \rangle$  will be also used.

*Notation.* If  $(X, S, t) \in \mathcal{R}$ , we say that  $X$  is a  $t$ -relevant contribution to  $S$  under  $\mathcal{R}$ , and we note it  $X\mathcal{R}_tS$ . When it is clear what relevance notion is being used, we just say that  $X$  is a  $t$ -relevant contribution to  $S$ . For individual sentences  $\alpha$  in  $\mathcal{L}$ , we also use the simpler notation  $\alpha\mathcal{R}_tS$  meaning that  $\{\alpha\}\mathcal{R}_tS$ .

Throughout this work we will refer to the partially instantiated  $\mathfrak{D}\mathfrak{F}$   $\mathfrak{F}^{\text{ar}} = \langle \mathbf{L} \cup \mathbf{L}_A, \mathbf{L}, \{\text{Yes}, \text{No}\}, \mathcal{L}_1, \mathcal{R}_\delta, \Psi, Ag \rangle$ . The languages  $\mathbf{L}$  and  $\mathbf{L}_A$  are the ones of Sect. 2. Hence, in this argumentation-based dialogue framework, the knowledge representation language consists of both evidence and arguments, topics correspond to claims, and outcomes might be Yes or No. As mentioned before, the reasoning model  $\Psi$  determines the justification status of a claim from a given set of evidence and arguments. That is,  $\Psi(K, \delta) = \text{Yes}$  if, and only if, the claim  $\delta$  is justified (under a certain argumentation semantics  $\mathbf{S}$ ) from the DAF  $\langle \mathbf{E}, (\mathbf{W}, \mathbf{R}) \rangle$ , with  $\mathbf{E} \cup \mathbf{W} = K$  and  $\mathbf{R}$  a certain attack relation between arguments. We will sometimes use the more specific notation  $\mathfrak{F}^{\text{ar}}(\mathbf{S}, \mathbf{R})$ . Particularly, the notation  $\mathfrak{F}^{\text{ar}}(\mathbf{G}, \mathbf{R})$  will be used for the instantiation with grounded semantics [3]. For the aim of simplicity, we will assume that the union of the evidence in all knowledge bases is consistent.

There are two different sets of knowledge involved in a dialogue: the *private knowledge* which is the union of the agents' knowledge bases, and the *public knowledge* which is the union of all the contributions already made, up to certain step. The former is a static set, whereas the latter grows as the dialogue progresses.

**Definition 5 (Public Knowledge).** Let  $d$  be a dialogue consisting of a sequence  $\langle \langle id_1, X_1 \rangle \dots \langle id_m, X_m \rangle \rangle$  of moves. The public knowledge associated to  $d$  at step  $j$  ( $j \leq m$ ) is the union of the first  $j$  contributions of the sequence and is noted  $\mathbf{PU}_d^j$  ( $\mathbf{PU}_d^j = X_1 \cup \dots \cup X_j$ ).

**Definition 6 (Private Knowledge).** Let  $\mathfrak{F}$  be a  $\mathcal{D}\mathfrak{F}$  including a set  $Ag$  of agents. The private knowledge associated to  $\mathfrak{F}$  (and to any admissible dialogue under  $\mathfrak{F}$ ) is the union of the knowledge bases of the agents in  $Ag$ , and is noted  $\mathbf{PR}_{\mathfrak{F}}$  ( $\mathbf{PR}_{\mathfrak{F}} = \bigcup_{K_{id} \in Ag} K$ ).

In order to restrict agents' contributions to be subsets of their private knowledge, we define next the set of *admissible dialogues* under a given  $\mathcal{D}\mathfrak{F}$ .

**Definition 7 (Admissible Dialogues).** Let  $\mathfrak{F} = \langle \mathcal{L}, \mathcal{L}_T, \mathcal{L}_O, \mathcal{L}_I, \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathcal{D}\mathfrak{F}$ ,  $t \in \mathcal{L}_T$  and  $o \in \mathcal{L}_O$ . A dialogue  $\langle t, \langle m_j \rangle, o \rangle$  is admissible under  $\mathfrak{F}$  if, and only if, for each move  $m = \langle id, X \rangle$  in the sequence, there is an agent  $K_{id} \in Ag$  such that  $X \subseteq K$ . The set of admissible dialogues under  $\mathfrak{F}$  is noted  $d(\mathfrak{F})$ .

*Remark 1.* For any step  $j$  of any dialogue  $d \in d(\mathfrak{F})$ , it holds that  $\mathbf{PU}_d^j \subseteq \mathbf{PR}_{\mathfrak{F}}$ .

Returning to the notions of relevance and reasoning, it was mentioned in Sec. 3 that these were not unattached concepts: a coherent dialogue must exhibit some connection between them. Assuming a contribution to be relevant whenever its addition alters the conclusion achieved by the reasoning model, as defined below, seems to be a natural connection.

**Definition 8 (Natural Relevance Notion).** Let  $\Phi$  be a reasoning model. The natural relevance notion associated to  $\Phi$  is a relation  $\mathcal{N}_t^\Phi$  such that:

- $X\mathcal{N}_t^\Phi S$  if, and only if,  $\Phi(S, t) \neq \Phi(S \cup X, t)$ .

When  $X\mathcal{N}_t^\Phi S$  we say that  $X$  is a natural  $t$ -relevant contribution to  $S$  under  $\Phi$ .

Hence, in the argumentative approach, the natural relevance notion  $\mathcal{N}_t^\Phi$  detects the change of the “justification status” for a given claim. It will be seen later that this connection can be relaxed, *i.e.*, other relevance notions which are not exactly the natural one, might also be accepted. We distinguish the subclass of  $\mathcal{D}\mathfrak{F}$ s in which the relevance notion is the natural one associated to the reasoning model. We refer to them as *Inquiry Dialogue Frameworks* ( $\mathcal{I}\mathcal{D}\mathfrak{F}$ ), and the relevance notion is omitted in their formal specification.

**Definition 9 (Inquiry Dialogue Framework).** An Inquiry Dialogue Framework ( $\mathcal{I}\mathcal{D}\mathfrak{F}$ ) is a  $\mathcal{D}\mathfrak{F} \langle \mathcal{R}_t, \Phi, Ag \rangle$  where  $\mathcal{R}_t = \mathcal{N}_t^\Phi$ . The brief notation  $\langle \Phi, Ag \rangle$  will be used.

Throughout this work we will refer to the partially instantiated  $\mathcal{I}\mathcal{D}\mathfrak{F}$   $\mathcal{I}^{\text{ar}} = \langle \Psi, Ag \rangle$ . As with  $\mathcal{D}\mathfrak{F}$ s, we will also use the notation  $\mathcal{I}^{\text{ar}}(\mathbf{S}, \mathbf{R})$  for specifying a particular argumentation semantics and a particular argument attack relation.

## 5 Utopian Collaborative Semantics

A *semantics* for a  $\mathcal{D}\mathfrak{F}$  is a subset of the admissible dialogues representing a particular dialogue behavior. We are interested in specifying which, from all the admissible dialogues under a given  $\mathcal{D}\mathfrak{F}$ , have an acceptable collaborative behavior. In Sec. 3 we

---

<sup>4</sup> The term *Inquiry* is inspired on the popularized typology of dialogues proposed in [14], since we believe that the natural relevance notion captures the essence of this type of interaction: *collaboration to answer some question*. However, the term will be used in a broader sense here, since nothing is assumed regarding the degree of knowledge of the participants.

identified three requirements,  $R_1$ - $R_3$ , to be ideally achieved by collaborative dialogue systems. In this section, we will define an *Utopian Collaborative Semantics* which gives a formal characterization of such ideal behavior. In order to translate requirements  $R_1$ - $R_3$  into a formal specification, some issues need to be considered first.

In particular, the notion of *relevant contribution* needs to be adjusted. On the one hand, there may be contributions which does not qualify as relevant but it would be adequate to allow. To understand this, it should be noticed that, since relevance notions are related to reasoning models, and reasoning models may be non-monotonic, then it is possible for a contribution to contain a relevant subset, without being relevant itself. For instance, in the context of the  $\mathcal{J}^{ar}(\mathbf{G}, \mathbf{R})$  framework, an active argumental structure  $\Sigma_1$  would be a natural  $c1(\Sigma_1)$ -relevant contribution to the empty set, but if we added an active structure  $\Sigma_2$  attacking  $\Sigma_1$ , then it would not. The possibility of some other agent having, for instance, an active structure  $\Sigma_3$  attacking  $\Sigma_2$ , explains why it would be useful to allow the whole contribution consisting of both  $\Sigma_1$  and its attacker  $\Sigma_2$  (and all the supporting evidence). In these cases, we say that the relevance notion fails to satisfy *left-monotonicity* and that the whole contribution is *weakly relevant*<sup>5</sup>. The formal definitions are given below.

**Definition 10 (Left Monotonicity).** *Let  $\mathcal{R}_t$  be a relevance notion. We say that  $\mathcal{R}_t$  satisfies left monotonicity if, and only if, the following condition holds:*

- if  $X\mathcal{R}_tS$  and  $X \subseteq Y$  then  $Y\mathcal{R}_tS$ .

**Definition 11 (Weak Contribution).** *Let  $\mathcal{R}_t$  be a relevance notion. We say that  $X$  is a weak  $t$ -relevant contribution to  $S$  if, and only if, there exists  $Y \subseteq X$  such that  $Y\mathcal{R}_tS$ .*

On the other hand, there may be contributions which qualify as relevant but they are not *purely* relevant. For example, the argument  $\langle \{b\}, a \rangle$  together with the set of evidence  $\{b, e\}$  constitute a natural ‘ $a$ ’-relevant contribution to the empty set, although the evidence ‘ $e$ ’ is clearly irrelevant. These impure relevant contributions must be avoided in order to obey requirement  $R_2$ . For that purpose, *pure relevant contributions* impose a restriction over weak relevant ones, disallowing absolutely irrelevant sentences within them, as defined below.

**Definition 12 (Pure Contribution).** *Let  $\mathcal{R}_t$  be a relevance notion, and  $X$  a weak  $t$ -relevant contribution to  $S$ . We say that  $X$  is a pure  $t$ -relevant contribution to  $S$  if, and only if, the following condition holds for all  $\alpha \in X$ :*

- there exists  $Y \subset X$  such that  $\alpha\mathcal{R}_t(S \cup Y)$ .

Finally, it has been mentioned that the relevance notion works under an assumption of *complete information*, and thus it will be necessary to inspect the private knowledge of the others for determining the actual relevance of a given move. Now we are able to give a formal interpretation of requirements  $R_1$ - $R_3$  in terms of the  $\mathcal{D}\mathcal{F}$  elements:

---

<sup>5</sup> The term *weak relevance* is used in [12] in a different sense, which should not be related to the one introduced here.

**Definition 13 (Utopian Collaborative Semantics).** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathfrak{D}\mathfrak{F}$ . A dialogue  $d = \langle t, \langle m_j \rangle, o \rangle \in d(\mathfrak{F})$  belongs to the Utopian Collaborative Semantics for  $\mathfrak{F}$  (noted  $Utopian(\mathfrak{F})$ ) if, and only if:

**Correctness:** if  $m_j$  is the last move in the sequence, then  $\Phi(\mathbf{PU}_d^j, t) = o$ .

**Global Progress:** for each move  $m_j = \langle id_j, X_j \rangle$  in the sequence, there exists  $Y \subseteq \mathbf{PR}_{\mathfrak{F}}$  such that  $X_j \subseteq Y$  and  $Y$  is a pure  $t$ -relevant contribution to  $\mathbf{PU}_d^{j-1}$ .

**Global Completeness:** if  $m_j$  is the last move in the sequence, then  $\mathbf{PR}_{\mathfrak{F}}$  is not a weak  $t$ -relevant contribution to  $\mathbf{PU}_d^j$ .

Requirement  $R_3$  is achieved by the *Correctness* condition, which states that the dialogue outcome coincides with the application of the reasoning model to the public knowledge at the final step of the dialogue (*i.e.*, the outcome of the dialogue can be obtained by reasoning from all that has been said). In the case of the  $\mathfrak{J}^{\text{ar}}$  framework, for instance, this means that the dialogue outcome is Yes if, and only if, the claim (topic) results justified considering all the arguments and evidence exposed during the dialogue. Requirement  $R_2$  is achieved by the *Global Progress* condition, which states that each move in the sequence is part of a distributed pure relevant contribution to the public knowledge generated so far. Finally, requirement  $R_1$  is achieved by the *Global Completeness* condition, which states that there are no more relevant contributions, not even distributed among different knowledge bases, after the dialogue ends. Notice that the three conditions are simultaneously satisfiable by any  $\mathfrak{D}\mathfrak{F}$  and topic, *i.e.*, there always exists at least one dialogue which belongs to this semantics, as stated in the following proposition.

**Proposition 1 (Satisfiability).** For any  $\mathfrak{D}\mathfrak{F} \mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$ , the set  $Utopian(\mathfrak{F})$  contains at least one dialogue over each possible topic in  $\mathcal{L}_T$ .

Furthermore, any sequence of moves satisfying *global progress* can be completed to a dialogue belonging to the semantics. This means that a system implementation under this semantics would not need to do *backtracking*. Although this property is useless for the case of the utopian semantics which, as will be seen in short, is not implementable in a distributed system, it will be useful in the case of the two practical semantics that will be presented in Sec. 6.

**Definition 14.** A dialogue  $d_2$  over a topic  $t$  is a continuation of a dialogue  $d_1$  over the same topic  $t$  if, and only if, the sequence of moves of  $d_2$  can be obtained by adding zero or more elements to the sequence of moves of  $d_1$ .

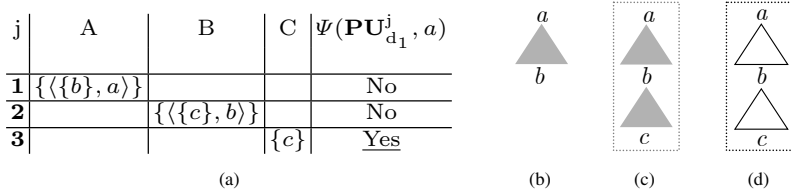
**Proposition 2 (No Backtracking).** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathfrak{D}\mathfrak{F}$ , and  $d_1 \in d(\mathfrak{F})$ . If  $d_1$  satisfies *global progress* under  $\mathfrak{F}$ , then there exists a dialogue  $d_2 \in Utopian(\mathfrak{F})$  which is a continuation of  $d_1$ .

Note that the truth of the previous statements (regarding *satisfiability* and *no backtracking*) comes from the following facts, which can be easily proven: (1) if *global completeness* is not achieved, then there exists at least one possible move that can be added to the sequence according to *global progress*, and (2) the *correctness* condition is orthogonal to the other two. Next, an illustrative example of the dialogues generated under the Utopian Semantics is given.

*Example 1.* Consider an instance of the  $\mathcal{J}^{\text{ar}}(\mathbf{G}, \mathbf{R})$  framework, where the knowledge bases of the agents in the set  $Ag$  are the following:

$$K_A = \{\langle\{b\}, a\rangle, e\rangle, \quad K_B = \{\langle\{c\}, b\rangle, \langle\{d\}, b\rangle, f\rangle \quad \text{and} \quad K_C = \{c, g\}.$$

The dialogue  $d_1$  shown in Fig. 1 over topic ‘ $a$ ’, and also all the permutations of its moves with the same topic and outcome, belong to the Utopian Semantics for the framework. The chart (a) traces the dialogue, showing the partial results of reasoning from the public knowledge so far generated. The last of these results (underlined) is the dialogue outcome. The evolution of the public knowledge is depicted in subfigures (b) through (d). At the first step of the dialogue, an inactive argument is added (b). The second step adds another inactive argument, supporting the first one (c). Finally, the supporting evidence is made available, and the whole structure becomes active (d), yielding to the justification of claim  $a$ .



**Fig. 1.** A dialogue under the Utopian Collaborative Semantics

An essential requirement of dialogue systems is ensuring the termination of the generated dialogues. This is intuitively related to requirement  $R_2$  (achieved by *global progress*) since it is expected that agents will eventually run out of relevant contributions, given that their private knowledge bases are finite. This is actually true as long as the relevance notion satisfies an intuitive property, defined below, which states that a relevant contribution must add some new information to the public knowledge.

**Definition 15 (Novelty).** A relevance notion  $\mathcal{R}_t$  satisfies novelty if, and only if, the following condition holds:

- if  $X\mathcal{R}_tS$  then  $X \not\subseteq S$ .

**Proposition 3 (Termination).** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathcal{D}\mathfrak{F}$ , and  $d = \langle t, \langle m_j \rangle, o \rangle \in d(\mathfrak{F})$ . If the notion  $\mathcal{R}_t$  satisfies novelty and dialogue  $d$  satisfies global progress under  $\mathfrak{F}$ , then  $\langle m_j \rangle$  is a finite sequence of moves.

It is easy to see that any natural relevance notion satisfies novelty, since it is not possible for the conclusion achieved by the reasoning model to change without changing the topic nor the knowledge base.

**Proposition 4.** For any reasoning model  $\Phi$ , it holds that its associated natural relevance notion,  $\mathcal{N}_t^\Phi$ , satisfies novelty.

Another desirable property of collaborative dialogue models is ensuring it is not possible to draw different conclusions, for the same set of agents and topic. In other words,

from the entirety of the information, it should be possible to determine the outcome of the dialogue, no matter what sequence of steps are actually performed<sup>6</sup>. Furthermore, this outcome should coincide with the result of applying the reasoning model to the private knowledge involved in the dialogue. We emphasize that this is required for *collaborative* dialogues (and probably not for non-collaborative ones). For instance, in Ex. 11 all the possible dialogues under the semantics end up justifying the claim, which is also justified from  $K_A \cup K_B \cup K_C$ . This is intuitively related to requirements  $R_1$  (achieved by *global completeness*) and  $R_3$  (achieved by *correctness*) since it is expected that the absence of relevant contributions implies that the current conclusion cannot be changed by adding more information. This is actually true as long as the relevance notion is the natural one associated to the reasoning model, or a *weaker* one, as stated below.

**Definition 16 (Stronger Relevance Notion).** Let  $\mathcal{R}_t$  and  $\mathcal{R}'_t$  be relevance notions. We say that the notion  $\mathcal{R}_t$  is stronger or equal than  $\mathcal{R}'_t$  if, and only if, the following holds:

- if  $X\mathcal{R}_tS$  then  $X\mathcal{R}'_tS$  (i.e.,  $\mathcal{R}_t \subseteq \mathcal{R}'_t$ ).

We will also say that  $\mathcal{R}'_t$  is weaker or equal than  $\mathcal{R}_t$ .

Observe that here we use the term *weaker*, as the opposite of *stronger*, denoting a binary relation between relevance notions, and this should not be confused with its previous use in Def. 11 of *weak relevant contribution*.

**Proposition 5 (Outcome Determinism).** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathfrak{D}\mathfrak{F}$  and  $d = \langle t, \langle m_j \rangle, o \rangle \in d(\mathfrak{F})$ . If  $d$  satisfies correctness and global completeness under  $\mathfrak{F}$ , and  $\mathcal{R}_t$  is weaker or equal than  $\mathcal{N}_t^\Phi$ , then  $o = \Phi(\mathbf{PR}_{\mathfrak{F}}, t)$ .

For example, a relevance notion which detects the generation of new justified arguments (in the usual sense) for a given claim, would be *weaker* than the natural one. It is easy to see that this weaker relevance notion would also achieve *outcome determinism*.

The following corollaries summarize the results regarding the utopian semantics for  $\mathfrak{D}\mathfrak{F}$ s, and also for the particular case of  $\mathfrak{I}\mathfrak{D}\mathfrak{F}$ s. Clearly, these results are inherited respectively by  $\mathfrak{F}^{\text{ar}}$  and  $\mathfrak{J}^{\text{ar}}$ .

**Corollary 1.** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathfrak{D}\mathfrak{F}$ . The dialogues in  $\text{Utopian}(\mathfrak{F})$  satisfy termination and outcome determinism, provided that the relevance notion  $\mathcal{R}_t$  satisfies novelty and is weaker or equal than  $\mathcal{N}_t^\Phi$ .

**Corollary 2.** Let  $\mathfrak{J}$  be an  $\mathfrak{I}\mathfrak{D}\mathfrak{F}$ . The dialogues in  $\text{Utopian}(\mathfrak{J})$  satisfy termination and outcome determinism.

It is clear that Def. 13 of the Utopian Collaborative Semantics is not constructive, since both *global progress* and *global completeness* are expressed in terms of the private knowledge  $\mathbf{PR}_{\mathfrak{F}}$ , which is not entirely available to any of the participants. The following example shows that, it is not only not constructive, but also in many cases not even implementable in a distributed MAS.

<sup>6</sup> This property, which we will call *outcome determinism*, has been studied in various works under different names. For instance in [9] it was called *completeness*. Notice that we use that term for another property, which is not the same but is related to the one under discussion.



*Example 2.* Consider the inquiry framework instance from Ex. 1. The dialogue  $d_2$  shown in Fig. 2 does not belong to the Utopian Semantics, since step 2 violates *global progress*. However, it would not be possible to design a dialogue system which allows  $d_1$  (from Ex. 1) but disallows  $d_2$ , since agent B cannot know in advance that ‘c’, rather than ‘d’, holds.

The undesired situation is caused by a relevant contribution distributed among several agents, in such a way that none of the parts is relevant by itself, leading to a tradeoff between requirements  $R_1$  and  $R_2$  (i.e., between *global progress* and *global completeness*). In the worst case, each sentence of the contribution resides in a different agent. Thus, to avoid such situations, it would be necessary for the relevance notion to warrant that every relevant contribution contains at least one individually relevant sentence. When this happens, we say that the relevance notion satisfies *granularity*, defined next.

**Definition 17 (Granularity).** Let  $\mathcal{R}_t$  be a relevance notion. We say that  $\mathcal{R}_t$  satisfies granularity if, and only if, the following holds:

- if  $X\mathcal{R}_tS$  then there exists  $\alpha \in X$  such that  $\alpha\mathcal{R}_tS$ .

Unfortunately, the relevance notions we are interested in, fail to satisfy granularity. It does not hold in general for the natural notions associated to deductive inference mechanisms. In particular, it has been shown in Ex. 2 that it does not hold for  $\mathcal{N}_\delta^\Psi$ .

## 6 Practical Collaborative Semantics

The lack of granularity of relevance notions motivates the definition of alternative semantics which approach the utopian one, and whose distributed implementation is viable. The simplest approach is to relax requirement  $R_1$  by allowing distributed relevant contributions to be missed, as follows.

**Definition 18 (Basic Collaborative Semantics).** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathfrak{D}\mathfrak{F}$ . A dialogue  $d = \langle t, \langle m_j \rangle, o \rangle \in d(\mathfrak{F})$  belongs to the Basic Collaborative Semantics for  $\mathfrak{F}$  (noted *Basic*( $\mathfrak{F}$ )) if, and only if, the following conditions, as well as **Correctness** (Def. 13), hold:

**Local Progress:** for each move  $m_j = \langle id_j, X_j \rangle$  in the sequence,  $X_j$  is a pure  $t$ -relevant contribution to  $\mathbf{PU}_d^{j-1}$ .

**Local Completeness:** if  $m_j$  is the last move in the sequence, then it does not exist an agent  $K_{id} \in Ag$  such that  $K$  is a weak  $t$ -relevant contribution to  $\mathbf{PU}_d^j$ .

In the above definition, requirement  $R_2$  is achieved by the *local progress* condition which states that each move in the sequence constitutes a pure relevant contribution to the public knowledge generated so far. Notice that this condition implies *global progress* (enunciated in Sec. 5), as stated below.

j	A	B	C	$\Psi(\mathbf{PU}_{d_2}^j, a)$
1	{ $\{b\}, a$ }			No
2		{ $\{d\}, b$ }		No
3		{ $\{c\}, b$ }		No
4			{ $c$ }	Yes

Fig. 2. A dialogue violating the Utopian Collaborative Semantics

**Proposition 6.** *Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathcal{D}\mathfrak{F}$ , and  $d \in d(\mathfrak{F})$ . If the dialogue  $d$  satisfies local progress, then it satisfies global progress under  $\mathfrak{F}$ .*

Requirement  $R_1$  is now compromised. The *local completeness* condition states that each agent has no more relevant contributions to make after the dialogue ends. Unless the relevance notion satisfies granularity, this is not enough for ensuring global completeness (enunciated in Sec. 5), since there could be a relevant contribution distributed among several agents, in such a way that none of the parts is relevant by itself.

**Proposition 7.** *Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathcal{D}\mathfrak{F}$ , and  $d \in d(\mathfrak{F})$ . If the dialogue  $d$  satisfies global completeness, then it satisfies local completeness under  $\mathfrak{F}$ . The reciprocal holds if, and only if, the relevance notion  $\mathcal{R}_t$  satisfies granularity.*

As a result, requirement  $R_4$  (termination) is achieved, given the same condition as in Sec. 5, whereas requirement  $R_5$  (outcome determinism) cannot be warranted. These results are summarized in the corollary below. Clearly, these results are inherited by  $\mathfrak{F}^{\text{ar}}$  and  $\mathcal{J}^{\text{ar}}$ .

**Corollary 3.** *Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathcal{D}\mathfrak{F}$ . The dialogues in  $\text{Basic}(\mathfrak{F})$  satisfy termination, provided that the relevance notion  $\mathcal{R}_t$  satisfies novelty.*

**Corollary 4.** *Let  $\mathcal{J}$  be an  $\mathcal{J}\mathcal{D}\mathfrak{F}$ . The dialogues in  $\text{Basic}(\mathcal{J})$  satisfy termination.*

Considering the same scenario as in Ex. 1 it is easy to see that the only possible dialogue under the Basic Semantics is the empty one (*i.e.*, no moves are performed), with outcome No. A more interesting example is shown next.

*Example 3.* Consider an instance of the  $\mathcal{J}^{\text{ar}}(\mathbf{G}, \mathbf{R})$  framework, where the knowledge bases of the agents in the set  $Ag$  are the following:

$$K_A = \{\langle \{b\}, a \rangle, b, g\}, K_B = \{\langle \{e\}, \bar{a} \rangle, \langle \{f\}, e \rangle, f, g\}, \text{ and } K_C = \{\langle \{g\}, \bar{a} \rangle, \bar{e}\}.$$

Also consider that both  $\langle \{e\}, \bar{a} \rangle$  and  $\langle \{g\}, \bar{a} \rangle$  attack  $\langle \{b\}, a \rangle$ , but not viceversa. The private knowledge is depicted in Fig. 3(a). The dialogue  $d_3$  traced in Fig. 3(b), over topic  $a$ , belongs to the Basic Semantics for the  $\mathcal{J}\mathcal{D}\mathfrak{F}$  instantiated above. The evolution of the public knowledge is depicted in figures 3(c) through 3(e). At the first step, an active argument for ‘ $a$ ’ is added (c). At the second step, an attacking structure is added (d). Finally, the attacking structure is deactivated due to a supporting argument becoming inconsistent wrt. new evidence (e). Note that *global completeness* is not achieved, since there still exists a distributed relevant contribution when the dialogue ends:  $\{\langle \{g\}, \bar{a} \rangle, g\}$ . Consequently, outcome determinism is not achieved: the outcome is Yes whereas the result of reasoning from the private knowledge is No.

In Sec. 3 we argued that requirement  $R_1$  may be mandatory in many domains, but the Basic Semantics does not achieve it unless the relevance notion satisfies granularity, which does not usually happen. In order to make up for this lack of granularity, we propose to build a new notion (say  $\mathcal{P}$ ) based on the original one (say  $\mathcal{R}$ ) which ensures that, in the presence of a distributed relevant contribution under  $\mathcal{R}$ , *at least one* of the parts will be relevant under  $\mathcal{P}$ . We will say that  $\mathcal{P}$  is a *potential relevance notion* for  $\mathcal{R}$ , since its aim is to detect contributions that could be relevant within certain *context*,

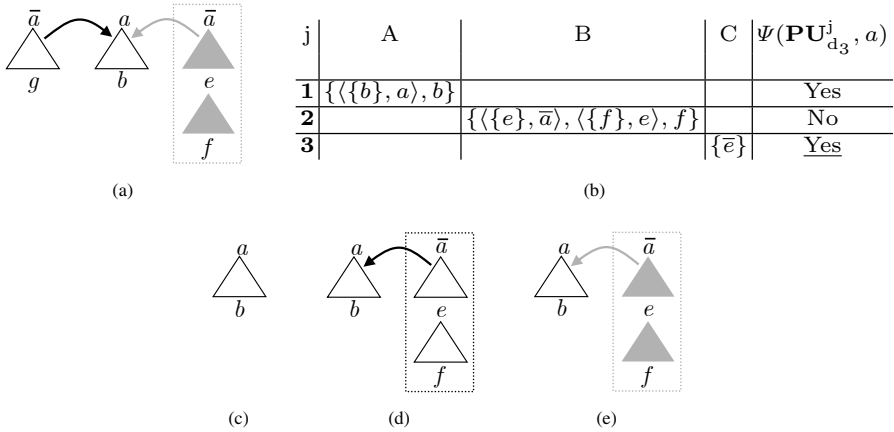


Fig. 3. A dialogue under the Basic Collaborative Semantics

but it is uncertain whether that context actually exists or not. Observe that the context is given by other agents’ private knowledge, which has not been exposed yet.

Below we define the binary relation (“is a potential for”) between relevance notions, and also its propagation to  $\mathcal{D}\mathfrak{F}$ s. Clearly, if a relevance notion already satisfies granularity then nothing needs to be done. Indeed, it would work as a potential relevance notion for itself.

**Definition 19 (Potential Relevance Notion).** Let  $\mathcal{R}_t$  and  $\mathcal{P}_t$  be relevance notions. We say that  $\mathcal{P}_t$  is a potential (relevance notion) for  $\mathcal{R}_t$  if, and only if, the following conditions hold:

1.  $\mathcal{R}_t$  is stronger or equal than  $\mathcal{P}_t$ , and
2. if  $X\mathcal{R}_tS$  then there exists  $\alpha \in X$  such that  $\alpha\mathcal{P}_tS$ .

If  $X\mathcal{P}_tS$  and  $\mathcal{P}_t$  is a potential for  $\mathcal{R}_t$ , we say that  $X$  is a potential  $t$ -relevant contribution to  $S$  under  $\mathcal{R}_t$ .

**Definition 20 (Potential Dialogue Framework).** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  and  $\mathfrak{F}^* = \langle \mathcal{P}_t, \Phi, Ag \rangle$  be  $\mathcal{D}\mathfrak{F}$ s. We say that  $\mathfrak{F}^*$  is a potential (framework) for  $\mathfrak{F}$  if, and only if,  $\mathcal{P}_t$  is a potential relevance notion for  $\mathcal{R}_t$ .

**Proposition 8.** If the relevance notion  $\mathcal{R}_t$  satisfies granularity, then  $\mathcal{R}_t$  is a potential relevance notion for itself.

Now we will show a more interesting potential relevance notion, in the context of the  $\mathfrak{J}^{\text{ar}}$  framework. The basic idea is to detect contributions that would be relevant given a certain situation (i.e., a certain set of evidence). To that end, we first introduce the concept of *abduction set* associated to a given claim  $\delta$  and a given set  $K$ . This abduction set reflects how the current situation (represented by the evidence in  $K$ ) could be minimally expanded in order to change the justification status of  $\delta$ .

**Definition 21 (Abduction Set).** Let  $K \subseteq \mathbf{L} \cup \mathbf{L}_A$  and  $\delta \in \mathbf{L}$ . The abduction set of  $\delta$  from  $K$ , noted  $\mathcal{AB}(K, \delta)$ , is defined as:

$$\mathcal{AB}(K, \delta) = \left\{ E \subseteq \mathbf{L} : E \text{ is consistent wrt. the evidence in } K, \text{ and } E \text{ is a minimal natural } \delta\text{-relevant contribution to } K. \right\}$$

*Example 4.* Consider the  $\mathfrak{J}^{\text{ar}}$  framework. In the chart of Fig. 4 the second column shows the abduction set of claim “ $a$ ”, from the set  $K$  on the first column. In the last case, assume that the argument  $\langle \{e\}, \bar{a} \rangle$  attacks the argument  $\langle \{b\}, a \rangle$ , but not viceversa.

Now we are able to introduce an *abductive relevance notion*  $\mathcal{A}_\delta^\Psi$ . Under this notion, a set  $X$  is considered an  $\delta$ -relevant contribution to  $K$  if, and only if, its addition generates a new element in the abduction set of  $\delta$  from  $K$ . This means that a new potential situation in which the justification status of  $\delta$  would change has arisen. It can be shown (proof is omitted due to space reasons) that  $\mathcal{A}_\delta^\Psi$  is a potential relevance notion for  $\mathcal{N}_\delta^\Psi$ .

$K$	$\mathcal{AB}(K, a)$
$\{\}$	$\{\{a\}\}$
$\{\langle \{b\}, a \rangle\}$	$\{\{a\}\{b\}\}$
$\{\langle \{b\}, a \rangle, b\}$	$\{\{\}\}$
$\{\langle \{b\}, a \rangle, b, \langle \{e\}, \bar{a} \rangle, \langle \{f\}, e \rangle, f\}$	$\{\{a\}\{\bar{e}\}\}$

**Fig. 4.** Some abduction set examples

**Definition 22 (Abductive Relevance).** Let  $K \subseteq \mathbf{L} \cup \mathbf{L}_A$  and  $\delta \in \mathbf{L}$ . A set  $X \subseteq \mathbf{L} \cup \mathbf{L}_A$  is an  $\delta$ -relevant contribution to  $K$  under  $\mathcal{A}_\delta^\Psi$  if, and only if, there exists  $E \subseteq \mathbf{L}$  such that the following conditions hold:

1.  $E \in \mathcal{AB}(K \cup X, \delta)$ , and
2.  $E \notin \mathcal{AB}(K, \delta)$ .

**Proposition 9.** The notion  $\mathcal{A}_\delta^\Psi$  is a potential relevance notion for  $\mathcal{N}_\delta^\Psi$ .

Returning to the semantics definition, the idea is to use the potential framework under the Basic Semantics, resulting in a new semantics for the original framework. Next we introduce the *Full Collaborative Semantics*, which is actually a family of semantics: each possible potential  $\mathfrak{D}\mathfrak{F}$  defines a different semantics of the family.

**Definition 23 (Full Collaborative Semantics).** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathfrak{D}\mathfrak{F}$ . A dialogue  $d = \langle t, \langle m_j \rangle, o \rangle \in d(\mathfrak{F})$  belongs to the Full Collaborative Semantics for  $\mathfrak{F}$  (noted  $\text{Full}(\mathfrak{F})$ ) if, and only if,  $d \in \text{Basic}(\mathfrak{F}^*)$  for some  $\mathfrak{D}\mathfrak{F} \mathfrak{F}^* = \langle \mathcal{P}_t, \Phi, Ag \rangle$  which is a potential for  $\mathfrak{F}$ . We will also use the more specific notation  $d \in \text{Full}(\mathfrak{F}, \mathcal{P}_t)$ .

In this way, each agent would be able to autonomously determine that she has no more potential relevant contributions to make, ensuring there cannot be any distributed relevant contribution when the dialogue ends, and hence achieving  $R_1$ . In other words, achieving local completeness under the potential relevance notion implies achieving global completeness under the original one, as stated below.

**Proposition 10.** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  and  $\mathfrak{F}^* = \langle \mathcal{P}_t, \Phi, Ag \rangle$  be  $\mathfrak{D}\mathfrak{F}$ s such that  $\mathfrak{F}^*$  is a potential for  $\mathfrak{F}$ , and  $d \in d(\mathfrak{F})$ . If dialogue  $d$  satisfies local completeness under  $\mathfrak{F}^*$ , then it satisfies global completeness under  $\mathfrak{F}$ .

Requirement  $\mathcal{R}_2$  is now compromised, since the context we have mentioned may not exist. In other words, achieving local progress under the potential relevance notion does not ensure achieving global progress under the original one. The challenge is to design *good* potential relevance notions which considerably reduce the amount of cases in which a contribution is considered potentially relevant but, eventually, it is not. Observe that a relevance notion which considers any sentence of the language as relevant, works as a potential for any given relevance notion, but it is clearly not a good one.

Next we summarize the results for the dialogues generated under the Full Collaborative Semantics. By achieving global completeness these dialogues achieve outcome determinism under the same condition as before. Although global progress is not achieved under the original relevance notion, it is achieved under the potential one, and thus termination can be ensured as long as the latter satisfies novelty. Clearly, these results are inherited by  $\mathfrak{F}^{\text{ar}}$  and  $\mathfrak{J}^{\text{ar}}$ .

**Corollary 5.** *Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathfrak{D}\mathfrak{F}$ , and  $\mathcal{P}_t$  a potential for  $\mathcal{R}_t$ . The dialogues in  $\text{Full}(\mathfrak{F}, \mathcal{P}_t)$  satisfy termination and outcome determinism, provided that  $\mathcal{P}_t$  satisfies novelty and  $\mathcal{R}_t$  is weaker or equal than  $\mathcal{N}_t^\Phi$ .*

**Corollary 6.** *Let  $\mathfrak{J} = \langle \Phi, Ag \rangle$  be an  $\mathfrak{J}\mathfrak{D}\mathfrak{F}$ , and  $\mathcal{P}_t$  a potential for  $\mathcal{N}_t^\Phi$ . The dialogues in  $\text{Full}(\mathfrak{J}, \mathcal{P}_t)$  satisfy termination and outcome determinism, provided that  $\mathcal{P}_t$  satisfies novelty.*

*Example 5.* Both dialogues  $d_1$  and  $d_2$ , presented in Ex. 1 and Ex. 2 respectively, belong to  $\text{Full}(\mathfrak{J}^{\text{ar}}, \mathcal{A}_\delta^{\text{ar}})$ . Also belongs to this semantics the dialogue which results from  $d_2$  by interchanging steps 2 and 3, or by merging these two steps together in a single one. Note that all these dialogues achieve *global completeness*, although *global progress* is achieved only by dialogue  $d_1$ .

*Example 6.* The dialogue  $d_3$  from Ex. 3 can be completed according to  $\text{Full}(\mathfrak{J}^{\text{ar}}, \mathcal{A}_\delta^{\text{ar}})$ , as shown in Fig. 5(a). The fifth column of the chart shows the evolution of the abduction set of the claim “ $a$ ” from the generated public knowledge. An additional step 0 is added, in order to show the initial state of this abduction set. At step 4 an attacking, for the meantime inactive, argument is added (5(b)). This generates a new potential situation in which the claim ‘ $a$ ’ would not be justified any more. At step 5 the previous situation is realized, activating the attack and leaving the claim ‘ $a$ ’ not justified (5(c)). Note that other dialogues also belong to the Full Collaborative Semantics, since the first three steps do not actually need to be natural relevant contributions. For instance, agent A could expose the argument  $\langle \{b\}, a \rangle$  and then, in the next step, the supporting evidence. Moreover, agent B could make her attack while agent A’s argument is still inactive. In that moment, the element  $\{b, \bar{e}\}$  would be added to the abduction set.

It is important to note the existence of alternative potential relevance notions which may be also adequate, and which may cause variations in the behavior of the dialogue. For instance, a variant of the abductive relevance notion defined earlier is to consider a contribution as relevant if its addition either adds or deletes an element of the abduction set. The latter case, deletion, could be seen as discarding a possible explanation before it is actually realized (or activated). For instance, assume from Ex. 6 that agent A exposes just the argument  $\langle \{b\}, a \rangle$  without the activating evidence. Before the activation,

j	A	B	C	$\mathcal{AB}(\mathbf{PU}_{d_4}^j, a)$	$\Psi(\mathbf{PU}_{d_4}^{\text{step}}, a)$
0				$\{\{a\}\}$	No
1	$\{\{\{b\}, a\}, b\}$			$\{\{\bar{a}\}\}$	Yes
2		$\{\{\{e\}, \bar{a}\}, \{\{f\}, e\}, f\}$		$\{\{a\}, \{\bar{e}\}\}$	No
3			$\{\bar{e}\}$	$\{\{\bar{a}\}\}$	Yes
4			$\{\{\{g\}, \bar{a}\}\}$	$\{\{\bar{a}\}, \{g\}\}$	Yes
5		$\{g\}$		$\{\{a\}\}$	No

(a)

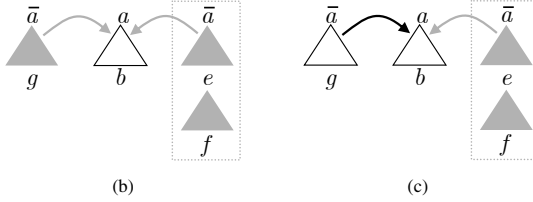


Fig. 5. A dialogue under the Full Collaborative Semantics

it would be possible for agent C to make an attack by exposing the argument  $\langle\{g\}, \bar{a}\rangle$  together with the supporting evidence  $\{g\}$ . Observe that that exact sequence of steps is not allowed under the abductive notion defined earlier, in Def. 22 since no element is added to the abduction set, instead the element  $\{b\}$  is deleted. Under that notion, it would be necessary for agent A to activate her argument before agent C can attack it. It is easy to see that the alternative notion which considers not only the expansion but also the reduction of the abduction set, may in some cases lead to shorter dialogues.

Results regarding *satisfiability* and *no-backtracking* also hold under the two practical semantics we have presented in this section, as stated below.

**Proposition 11.** For any  $\mathcal{DF} \mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$ , each one of the sets  $Basic(\mathfrak{F})$  and  $Full(\mathfrak{F}, \mathcal{P}_t)$ , contains at least one dialogue over each possible topic in  $\mathcal{L}_T$ .

**Proposition 12.** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  and  $\mathfrak{F}^* = \langle \mathcal{P}_t, \Phi, Ag \rangle$  be  $\mathcal{DF}$ s such that  $\mathfrak{F}^*$  is a potential for  $\mathfrak{F}$ , and let  $d_1 \in d(\mathfrak{F})$ . If  $d_1$  satisfies local progress under  $\mathfrak{F}$  ( $\mathfrak{F}^*$ ), then there exists a dialogue  $d_2 \in Basic(\mathfrak{F})$  ( $d_2 \in Full(\mathfrak{F}, \mathcal{P}_t)$ ) which is a continuation of  $d_1$ .

Finally, a result showing the relation among the three collaborative semantics, for the case in which the relevance notion satisfies *granularity*, is stated.

**Proposition 13.** Let  $\mathfrak{F} = \langle \mathcal{R}_t, \Phi, Ag \rangle$  be a  $\mathcal{DF}$ . If the relevance notion  $\mathcal{R}_t$  satisfies granularity, then it holds that:  $Basic(\mathfrak{F}) = Full(\mathfrak{F}, \mathcal{R}_t) \subseteq Utopian(\mathfrak{F})$ .

To sum up, we have defined three collaborative semantics for a  $\mathcal{DF}$ . The Utopian Semantics describes an idealistic, in most cases impractical behavior of a collaborative dialogue. Its usefulness is theoretical. It is approximated, in different ways, by the other two practical semantics. The Basic Semantics, on the other side, describes a straightforward implementable behavior of a collaborative dialogue. The weak point of this semantics is not ensuring *global completeness* (neither *outcome determinism*, thus).

The Full Collaborative Semantics is actually a family of semantics: each potential relevance notion  $\mathcal{P}_t$  associated to  $\mathcal{R}_t$  defines a semantics of the family. Thus, the constructiveness of these semantics is reduced to the problem of finding a potential relevance notion for  $\mathcal{R}_t$ . These semantics succeed in achieving *global completeness*, at the price of allowing moves which may not be allowed by the Utopian Semantics. The goodness of a given potential relevance notion increases as it minimizes the amount of such moves.

## 7 Related Work

There are some works particularly related to our proposed approach, due to any of the following: (a) an explicit treatment of the notion of relevance in dialogue, (b) the search of the *global completeness* property, as we called it in this work, or (c) a tendency to examine general properties of dialogues rather than designing particular systems.

Regarding category (a), in [10], [11] and [12], the importance of a precise relevance notion definition is emphasized. However, these works focus on argumentation-based persuasion dialogues (actually a subset of those, which the author called *disputes*), which belong to the non-collaborative class, and thus *global completeness* is not pursued. Instead, the emphasis is put on properties with similar spirit to our properties of *correctness* and *local progress* (i.e., only the *public knowledge* involved in the dialogue is given importance). In [11] the author considers dynamic disputes in which two participants (proponent and opponent) interchange arguments and counter-arguments, and studies two properties of protocols (namely *soundness* and *fairness*) regarding the relation between the generated public knowledge and the conclusion achieved (in this case, the *winner* of the dispute). The author also gives a natural definition of when a move is relevant: “*iff it changes the status of the initial move of the dispute*” whose spirit is similar to our definition of *natural relevance notion* but taken to the particular case in which the reasoning model is a logic for defeasible argumentation. In [12] the author considers more flexible protocols for disputes, allowing alternative sets of locutions, such as *challenge* and *concede*, and also a more flexible notion of relevance.

Another work in which *relevance* receives an explicit treatment is [9], where the authors investigate the relevance of utterances in an argumentation-based dialogue. However, our *global completeness* property is not pursued, so they do not consider the problematic of distributed contributions (distributed arguments in this case). They study three notions of relevance showing how they can affect the dialogue outcome.

Regarding category (b), in [2] an inquiry dialogue protocol which successfully pursues our idea of *global completeness* is defined. However, the protocol is set upon a particular argumentative system, with the design methodology implicit. They take a simplified version of the *DeLP* system [5], and define an *argument inquiry dialogue* which allows exactly two agents to jointly construct arguments for a given claim. In the present work, we not only explicitly and abstractly analyze the distributed relevance issue, but also consider the complete panorama of collaborative dialogue system behavior, including *correctness* and *progress* properties.

Regarding category (c), different measures for analyzing argumentation-based persuasion are proposed in [1]: measures of the quality of the exchanged arguments, of the

behavior of each agent, and of the quality of the dialogue itself in terms of the relevance and usefulness of its moves. The analysis is done from the point of view of an external agent (*i.e.*, *private knowledge* is not considered), and it is focused in a non-collaborative dialogue type, so they are not concerned with our main problematic.

## 8 Conclusions

We have shown how an existent abstract dialogue framework can be instantiated for modeling argumentation-based dialogues. This new instance, in contrast with a previous one in terms of Propositional Logic Programming [6], naturally deals with possible differences of opinion that can emerge among participants in a dialogue. Also the versatility of the abstract framework is shown through this new instantiation based on a non-monotonic reasoning model.

The obtained framework instance is capable of modeling collaborative argumentation-based dialogues among any number of participants, each of them exposing indifferently either type of argument ('for' and 'against'), and also building arguments together. The model inherits the chance of parametrization and properties from the abstract framework. The most appropriate relevance notion can be chosen according to the dialogue purpose, *e.g.* all the possible justifications for a given claim could be searched, or just one. Also the degree of collaboration is selectable by picking a certain semantics, either basic or full collaborative, according to domain requirements.

In particular, by picking the natural relevance notion and the full collaborative semantics, a model for argumentation-based inquiry has been provided, which ensures a sound and complete (in the usual sense) distributed reasoning. This model is still parametrizable, since different potential relevance notions could be investigated, for instance trying to enhance efficiency. This last issue has been left for future research.

Finally, another item left as future work is the consideration of the case in which the agents can disagree also about evidence. This would imply redefining the reasoning model in order to deal with inconsistencies in the set of evidence.

## References

1. Amgoud, L., de Saint-Cyr, F.D.: On measuring persuasion dialogs quality. In: Besnard, P., Doutre, S., Hunter, A. (eds.) COMMA., vol. 172, pp. 13–24. IOS Press, Amsterdam (2008)
2. Black, E., Hunter, A.: A generative inquiry dialogue system. In: Durfee et al [4], p. 241
3. Dung, P.M.: On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artif. Intell.* 77(2), 321–358 (1995)
4. Durfee, E.H., Yokoo, M., Huhns, M.N., Shehory, O. (eds.): 6th International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS 2007), IFAAMAS, Honolulu, Hawaii, USA, May 14–18 (2007)
5. García, A.J., Simari, G.R.: Defeasible logic programming: An argumentative approach. *Theory and Practice of Logic Programming* 4(1-2), 95–138 (2004)
6. Marcos, M.J., Falappa, M.A., Simari, G.R.: Semantically characterizing collaborative behavior in an abstract dialogue framework. In: Link, S., Prade, H. (eds.) FoIKS 2010. LNCS, vol. 5956, pp. 173–190. Springer, Heidelberg (2010)



7. McBurney, P., Parsons, S.: Dialogue games in multi-agent systems. *Informal Logic. Special Issue on Applications of Argumentation in Computer Science* 22(3), 257–274 (2002)
8. McBurney, P., Parsons, S., Wooldridge, M.: Desiderata for agent argumentation protocols. In: *AAMAS*, pp. 402–409. ACM Press, New York (2002)
9. Parsons, S., McBurney, P., Sklar, E., Wooldridge, M.: On the relevance of utterances in formal inter-agent dialogues. In: Durfee et al [4], p. 240
10. Prakken, H.: On dialogue systems with speech acts, arguments, and counterarguments. In: Brewka, G., Moniz Pereira, L., Ojeda-Aciego, M., de Guzmán, I.P. (eds.) *JELIA 2000. LNCS (LNAI)*, vol. 1919, pp. 224–238. Springer, Heidelberg (2000)
11. Prakken, H.: Relating protocols for dynamic dispute with logics for defeasible argumentation. *Synthese* 127, 187–219 (2001)
12. Prakken, H.: Coherence and flexibility in dialogue games for argumentation. *J. Log. Comput.* 15(6), 1009–1040 (2005)
13. Rotstein, N., Moguillansky, M., García, A., Simari, G.: An abstract argumentation framework for handling dynamics. In: Pagnucco, M., Thielscher, M. (eds.) *NMR*, pp. 131–139 (2008)
14. Walton, D., Krabbe, E.C.W.: *Commitment in Dialogue: Basic Concepts of Interpersonal Reasoning*. State University of New York Press, Albany (1995)

# Argumentation System Allowing Suspend/Resume of an Argumentation Line

Kenichi Okuno\* and Kazuko Takahashi

School of Science&Technology, Kwansei Gakuin University,  
2-1, Gakuen, Sanda, 669-1337, Japan  
o.kenichi@gmail.com, ktaka@kwansei.ac.jp

**Abstract.** This paper discusses an argumentation system that treats argumentation dynamically. We previously proposed a model for dynamic treatment of argumentation in which all lines of argumentation are executed in succession, with the change of the agent's knowledge base. This system was designed for grasping the behaviour of actual argumentation, but it has several limitations. In this paper, we propose an extended system in which these points are revised so that the model can more precisely simulate actual argumentation. In addition, we provide a simpler algorithm for judgement of given argumentation, which can be applied to make a strategy to win.

**Keywords:** computational model for argumentation, belief change, agent communication.

## 1 Introduction

Argumentation is a model that evaluates arguments. It was originally investigated in legal reasoning. Dung's work on constructing a logical framework for argumentation and showing the relationships with nonmonotonic reasoning and logic programming [8] enlarged the possibility of application area of argumentation to the field of artificial intelligence (AI). As a result, formal models of argumentation have received much attention by AI researchers [4,22]. These works include applications for defeasible logic programming [10,6,18,17,15], belief revision [9,19] and so on. Argumentation is considered to be a powerful tool to logically analyse significant phenomena that appear in multiagent systems such as negotiation, agreement and persuasion [14,11], and to make a computational model for a behavior of multiagents [12]. Generally, argumentation proceeds between two agents by giving arguments in turn that attacks the opponent's argument until one of them cannot attack any more. Finally, the loser accepts the winner's proposal. This process is usually represented in the tree form [11,10]. The root node is a proposed formula and each branch corresponds to a single argumentation line, namely, a sequence of arguments. Lots of argumentation systems have proposed so far [4,22], but they considered evaluation of a single argumentation line, and it cannot handle the dynamic properties of actual

---

\* Currently, JSOL Corporation.

argumentation. On the other hand, we have proposed a system that can treat continuous evaluation of multiple argumentation lines [20,21].

Let us consider an example of argumentation. According to many argumentation systems, a proposer P makes the first argument and a defeater C makes counterarguments. We suppose a situation in which a murderer P tells a lie: "I did not commit murder". A policeman C argues that P's statement is a lie.  $P_i$  and  $C_i$  represent P's and C's  $i$ -th utterances, respectively.

- P<sub>1</sub>: "I did not commit murder! There is no evidence!"
- C<sub>1</sub>: "There is evidence. We found your license near the scene."
- P<sub>2</sub>: "It's not evidence! I had my license stolen!"
- C<sub>2</sub>: "It is you who killed the victim. Only you were near the scene at the time of the murder."
- P<sub>3</sub>: "I didn't go there. I was at facility A at that time."
- C<sub>3</sub>: "At facility A? No, that's impossible. Facility A does not allow a person to enter without a license. You said that you had your license stolen, didn't you?"

Figure 1 shows the structure of this argumentation.

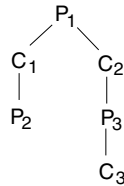


Fig. 1. Structure of argumentation

In this example, if argumentation proceeds along the left branch, and if C has no counterargument to P<sub>2</sub>, then C continues a counterargument in the right branch which attacks P<sub>1</sub> from another side. Finally, C points out the contradiction between P's utterances and wins. P's utterance P<sub>2</sub> gives C new information and causes C to generate C<sub>3</sub>.

To capture the behaviour in this example, we have proposed an argumentation system incorporating changes in an agent's knowledge base caused by the exchange of arguments and defined a new concept of "dynamic win" which is different from the usual concept of "win" obtained by static analysis [20,21]. The goal of this system is to dynamically grasp argumentation by providing a model for actual argumentation. However, several points exist in which this earlier system does not reflect actual argumentation.

The first limitation is on the mechanism that brings up the settled matter in some argumentation line. Consider another argumentation:

- P<sub>1</sub>: "I did not commit murder! There is no evidence!"
- C<sub>2</sub>: "It is you who killed the victim. Only you were near the scene at the time of the murder."

P<sub>3</sub>: “I didn’t go there. I was at facility A at that time.”

C<sub>1</sub>: “There is evidence. We found your license near the scene.”

P<sub>2</sub>: “That’s not evidence! I had my license stolen!”

C’<sub>3</sub>: “That’s strange. Facility A does not allow a person to enter without a license. You said that you were at facility A when the murder occurred. How did you enter?”

In this case, an argumentation first proceeds along the right branch, and then continues to the left branch, P’s utterance P<sub>2</sub> gives C new information and causes C to generate C’<sub>3</sub> as a counterargument to P<sub>3</sub>. C also points out the contradiction between P’s utterances, and wins. Such a phenomenon frequently occurs in actual argumentation when each argumentation line is not so long. In our earlier system, this mechanism could not be handled. In this paper, we present a revised system in which each argumentation line is considered as suspended but may be resumed afterward.

A second shortcoming is the inequivalent rights of agents. In the earlier version, the defeater could continue an argumentation with the revised knowledge base after he/she loses one argumentation line, leading him/her to ultimately win the argumentation tree. However, the proposer loses the whole argumentation tree if he/she loses one argumentation line. In the revised version, we also allow the proposer to continue an argumentation after he/she loses one argumentation line.

The third point is also related to the equivalent rights of agents. In the earlier version, a proposer could not use disclosed information whereas a defeater could. This condition is unfair and unnatural. In the revised version, we adopt *commitment store* [13], a common knowledge base to store all the disclosed information, and both agents can use this knowledge base.

We extend the earlier system by addressing these three points so that it can more precisely simulate an actual argumentation and redefine the dynamic win/lose of an argumentation tree. We show that how this extension improves the treatment of dynamic argumentation.

Moreover, we propose an algorithm for judging the result of an argumentation tree. This algorithm is simpler and easier to implement and it can be applied to formulate an argumentation strategy.

This paper is organised as follows. Section 2 provides the definitions of basic concepts such as argumentation and the argumentation tree. Section 3 proposes an extended model for argumentation incorporating changes in an agent’s knowledge base and also presents an algorithm for the judgement of an argumentation tree. Section 4 provides an example of this algorithm. Section 5 outlines the major changes from our previous work and compares the proposed approach with related works. Finally, section 6 presents conclusions.

## 2 Argumentation

### 2.1 Argumentation Framework

We define an argumentation framework based on Dung [8].

**Definition 1 (consistent).** Let  $\Psi$  be a set of formulas in propositional logic. If there does not exist  $\psi$  that satisfies both  $\psi \in \Psi$  and  $\neg\psi \in \Psi$ ,  $\Psi$  is said to be consistent.

The knowledge base  $\mathbf{K}_a$  for each agent  $a$  is a finite set of propositional formulas. Note that  $\mathbf{K}_a$  is not necessarily consistent and may have no deductive closure; that is, a case may exist in which  $\phi, \phi \rightarrow \psi \in \mathbf{K}_a$  and  $\psi \notin \mathbf{K}_a$  hold. An agent  $a$  participates in argumentation using elements of  $\mathbf{K}_a$ .

**Definition 2 (support).** For a nonempty set of formulas  $\Psi$  and a formula  $\psi$ , if there exist  $\phi, \phi \rightarrow \psi \in \Psi$ , then  $\Psi$  is said to be a support for  $\psi$ .

**Definition 3 (argument).** Let  $\mathbf{K}_a$  be a knowledge base for an agent  $a$ . An argument of  $a$  is a pair  $(\Psi, \psi)$  where  $\Psi$  is a subset of  $\mathbf{K}_a$ , and  $\psi \in \mathbf{K}_a$  such that  $\Psi$  is the empty set or a consistent support for  $\psi$ ,

For an argument  $A = (\Psi, \psi)$ ,  $\Psi$  and  $\psi$  are said to be *grounds* and a *sentence* of  $A$ , respectively. They are denoted by  $Grounds(A)$  and  $Sentence(A)$ , respectively.  $S(A)$  denotes  $Grounds(A) \cup \{Sentence(A)\}$ . If  $\psi \in S(A)$ , it is said that a formula  $\psi$  is contained in an argument  $A$ .

Similar to many argumentation systems, we adopt the concept of *preference* [3,16]. Preferences are assigned to formulas depending on their strength, certainty and stability to avoid loops in the argumentation. Here, we assume that a formula is given a preference value based on some basic rules in advance regardless of the knowledge base in which it is contained, and adopt a simple definition for computing the preference of an argument. Although these definitions affect the result of argumentation, we do not discuss the definitions here, since this aspect of argumentation is out of the scope of this paper.

**Definition 4 (preference).** Each formula is assigned a preference value. Let  $\nu(\psi)$  be the preference for a formula  $\psi$ . Then, the preference of an argument  $A$  is defined by  $\sum_{\psi \in S(A)} \nu(\psi)$ .

**Definition 5 (attack).** Let  $AR_{\mathbf{K}_a}$  and  $AR_{\mathbf{K}_b}$  be sets of all possible arguments of agents  $a$  and  $b$ , respectively.

1. If  $Sentence(A_a) \equiv \neg Sentence(A_b)$  and  $\nu(A_a) \geq \nu(A_b)$ , then  $(A_a, A_b)$  is said to be a rebut from  $a$  to  $b$ .
2. If  $\neg Sentence(A_a) \in Grounds(A_b)$  and  $\nu(A_a) \geq \nu((\emptyset, \neg Sentence(A_a)))$ , then  $(A_a, A_b)$  is said to be an undercut from  $a$  to  $b$ .
3. An attack from  $a$  to  $b$  is either a rebut or an undercut from  $a$  to  $b$ .

When  $(A_a, A_b)$  is an attack from  $a$  to  $b$ , it is said that  $A_a$  attacks  $A_b$ .

Based on Dung [8], in an argumentation framework between two agents, a proposer  $P$  makes the first argument and a defeater  $C$  makes counterarguments. Hereafter,  $\mathbf{K}_P$  and  $\mathbf{K}_C$  denote their knowledge bases, respectively.

**Definition 6 (argumentation framework).** Let  $AR_{\mathbf{K}_P}$  and  $AR_{\mathbf{K}_C}$  be sets of all possible arguments of  $P$  and  $C$ , respectively, with preferences  $\nu$ . Let  $AT_{\mathbf{K}_P \rightarrow \mathbf{K}_C}$  and  $AT_{\mathbf{K}_C \rightarrow \mathbf{K}_P}$  be sets of attacks from  $P$  to  $C$  and from  $C$  to  $P$ , respectively. An argumentation framework between  $P$  and  $C$ ,  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  is defined as a quadruple  $\langle AR_{\mathbf{K}_P}, AR_{\mathbf{K}_C}, AT_{\mathbf{K}_P \rightarrow \mathbf{K}_C}, AT_{\mathbf{K}_C \rightarrow \mathbf{K}_P} \rangle$ .

## 2.2 Argumentation Tree

**Definition 7 (move).** A move is a pair of a player (an agent)  $P/C$  and an argument  $A$  in which  $A \in AR_{\mathbf{K}_P}/AR_{\mathbf{K}_C}$ . If player is  $P/C$ , then it is said to be  $P/C$ 's move. For a move  $M = (\text{player}, \text{argument})$ , we denote player and argument by  $\text{Ply}(M)$  and  $\text{Arg}(M)$ , respectively.

**Definition 8 (move's attack).**  $M$  is said to be an attack to  $M'$ , if  $(\text{Arg}(M), \text{Arg}(M'))$  is an attack from  $\text{Ply}(M)$  to  $\text{Ply}(M')$ .

**Definition 9 (argumentation line, argument set).** Let  $P$  and  $C$  denote a proposer of a formula  $\varphi$  and its defeater, respectively. Let  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  be an argumentation framework between  $P$  and  $C$ . An argumentation line  $D$  for  $\varphi$  on  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  is a finite nonempty sequence of moves  $[M_1, \dots, M_n]$  ( $i = 1, \dots, n$ ) that satisfies the following:

1.  $\text{Ply}(M_1) = P$ , where  $\text{Sentence}(\text{Arg}(M_1)) = \varphi$ .
2. If  $i$  is odd, then  $\text{Ply}(M_i) = P$ , and if  $i$  is even, then  $\text{Ply}(M_i) = C$ .
3.  $M_{i+1}$  is an attack to  $M_i$  for each  $i$  ( $1 \leq i \leq n - 1$ ).
4. No attack occurs against  $\text{Arg}(M_n)$ .
5.  $M_i \neq M_j$  for each pair of  $i, j$  ( $1 \leq i \neq j \leq n$ ).
6. Both  $S(\text{Arg}(M_1)) \cup S(\text{Arg}(M_3)) \cup S(\text{Arg}(M_5)) \cup \dots \cup S(\text{Arg}(M_o))$  and  $S(\text{Arg}(M_2)) \cup S(\text{Arg}(M_4)) \cup S(\text{Arg}(M_6)) \cup \dots \cup S(\text{Arg}(M_e))$  are consistent, where  $o$  and  $e$  are the largest odd number and the largest even number less than or equal to  $n$ , respectively.

The above  $S(\text{Arg}(M_1)) \cup S(\text{Arg}(M_3)) \cup S(\text{Arg}(M_5)) \cup \dots \cup S(\text{Arg}(M_o))$  and  $S(\text{Arg}(M_2)) \cup S(\text{Arg}(M_4)) \cup S(\text{Arg}(M_6)) \cup \dots \cup S(\text{Arg}(M_e))$  are said to be  $P$ 's argument set on  $D$  and  $C$ 's argument set on  $D$ , and they are denoted by  $S_P(D)$  and  $S_C(D)$ , respectively.

This definition puts the constraints of loop-freeness and consistency of each agent's arguments on an argumentation line.

**Definition 10 (win of an argumentation line).** If the last element of an argumentation line  $D$  is  $P$ 's move, then it is said that  $P$  wins  $D$ ; otherwise,  $P$  loses  $D$ .

**Definition 11 (argumentation tree).** An argumentation tree for  $\varphi$  on  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  is a tree in which the root node at depth 0 is empty and all the branches<sup>1</sup> starting from the node of depth 1 are different argumentation lines for  $\varphi$  on  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$ .

<sup>1</sup> Here, a branch is a path from the designated node to a leaf node.

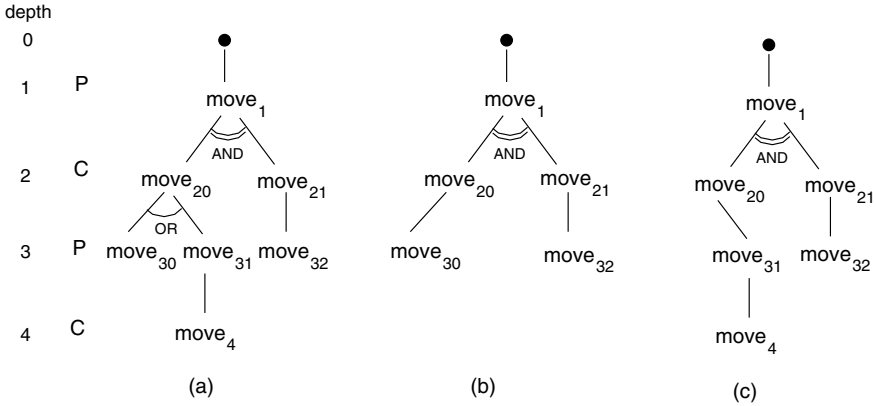


Fig. 2. An argumentation tree and its candidate subtrees

**Definition 12 (candidate subtree).** A candidate subtree is a subtree of an argumentation tree that selects only one child node for each node corresponding to C’s move in the original tree, and selects all child nodes for each node corresponding to P’s move.

**Definition 13 (solution subtree).** A solution subtree is a candidate subtree in which P wins all of the argumentation lines in the tree.

Each candidate subtree corresponds to P’s selection of an argument, and the solution subtree indicates the case in which P takes a winning strategy. In Figure 2, (a) is an argumentation tree, (b) and (c) are its candidate subtrees, and (b) is the solution subtree.

In general, judgement of an argumentation tree is defined as follows.

**Definition 14 (static win of an argumentation tree).** If an argumentation tree has a solution subtree, then P statically wins the argumentation tree; otherwise, P statically loses it.

### 3 Argumentation with Changes in the Knowledge Base

#### 3.1 Execution of Argumentation

We propose a dynamic argumentation system that considers the successive executions of all possible argumentation lines, whilst the usual ones consider only a single argumentation line. In a dynamic argumentation, we have to consider the interaction of argumentation lines.

We introduce the commitment store and a history. The commitment store is a set of all the formulas contained in all arguments given so far. History  $H_a$  is prepared for each agent  $a$  to preserve the coherence of each agent’s arguments.  $H_a$  is a set of all the formulas contained in  $a$ ’s arguments in the argumentation lines in which  $a$  wins. We, however, ignore the coherence of the loser’s side. This

is based on the idea that the winner should be responsible for his/her arguments, but the loser can make an attack from a different side.

A dynamic argumentation line is defined by extending a static argumentation line with history.

**Definition 15 (dynamic argumentation line).** *Let  $P, C$  denote a proposer of a formula  $\varphi$  and its defeater. Let  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  be an argumentation framework between  $P$  and  $C$ . A dynamic argumentation line  $D = [M_1, \dots, M_n]$  for  $\varphi$  on  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  with histories  $\mathbf{H}_P$  and  $\mathbf{H}_C$  is defined as the extension of the (static) argumentation line by adding the following additional condition.*

7.  $\mathbf{H}_{Ply(M_i)} \cup S(\text{Arg}(M_i))$  is consistent for each  $i$  ( $1 \leq i \leq n$ ).

If no misleading is involved, a dynamic argumentation line for  $\varphi$  on  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  with a history  $\mathbf{H}_P$  and  $\mathbf{H}_C$ , is said to be just an argumentation line on  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$ .

Here, we present a dynamic argumentation model. We consider the execution of an argumentation as selecting a branch, updating the commitment store and agents' histories and modifying a tree.

An argumentation starts by selecting a branch of an initial argumentation tree. It proceeds along the branch and when the execution reaches the leaf node, the branch is suspended. At that time, the commitment store is updated and agents can make new arguments using the commitment store in addition to their own knowledge bases. New nodes are added to the argumentation tree if new arguments are generated due to this change of knowledge base. Next, another branch is selected.

On the execution procedure, the executed node is marked and the branch containing unmarked nodes can be selected. The suspended branch may be resumed if a new unmarked node is added to it. On the selection of a branch, the turn of an utterance should be kept. This means that if one branch is suspended at the node that corresponds to one agent's argument, then a next branch should start with the node that corresponds to the other agent's argument.

**Definition 16 (executable node).** *For a node  $M_i$  ( $1 \leq i \leq n$ ) in a branch  $D = [M_1, \dots, M_n]$  and a current turn  $t$ , if  $M_1, \dots, M_{i-1}$  are marked and  $M_i, \dots, M_n$  are unmarked, and  $Ply(M_i) = t$ , then the node  $M_i$  is said to be executable.*

**Definition 17 (execution of a branch).** *For a branch  $D = [M_1, \dots, M_n]$ , histories  $\mathbf{H}_P, \mathbf{H}_C$  and the commitment store  $\mathbf{K}$ , execution of  $D$  from  $i$  ( $1 \leq i \leq n$ ) with  $\mathbf{H}_P$  and  $\mathbf{H}_C$  is defined as follows.*

1. Mark  $M_i, \dots, M_n$ .
2. Set  $\mathbf{K} = \mathbf{K} \cup \bigcup_{k=i}^n S(\text{Arg}(M_k))$ .
3. **if**  $Ply(M_n) = P$ ,  
     **then** set the current turn to  $C$  and  $\mathbf{H}_P = \mathbf{H}_P \cup S_P(D)$ .  
   **if**  $Ply(M_n) = C$ ,  
     **then** set the current turn to  $P$  and  $\mathbf{H}_C = \mathbf{H}_C \cup S_C(D)$ .



**Definition 18 (suspend/resume).** *After the execution of all nodes in a branch,  $D$  is said to be suspended. For a suspended branch  $D$ , if an executable node is added to its leaf on the modification of a tree, and  $D$  is selected, then  $D$  is said to be resumed.*

This Argumentation Procedure with Knowledge Change is formalised as follows.

Argumentation Procedure with Knowledge Change ( $APKC2$ )

Let  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  be an argumentation framework, and  $\varphi$  be a proposed formula.

[STEP 1(initialisation)]

Set  $\mathbf{K} = \emptyset$ ,  $\mathbf{H}_P = \emptyset$ ,  $\mathbf{H}_C = \emptyset$ ,  $turn = P$ . Construct an initial argumentation tree for  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  on  $\varphi$  with  $\mathbf{H}_P, \mathbf{H}_C$  with all the nodes unmarked.

[STEP 2(execution of an argumentation)]

**if** no branch has an executable node,  
     **if**  $turn=P$ , **then** terminate with P's lose.  
     **else**  $turn=C$ , **then** terminate with P's win.  
     **else** select a branch and execute it from the executable node.

[STEP 3(modification of a tree)]

Reconstruct an argumentation tree for  $AF(\mathbf{K}_P \cup \mathbf{K}, \mathbf{K}_C \cup \mathbf{K}, \nu)$  on  $\varphi$  with  $\mathbf{H}_P, \mathbf{H}_C$ .

**if** for any pair of node  $N$  and  $M$  in the tree  
     where  $Ply(N) = Ply(M)$  and  $Arg(N) = Arg(M)$ ,  
      $N$  is marked whilst  $M$  is unmarked,  
     **then** mark  $M$ .  
 go to STEP 2.

The elements of  $\mathbf{K}$  are included either by  $\mathbf{K}_P$  or  $\mathbf{K}_C$ , which are both finite sets. It follows that finite kinds of moves can be generated. Therefore,  $APKC2$  terminates.

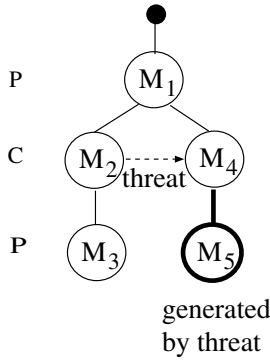
In the modification of a tree in  $APKC2$ , a new node may be added. An idea of *threat* is introduced to explain this situation.

**Definition 19 (threat).** *Let  $M$  and  $M'$  be moves in an argumentation tree  $T$  on  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$ . If  $S(Arg(M))$  generates more than one new move that attacks  $M'$ , then it is said that  $M$  is a threat to  $M'$ , and that  $T$  contains a threat.  $M$  and  $M'$  are said to be a threat resource and a threat destination, respectively.*

Intuitively, a threat is a move that may provide information advantageous to the opponent. A move may be a threat to a move in the same branch.

**Definition 20 (continuous candidate subtree).** *For a candidate subtree  $CT$ , if at least one candidate subtree is generated by the addition of nodes, then these subtrees are said to be continuous candidate subtrees of  $CT$ .*

Note that nondeterminism is involved in the selection of a branch in  $APKC2$ , and finally obtained trees are different depending on the selection.



**Fig. 3.** Argumentation affected on the execution order of branches

Consider an argumentation tree in Figure 3. In this figure,  $M_2$  and  $M_4$  are a threat resource and a threat destination, respectively, and  $M_5$  is a newly generated node by this threat. If we execute from the left branch, then *APKC2* proceeds by executing  $M_1, M_2, M_3, M_4, M_5$ , and terminates with P’s win. On the other hand, if we execute from the right branch, then *APKC2* proceeds by executing  $M_1, M_4$  and suspends. The next turn is P. If there exists no branch in the other candidate trees that starts with P and ends with P, the suspended branch never resumes, and *APKC2* terminates with P’s lose.

We define a dynamic win/lose of an argumentation tree according to *APKC2*.

**Definition 21 (dynamic solution subtree).** *Let  $CT$  be a candidate subtree of an initial argumentation tree. For any execution order of branches of  $CT$ , if *APKC2* terminates with P’s win or  $CT$  has a continuous subtree such that P wins, then  $CT$  is said to be a dynamic solution subtree.*

**Definition 22 (dynamic win of an argumentation tree).** *If an argumentation tree has a dynamic solution subtree, then P dynamically wins the argumentation tree; otherwise, P dynamically loses it.*

### 3.2 Judgement of Dynamic Win/Lose

*APKC2* gives an execution model for an argumentation procedure. If we only want to judge the result of an argumentation and not simulate the procedure, then there exists a simpler algorithm.

**Definition 23 (consistent candidate subtree).** *Let  $CT$  be a candidate subtree. If there does not exist moves  $M, M'$  and a formula  $\psi$  that satisfy  $\text{Ply}(M) = \text{Ply}(M') = P, \psi \in S(\text{Arg}(M))$  and  $\neg\psi \in S(\text{Arg}(M'))$ , then  $CT$  is said to be a consistent candidate subtree.*

Let  $CT$  be a candidate subtree of an argumentation tree of  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$ . Then we can judge a proposer P’s win/lose of  $CT$  by the following algorithm. Hereafter,  $D \in T$  denotes that a branch  $D$  in a tree  $T$ .

Judgement of Win/Lose of a Candidate Subtree (*JC*)

[STEP 1]

if *CT* is not consistent, **then** terminate with P's lose.

**else if** there exists a leaf corresponding C's move in *CT*,

**then** terminate with P's lose.

**else** set  $\mathcal{K} = \bigcup_{D \in CT} S_P(D) \cup \bigcup_{D \in CT} S_C(D)$ .

[STEP 2]

Reconstruct *CT* on  $AF(\mathbf{K}_P \cup \mathcal{K}, \mathbf{K}_C \cup \mathcal{K}, \nu)$ , and let the resultant tree be *CT'*.

[STEP 3]

if  $CT' = CT$ , **then** terminate with P's win.

**else** select a new continuous candidate subtree and go to STEP 1.

The algorithm *JC* terminates by the same reason as that for termination of *APKC2*.

We show the relationship of dynamic win of an argumentation tree and the judgement by *JC*.

First, we show that P dynamically wins an argumentation tree *T* if there exists a candidate subtree in *T* for which *JC* terminates with P's win.

**Theorem 1.** *Let T be an argumentation tree which includes no threat over different candidate subtrees. P dynamically wins T if there exists a candidate subtree in T for which JC terminates with P's win.*

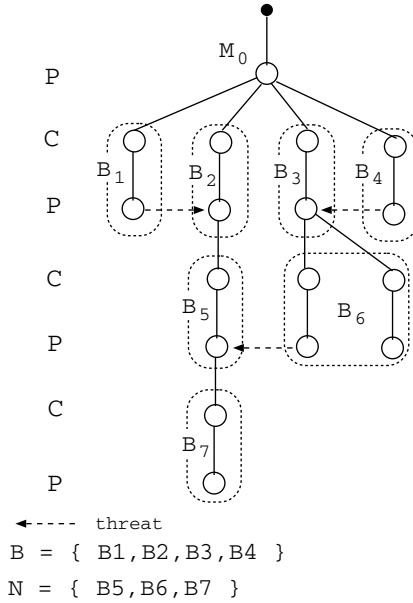
Proof

Let *CT* be a candidate subtree of an argumentation tree of  $AF(\mathbf{K}_P, \mathbf{K}_C, \nu)$  for which *JC* terminates with P's win. We show that for any execution order of branches *APKC2* terminates with P's win.

Let  $M_0$  be a node nearest to the root node of *CT* that corresponds to P's move. Let  $\mathcal{B}$  be a set of sequences each of which consists of nodes except for  $M_0$  in each branch of *CT*. Let  $\mathcal{N}$  be a set of sequences each of which consists of nodes added on STEP2 of *JC*. And let  $Nodes = \mathcal{B} \cup \mathcal{N}$  (Figure 4). Every element of *Nodes* is a sequence of nodes  $M_1 \dots, M_h$  where  $Ply(M_1) = C, Ply(M_h) = P$  and  $M_{i+1}$  attacks  $M_i$  ( $1 \leq i \leq h - 1$ ). Then, any execution order of branches can be represented as a finite sequence of elements of *Nodes* following  $M_0$ . For example,  $M_0 \rightarrow B_1 \rightarrow B_2 \rightarrow B_5$  in Figure 4 is such a sequece. Then, its final node is P's move. Moreover, consistency of  $\mathbf{H}_P$  and  $\mathbf{H}_C$  in *APKC2* are preserved since all reconstructed candidate trees in *JC* are consistent. Therefore, for any execution order of branches *APKC2* terminates with P's win.  $\square$

Next, we show the opposite direction of this theorem, that is, there exists a candidate subtree in *T* for which *JC* terminates with P's win, if P dynamically wins *T*. First, we prove it for a simple case, then for a general case.

**Lemma 1.** *Let T be an argumentation tree which includes no threat over different candidate subtrees. Assume that all the branches selected in APKC2 belong to a single candidate subtree. There exists a candidate subtree in T for which JC terminates with P's win if P dynamically wins T.*



**Fig. 4.** Illustration of the proof for Theorem □

**Proof**

In this case, we show that P dynamically loses  $T$  if  $JC$  terminates with P's lose for all candidate subtrees in  $T$ .

For a candidate subtree  $CT$ , if  $JC$  terminates with P's lose for  $CT$ , there exists a branch whose leaf node is C's move in some step of  $JC$  procedure.

Assume that there exists such a branch at an initial state.  $APKC2$  terminates with P's lose immediately if this branch is selected, since there is no branch beginning from P's move and there is no executable node.

Let  $CT_0, \dots, CT_k$  be a sequence of reconstructed candidate subtrees of  $CT$  in  $JC$  procedure. Assume that all the leaves are P's moves in  $CT_i$  ( $0 \leq i \leq k - 1$ ) and there exists a branch  $D$  whose leaf node is C's move in  $CT_{i+1}$ . The leaf node  $M$  of  $D$  in  $CT_i$  is P's move (Figure 5). There exists threat resource in the nodes in  $CT_i$  whose threat destination is  $M$ . It follows that new nodes are added to  $D$  in  $CT_{i+1}$ . Let  $D'$  be the branch that contains the threat resource. (Note that  $D'$  may be  $D$ .) If  $D'$  and  $D$  are executed in this order in  $APKC2$ ,  $APKC2$  terminates with C's move after executing all the nodes including new nodes. Then,  $APKC2$  terminates with P's lose since there is no branch beginning from P's move and there is no executable node.

Therefore, P dynamically loses  $T$ . □

In lemma □, we assume that all the branches selected in  $APKC2$  belong to a single candidate subtree. However, branches in multiple candidate subtrees may be selected, since  $APKC2$  allows selection of any branch. Generally, there should

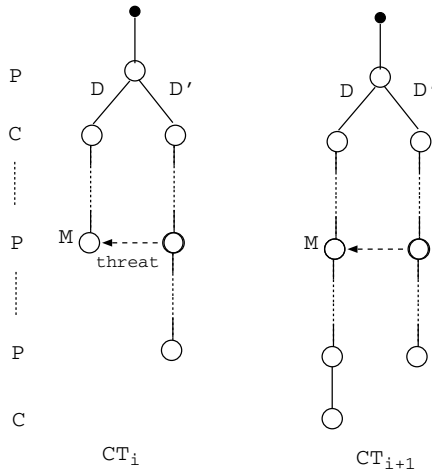


Fig. 5. Illustration of the proof for Lemma 1

exist a (static) solution subtree that is included by a set of all executed branches. We show this by the following two lemmas.

**Lemma 2.** *Let  $T_f$  be a finally obtained tree when APKC2 terminates with  $P$ 's win. For a subtree  $T$  whose root node  $M$  is  $C$ 's move in  $T_f$ , let  $M_{P_1}, \dots, M_{P_n}$  be  $M$ 's child nodes, and let  $T_1, \dots, T_n$  be subtrees whose root nodes are  $M_{P_1}, \dots, M_{P_n}$ , respectively. If  $T_1, \dots, T_n$  are all candidate subtrees, then there exists a (static) solution subtree  $T_i$  ( $1 \leq i \leq n$ ).*

Assume that  $T_i$  is not a solution subtree for some  $i$ . Then,  $T_i$  includes  $C$ 's move as a leaf node. Let  $D$  be a branch of  $T_f$  that contains this node. There should exist another branch as the successive execution of  $D$ , since  $APKC2$  terminates with  $P$ 's move. On the other hand, when the leaf node of  $D$  has been executed, the unmarked nodes nearest to the root node of  $T_f$  in every branch of  $T_i$  that includes unmarked nodes are  $C$ 's moves. They are unexecutable. Therefore, a branch in subtrees other than  $T_i$  should be selected as  $D$ 's successive execution. If none of  $T_1, \dots, T_n$  is a solution subtree, it is impossible to terminate  $APKC2$  with  $P$ 's move. Hence, one of them should be a solution subtree.  $\square$

We can take such  $T_i$  as  $T$ 's candidate subtree, and obtain the following lemma.

**Lemma 3.** *Let  $T_f$  be a finally obtained tree when APKC2 terminates with  $P$ 's win.  $T_f$  includes a (static) solution subtree.*

Proof

For a subtree whose root node  $M$  is  $C$ 's move in  $T_f$ , let  $M_{P_1}, \dots, M_{P_n}$  be  $M$ 's child nodes, and let  $T_1, \dots, T_n$  be subtrees whose root nodes are  $M_{P_1}, \dots, M_{P_n}$ , respectively. For each  $i$  ( $1 \leq i \leq n$ ), if  $T_i$  is not a candidate subtree, then replace it by its candidate subtree  $T'_i$  from lemma 2; otherwise, set  $T'_i$  be  $T_i$ . There

exists a solution subtree  $T'_i$  ( $1 \leq i \leq n$ ), since all of them are candidate subtrees. Repeating this procedure, it is proved that  $T_f$  includes a solution subtree.  $\square$

**Theorem 2.** *Let  $T$  be an argumentation tree which includes no threat over different candidate subtrees. There exists a candidate subtree in  $T$  for which  $JC$  terminates with  $P$ 's win if  $P$  dynamically wins  $T$ .*

Proof

Let  $CT'$  is a finally obtained tree for a candidate subtree  $CT$  in  $JC$ . From lemma 3, the finally obtained tree  $T_f$  in  $APKC2$  includes a (static) solution subtree. There exists  $CT$  that contains this solution subtree, since both threat resource and threat desitiation are in the same candidate subtree from the condition. Moreover,  $\bigcup_{D \in T_f} S_P(D)$  is consistent because of the constraints on  $\mathbf{H}_P$ . Therefore, there exists a candidate subtree for which  $JC$  terminates with  $P$ 's win.  $\square$

## 4 An Example

Consider the example shown in Section 1. We illustrate various properties of  $APKC2$  and  $JC$  using this example.

### 4.1 Formalisation

The knowledge bases of a proposer  $P$  and a defeater  $C$  are shown below. The number attached to each formula shows its preference. We assume that the facts and rules are all represented in the knowledge base and the agents have no other knowledge.

$$\mathbf{K}_P = \left\{ \begin{array}{l} \neg m[1], \neg e[2], (\neg e \rightarrow \neg m)[1], \neg(la \rightarrow e)[1], \\ ls[1], (ls \rightarrow \neg(ls \rightarrow e))[1], \neg n[1], a[2], \\ (a \rightarrow \neg n)[1] \end{array} \right\}$$

$$\mathbf{K}_C = \left\{ \begin{array}{l} e[1], la[1], (la \rightarrow e)[2], m[2], n[2], \\ (n \rightarrow m)[1], \neg a[1], (ls \rightarrow \neg a)[1] \end{array} \right\}$$

The propositions have the following meanings:

- $m$ :  $P$  commits murder.
- $e$ : there is evidence.
- $la$ :  $P$ 's license was left at the scene of the murder.
- $ls$ :  $P$ 's license was stolen.
- $n$ :  $P$  was near the scene when the murder was committed.
- $a$ :  $P$  was at facility A when the murder was committed.

### 4.2 The Case of Changing from Static Win to Dynamic Lose

Figure 6 shows a relevant part of an initial argumentation tree and a final argumentation tree in  $APKC2$ . The argumentation starts with the murderer's utterance. The nodes  $M_1, M_2, M_3, M_4, M_5$ , and  $M_6$  correspond to the utterances  $P_1, C_1, P_2, C_2, P_3$ , and  $C_3$ , respectively.

This example shows the case in which a proposer statically wins but dynamically loses the argumentation tree.

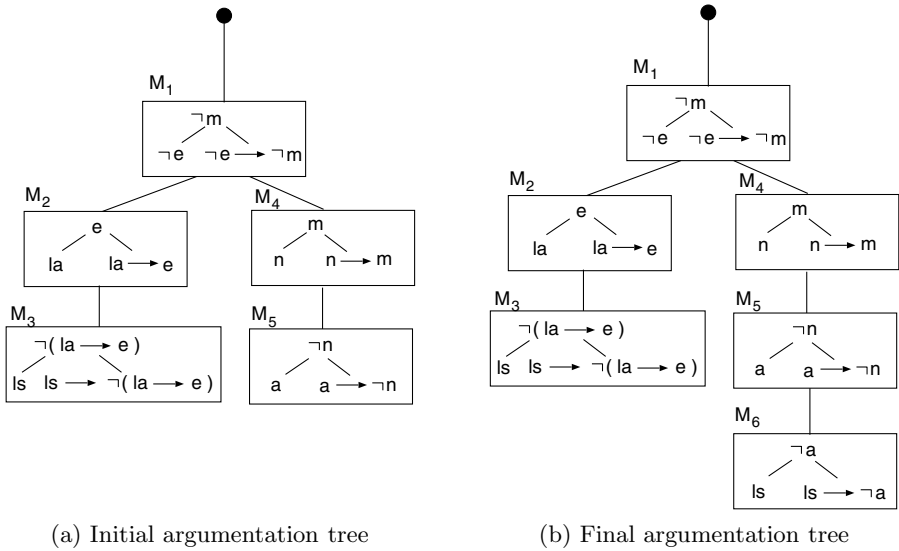


Fig. 6. The argumentation trees starting from the murderer

### 4.3 Behaviour of Suspend/Resume

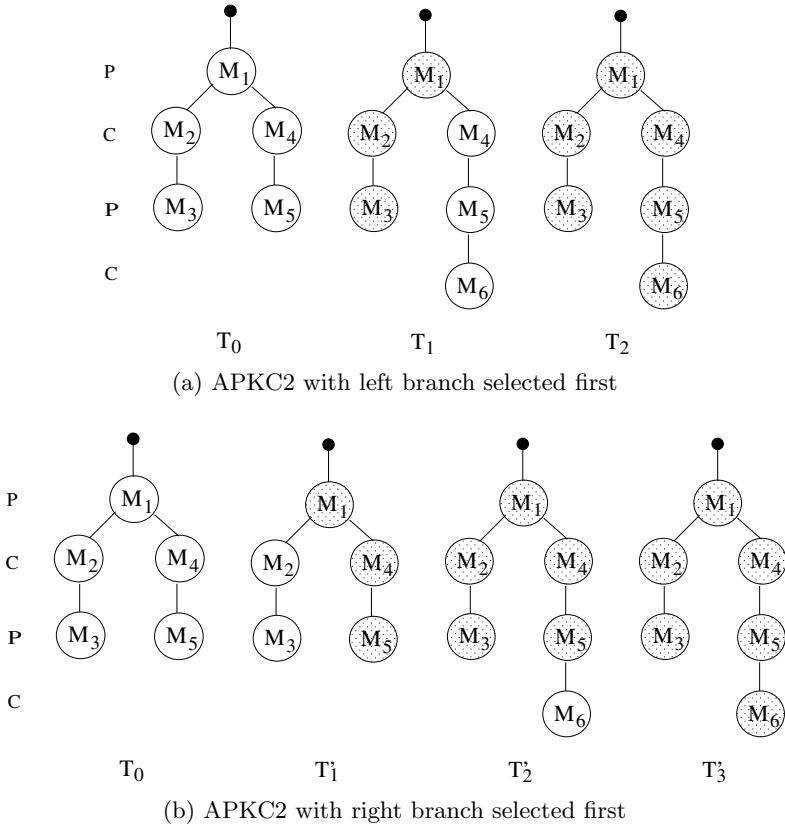
Figure 7(a) shows the trees at each step of *APKC2* procedure in case the left branch of the tree in Figure 6(a) is selected first.  $T_0$  is the initial argumentation tree, and  $T_1$  is the modified tree based on the knowledge bases obtained after the execution of the left branch. In these trees, the hatched nodes are marked.  $T_2$  is the tree modified based on the knowledge bases after the execution of the right branch afterward. No more attacks to the leaf nodes exist. No unmarked node exists in  $T_2$ , which indicates the absence of a counterargument. Then, the procedure terminates. The winner is C, who gives the final argument.

Figure 7(b) shows the trees at each step in case the right branch of the tree in Figure 6(a) is selected first.  $T'_1$  is the modified tree based on the knowledge bases obtained after the execution of the right branch. The right branch is suspended.  $T'_2$  is the modified tree based on the knowledge bases obtained after the execution of the left branch afterward. In this case, a new node  $M_6$ , which corresponds to the utterance  $C'_3$ , is added, and it is the only node that is unmarked in  $T'_2$ . To execute this node, the right branch is resumed.  $T'_3$  is the modified tree based on the knowledge bases obtained after this execution. No unmarked node exists in  $T'_3$ . Then, the procedure terminates. The winner is C, who gives the final argument.

This example shows the procedures with different branch selection orders, and illustrates how suspend/resume occurs.

### 4.4 The Case of Changing from Static Lose to Dynamic Win

Next, we show an argumentation that starts with the policeman's utterance  $C_0$  in the first example:



**Fig. 7.** Comparison of procedures on the order of selecting branches

C<sub>0</sub>: “You committed the murder.”

and continues to P<sub>1</sub>, C<sub>2</sub>, P<sub>3</sub>, C<sub>1</sub>, P<sub>2</sub>, similar to the first example. The argumentation trees are shown in Figure 8. M<sub>0</sub> is a node corresponding to C<sub>0</sub>.

The trees can be regarded as C’s argumentation trees because the roles of P and C are switched from the first example. C statically loses, since all the leaf nodes in the initial argumentation tree shown in Figure 8(a) are P’s move, but dynamically wins, since the final argumentation tree shown in Figure 8(b) is obtained by *APKC2*.

This example shows the case in which a proposer statically loses but dynamically wins the argumentation tree.

#### 4.5 Judgement of Dynamic Win/Lose

Here, we apply an algorithm *JC* to the first example, starting from the murderer’s utterance.



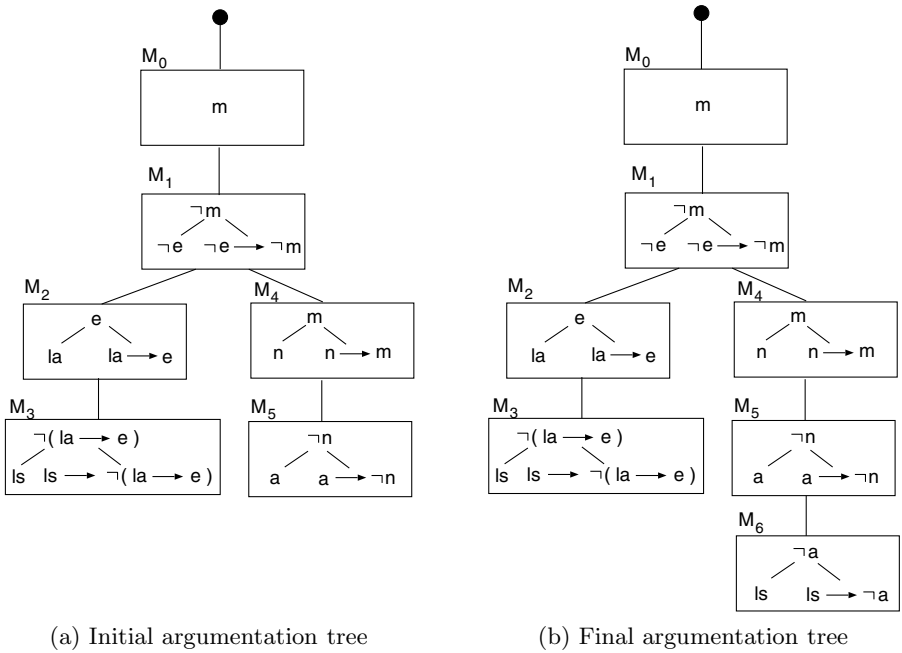


Fig. 8. The argumentation trees starting from the policeman

The initial argumentation tree is shown in Figure 6(a). It includes only one candidate subtree<sup>2</sup> and no threats over different candidate subtrees.

Figure 9 shows how *JC* works.

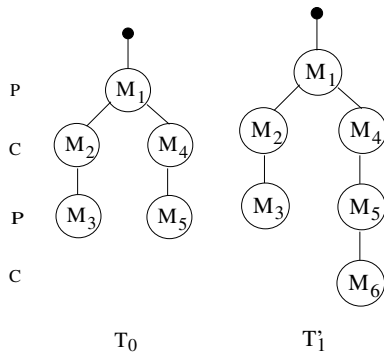


Fig. 9. The argumentation trees for judgement

<sup>2</sup> This figure shows only the relevant part, and it actually contains more candidate subtrees. Although we ignore them to make a description simple, the result is the same.

We take  $T_0$  as a candidate subtree, which is consistent.  
 First, we obtain  $\mathcal{K}$ , a set of all formulas in  $T_0$ .

$$\mathcal{K} = \left\{ \begin{array}{l} \neg m, \neg e, (\neg e \rightarrow \neg m), \neg(la \rightarrow e), \\ ls, (ls \rightarrow \neg(ls \rightarrow e)), \neg n, a, \\ (a \rightarrow \neg n) \\ e, la, (la \rightarrow e), m, n, \\ (n \rightarrow m) \end{array} \right\}$$

Reconstruct the tree, then a new node  $M_6$  is added to obtain the tree  $T'_1$ , which is consistent. Since one leaf node in  $T'_1$  is C's move, P loses this candidate subtree.

Since no other candidate subtrees exist, P dynamically loses the argumentation tree.

## 5 Discussion

### 5.1 Improvements on the Earlier Version

Three significant points distinguish the argumentation system proposed in this paper from the earlier version.

First, suspend/resume of a branch is enabled, allowing for the resumption of a settled matter. We mark the executed node instead of deleting it, and make it possible to add a new node to already executed ones. We also provide a simpler judgement algorithm of win/lose for a given candidate subtree. The method of selecting a candidate tree to win can contribute to argumentation strategy.

Second, both agents can continue an argumentation after he/she loses one argumentation line, whilst only the defeater could do so in the earlier version. This makes it possible to handle the case in which a proposer statically loses but dynamically wins.

Third, both P and C can use disclosed knowledge, whereas only C could do so in the earlier version. We prepare the commitment store for this purpose.

Due to these improvements, *APKC2* provides a more natural model for actual argumentations.

In addition, in the earlier version, we had to reconstruct an argumentation tree every time a branch was executed, since some formulas might be deleted from C's knowledge base. However, in the revised version, we do not need to reconstruct a tree, only add nodes to the existing tree, since the usable knowledge is monotonically increasing. This makes the implementation of *APKC2* easier.

### 5.2 Related Works

García et al. formalized argumentation based on Defeasible Logic Programming (DeLP) [11]. In DeLP, agent's knowledge base consists of two kinds of rules: strict rules and defeasible rules. The result of argumentation is different depending on which defeasible rules are used. Afterwards, Moguillansky discussed revision of

the knowledge base [17]. In his method, after constructing the initial argumentation tree called dialectical tree, knowledge base is changed by extracting defeasible rules and the tree is altered. The goal is to construct undefeated argumentation by selecting suitable defeasible rules. They presented an algorithm for this alteration of the tree and considered a strategy to get the undefeated argumentation. In a series of studies, they formalized several properties in argumentation based on this approach [15]. Similar to our approach, they consider multiple argumentation lines altogether. The different point is that they investigate the effect of the change of knowledge base not considering the change caused by the execution of argumentation, while we focus on the effect of the execution.

Argumentation-based approach is applied to formalize processes appeared in agents communication such as negotiation, persuasion, agreement and so on [14][19]. Considering the effect of the execution of arguments, agents communication are rather related issue, since belief of each agent is updated on receiving information from the other agent. Amgoud proposed the protocol that handles arguments and formalized the case in accepting/rejecting new information [1]. She also presented a general framework for argumentation-based negotiation in which agent has a theory and it evolves during a dialogue [2]. She considered the knowledge base for each agent separately, as well as its revision by exchanging arguments. The significant difference between her work and ours is that in her approach only a single argumentation line is considered, so only threats to the same branch are taken into account, whereas in our approach all argumentation lines are considered successively, so threats to the other branches are examined. Dunne proposed a “dispute tree” on which successive execution of all argumentation lines are considered [7]. However, the revision of agents’ knowledge base, allowing executed moves to add new information to the opponent’s knowledge base, is not considered.

Cayrol studied how acceptable arguments are changed when a new argument is added to an argumentation system based on Dung’s framework [5]. The aim of her research is a formal analysis on changes to argumentation, and the contents of the additional arguments and reasons for the addition are beyond its scope. In contrast, we focus specifically on the effect of knowledge gained by executing argumentation.

## 6 Conclusion

We have proposed an argumentation system *APKC2*, which is an extension of our earlier argumentation system *APKC*. *APKC* is a system in which multiple argumentation lines are executed in succession, and an agent’s knowledge base can change during argumentation. We have extended *APKC* so that the suspend/resume of an argumentation line can be processed, both agents can continue an argumentation after he/she loses one argumentation line and both can use information given in previous arguments. These extensions provide a more natural model of actual argumentation. In addition, we proposed a simpler

algorithm for the judgement of the win/lose result of an argumentation tree, and showed that its result is equivalent to that of *APKC2*.

In future, we are considering an extension of *APKC2* that can not only directly use new information, but also derive new facts from the new knowledge. We are also considering a strategy to win an argumentation.

## References

1. Amgoud, L., Parsons, S., Maudet, N.: Arguments, dialogue, and negotiation. In: ECAI 2000, pp. 338–342 (2000)
2. Amgoud, L., Dimopoulos, Y., Moraitis, P.: A general framework for argumentation-based negotiation. In: Rahwan, I., Parsons, S., Reed, C. (eds.) *Argumentation in Multi-Agent Systems*. LNCS (LNAI), vol. 4946, pp. 1–17. Springer, Heidelberg (2008)
3. Amgoud, L., Vesic, S.: Repairing preference-based argumentation frameworks. In: IJCAI 2009, pp. 665–670 (2009)
4. Bench-Capon, T., Dunne, P.: Argumentation in artificial intelligence. *Artificial Intelligence* 171, 619–641 (2007)
5. Cayrol, C., de St-Cyr, F.D., Lagasque-Shiex, M.-C.: Revision of an argumentation system. In: KR 2008, pp. 124–134 (2008)
6. Chesñevar, C.I., Maguitman, A., Loui, R.: Logical models of argument. *ACM Computing Surveys* 32(4), 337–383 (2005)
7. Dunne, P.E., Bench-Capon, T.J.M.: Two party immediate response disputes: properties and efficiency. *Artificial Intelligence* 149(2), 221–250 (2003)
8. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artificial Intelligence* 77, 321–357 (1995)
9. Falappa, M., Kern-Isberner, G., Simari, G.R.: Explanations, belief revision and defeasible reasoning. *Artificial Intelligence* 141(1-2), 1–28 (2002)
10. García, A., Simari, G.: Defeasible logic programming: an argumentative approach. *Theory and practice of logic programming* 4(1), 95–138 (2004)
11. García, A., Chesnevar, C., Rotstein, N., Simari, G.: An abstract presentation of dialectical explanations in defeasible argumentation. In: *ArgNMR 2007*, pp. 17–32 (2007)
12. Joseph, S., Prakken, H.: Coherence-driven argumentation to norm consensus. In: *ICAIL 2009*, pp. 58–67 (2009)
13. Hamblin, C.: *Fallacies*. Methuen (1970)
14. Kraus, S., Sycara, K., Evenchik, A.: Reaching agreements through argumentation: a logical model and implementation. *Artificial Intelligence* 104(1-2), 1–69 (1998)
15. Lucero, M.J.G., Chesnevar, C.I., Simari, G.R.: On the accrual of arguments in defeasible logic programming. In: IJCAI 2009, pp. 804–809 (2009)
16. Modgil, S.: Reasoning about preferences in argumentation frameworks. *Artificial Intelligence* 173(9-10), 901–1040 (2009)
17. Moguillansky, M.O., et al.: Argument theory change applied to defeasible logic programming. In: *AAAI 2008*, pp. 132–137 (2008)
18. Prakken, H.: Combining skeptical epistemic reasoning with credulous practical reasoning. In: *COMMA 2006*, pp. 311–322 (2006)

19. Paglieri, F., Castelfranchi, C.: Revising beliefs through arguments: Bridging the gap between argumentation and belief revision in MAS. In: Rahwan, I., Moraïtis, P., Reed, C. (eds.) ArgMAS 2004. LNCS (LNAI), vol. 3366, pp. 78–94. Springer, Heidelberg (2005)
20. Okuno, K., Takahashi, K.: Argumentation with a revision of knowledge base. In: EUMAS 2008, CD-ROM (2008)
21. Okuno, K., Takahashi, K.: Argumentation system with changes of an agent's knowledge base. In: IJCAI 2009, pp. 226–232 (2009)
22. Rahwan, I., Simari, G. (eds.): Argumentation in Artificial Intelligence. Springer, Heidelberg (2009)

# Computing Argumentation in Polynomial Number of BDD Operations: A Preliminary Report

Yuqing Tang<sup>1</sup>, Timothy J. Norman<sup>2</sup>, and Simon Parsons<sup>1,3</sup>

<sup>1</sup> Dept. of Computer Science, Graduate Center, City University of New York,  
365 Fifth Avenue, New York, NY 10016, USA  
ytang@gc.cuny.edu

<sup>2</sup> Dept of Computing Science, The University of Aberdeen,  
Aberdeen, AB24 3UE, UK  
t.j.norman@abdn.ac.uk

<sup>3</sup> Dept of Computer & Information Science, Brooklyn College, City University of New York,  
2900 Bedford Avenue, Brooklyn, NY 11210 USA  
parsons@sci.brooklyn.cuny.edu

**Abstract.** Many advances in argumentation theory have been made, but the exponential complexity of argumentation-based reasoning has made it impractical to apply argumentation theory. In this paper, we propose a binary decision diagram (BDD) approach to argumentation-based reasoning. In the approach, sets of arguments and defeats are encoded into BDDs so that an argumentation process can work on a set of arguments and defeats simultaneously in one BDD operation. As a result, the argumentation can be computed in polynomial number of BDD operations on the number of input sentences.

## 1 Introduction

Argumentation provides an elegant approach to nonmonotonic reasoning [15] and decision making [17,26], and now sees wide use as a mechanism for supporting dialogue in multiagent systems [32,33]. As an approach that has its roots in logic — in many systems of argument, the arguments are constructed using some form of logical inference — the efficiency of reasoning using argumentation is a topic of considerable interest [13,16,25] with a number of negative results that stress the fact that generating arguments and establishing properties of arguments can be very costly in computational terms.

In this paper we take a rather different look at the computation of arguments. We have been investigating the creation of multiagent plans [36,37,38], especially the construction of plans that take into account the communication between agents [34,35]. In doing so, we have been using a representation, that of quantified boolean formulae (QBFs) and binary decision diagrams (BDDs), which has been widely adopted in symbolic planning in non-deterministic domains. It turns out that this representation provides a way to compute arguments, and given the computational efficiency of planning based on QBFs and BDDs, it seems that it can provide an efficient way to compute arguments. In this paper, we investigate exactly how efficient this approach is, and conclude that we can carry out many of the basic operations needed to compute arguments in a polynomial number of operations.

Note that we are not claiming to be performing general logical inference in polynomial time. As we explain in detail later in the paper, the “polynomial number of operations” are operations on the BDD representation, and while this representation in many cases can be constructed compactly from a set of logical formulae, there are some cases in which the size of this representation is exponential in the number of formulae.

## 2 Background

This section gives the technical background needed by the remainder of the paper, a description of *quantified boolean formulae*, and *binary decision diagrams*.

### 2.1 Quantified Boolean Formulae

A propositional language  $\mathcal{L}$  based on a set of proposition symbols  $\mathcal{P}$  with quantification can be defined by allowing standard connectives  $\wedge, \vee, \rightarrow, \neg$  and quantifiers  $\exists, \forall$  over the proposition variables. The resulting language is a logic of quantified boolean formulae (QBF) [5]. A *symbol renaming operation*, which we use below, can be defined on  $\mathcal{L}$ , denoted by  $\mathcal{L}[\mathcal{P}/\mathcal{P}']$ , which means that a new language is obtained by substituting the symbols of  $\mathcal{P}$  with the symbols of  $\mathcal{P}'$  where  $\mathcal{P}'$  contains the same set of propositions as that of  $\mathcal{P}$  but using different symbol names (notice that  $|\mathcal{P}'| = |\mathcal{P}|$ ). Similarly, for a formula  $\xi \in \mathcal{L}$ , if  $\mathbf{x}$  is a vector of propositional variables for  $\mathcal{P}$ , then a variable renaming operation can be defined by  $\xi[\mathbf{x}/\mathbf{x}']$  which means that all the appearances of variables  $\mathbf{x} = x_1x_2 \dots x_n$  are substituted by  $\mathbf{x}' = x'_1x'_2 \dots x'_n$  which is a vector of the corresponding variables or constants in  $\mathcal{P}'$ . In a QBF, propositional variables can be universally and existentially quantified: if  $\phi[\mathbf{x}]$  is a QBF formula with propositional variable vector  $\mathbf{x}$  and  $x_i$  is one of its variables, the existential quantification of  $x_i$  in  $\phi$  is defined as  $\exists x_i \phi[\mathbf{x}] = \phi[\mathbf{x}][x_i/FALSE] \vee \phi[\mathbf{x}][x_i/TRUE]$  and the universal quantification of  $x_i$  in  $\phi$  is defined as  $\forall x_i \phi[\mathbf{x}] = \phi[\mathbf{x}][x_i/FALSE] \wedge \phi[\mathbf{x}][x_i/TRUE]$ . Here *FALSE* and *TRUE* are two propositional constants representing “true” and “false” in the logic. Quantifications over a set  $X = \{x_1, x_2, \dots, x_n\}$  of variables is defined as sequential quantifications over each variables  $x_i$  in the set:

$$Q_X \xi = Q_{x_n} Q_{x_{n-1}} \dots Q_{x_1} \xi$$

where  $Q$  is either  $\exists$  or  $\forall$ . The introduction of quantification doesn’t increase the expressive power of propositional logic but allows us to write concise expressions whose quantification-free versions have exponential sizes [11].

With the above language, we can encode sets and relations to manipulate sets of arguments and defeats. Let  $x$  be an element of a set  $X = 2^{\mathcal{P}}$ ,  $x$  can then be explicitly encoded by a conjunction composed of all proposition symbols in  $\mathcal{P}$  in either positive or negative form

$$\xi(x) = \bigwedge_{p_i \in x \text{ and } p_i \in \mathcal{P}} p_i \wedge \bigwedge_{p_j \notin x \text{ and } p_j \in \mathcal{P}} \neg p_j$$

where  $p_i \in x$  means that the corresponding bit  $p_i$  is set to be *TRUE* in the encoding of  $x$ , and  $p_j \notin x$  means that the corresponding bit  $p_j$  is set to be *FALSE* in the encoding

**Table 1.** The mapping between set operators and QBF operators

Set operator	QBF operator
$X_1 \cap X_2$	$\xi(X_1) \wedge \xi(X_2)$
$X_1 \cup X_2$	$\xi(X_1) \vee \xi(X_2)$
$X_1 \setminus X_2$	$\xi(X_1) \wedge \neg \xi(X_2)$
$x \in X$	$\xi(x) \rightarrow \xi(X)$
$X_1 \subseteq X_2$	$\xi(X_1) \rightarrow \xi(X_2)$

of  $x$ . We denote that a formula  $\gamma$  can be satisfied in an element  $x$  by  $x \models \gamma$ . Then a set of elements can be characterized by a formula  $\gamma \in \mathcal{L}$ , with the set denoted by  $X(\gamma)$ , where  $X(\gamma) = \{x | x \models \gamma\}$ . Two special sets, the empty set  $\emptyset$  and the universal set  $\mathcal{U}$ , are represented by *FALSE* and *TRUE* respectively.

With these notions we can have a mapping between the set operations on states and the boolean operations on formulae as shown in Table 1 when  $X_1$  and  $X_2$  are interpreted as two sets of states.

## 2.2 Binary Decision Diagrams

In the above, we have showed the natural connections between the set paradigm and its implicit representation using QBF formulae. Now we will briefly discuss how QBF formulae and the operations over them can be represented and efficiently computed using a data structure called Binary Decision Diagrams (BDD) [5]. In this way, the time and space complexity for exploring the space of arguments and defeats for acceptable arguments can be significantly reduced due to the compact representation provided by BDDs in comparison to explicit search techniques.

A BDD is a rooted directed acyclic graph. The terminal nodes are either *TRUE* or *FALSE*. Each non-terminal node is associated with a boolean variable  $x_i$ , and two BDDs, called *left* and *right*, corresponding to the values of the sub-formula when  $x_i$  is assign *FALSE* and *TRUE* respectively. The value of a QBF formula can be determined by traversing the graph from the root to the leaves following the boolean assignment given to the variables of the QBF formula. The advantage of using BDDs to represent QBF formulae is that most basic operations on QBFs can be performed in linear or quadratic time in terms of the number of nodes used in a BDD representation of the formulae if a special form of BDD, called Reduced Ordered Binary Decision Diagram (ROBDD) [5], is used. A ROBDD is a compact BDD which uses a fixed ordering over the variables from the root to the leaves in the BDD, merges duplicate subgraphs into one, and directs all their incoming edges into the merged subgraph. Following the notation traditionally used in symbolic model checking and AI planning, we will refer to an ROBDD simply as a BDD.

Let  $\xi, \xi_1, \xi_2$  be QBF formulae, let the number of nodes used in its BDD representation denoted by  $\|\cdot\|$ . With this BDD representation, the complexity of a QBF binary operator  $\langle op \rangle$  (e.g.  $\wedge, \vee, \rightarrow$ ) on two formulae  $\xi_1$  and  $\xi_2$ , namely  $\xi_1 \langle op \rangle \xi_2$ , is  $O(\|\xi_1\| \times \|\xi_2\|)$ , that of negation  $\neg \xi$  is  $O(\|\xi\|)$  (or  $O(1)$  if complement edges are introduced to the BDDs), and that of quantification  $Q_{x_i}(f[x])$ , where  $Q$  is either  $\exists$  or  $\forall$ , is  $O(\|f\|^2)$  [5][11] as summarized in Table 2.



**Table 2.** The mapping between QBF operators and BDD operators.  $\xi, \xi_1, \xi_2$  are formulae in QBF;  $G(\xi), G(\xi_1), G(\xi_2)$  are BDD representations for these formulae.

QBF/Set operator	BDD operator	Complexity
$\neg\xi$	$\neg G(\xi)$	$O(\ \xi\ )$
$\exists x_i(\xi)$	$G(\xi_{x_i=0}) \vee G(\xi_{x_i=1})$	$O(\ \xi\ ^2)$
$\forall x_i(\xi)$	$G(\xi_{x_i=0}) \wedge G(\xi_{x_i=1})$	$O(\ \xi\ ^2)$
$\xi_1 \wedge \xi_2$	$G(\xi_1) \wedge G(\xi_2)$	$O(\ \xi_1\  \cdot \ \xi_2\ )$
$\xi_1 \vee \xi_2$	$G(\xi_1) \vee G(\xi_2)$	$O(\ \xi_1\  \cdot \ \xi_2\ )$
$\xi_1 \rightarrow \xi_2$	$G(\xi_1) \rightarrow G(\xi_2)$	$O(\ \xi_1\  \cdot \ \xi_2\ )$
$ X $	$Sat-count(G(\xi(X)))$	$O(\ \xi(X)\ )$

The key advantage of using BDDs (and QBFs, which function as a front end language for BDDs) to represent sets and relations is that the complexity of the operations on those sets and relations will depend on the complexity of the BDD representation instead of the size of the sets and relations. Since the complexity of the BDD representation of doesn't necessarily depend on the size of those sets and relations either, it is possible to carry out operations in a time that is not directly a function of the size of the sets and relations. In fact, the operations on BDDs are polynomial in the size of the BDD, and so using BDDs we can compute operations on sets and relations in time polynomial in the size of their BDD representation. This can be a considerable improvement over a more direct implementation of the operations which is, of course, exponential in the size of the sets and relations.

### 3 Set-Theoretic Argumentation

Having introduced the ideas from QBFs and BDDs, in this section we give an overview of the argumentation system we will capture using them. The framework we use is mostly drawn from the work of Amgoud and her colleagues [11, 2] with some slight modifications. This framework will abstract away the inference procedure by which the arguments are created and only keep track of the premises the arguments are based on. In the next section, we will introduce the inference procedure back into the representation of arguments.

**Definition 1.** An argument based on  $\Sigma \subseteq \mathcal{L}$  is pair  $(H, h)$  where  $H \subseteq \Sigma$  and  $H \neq \emptyset$  such that

1.  $H$  is consistent with respect to  $\mathcal{L}$ ,
2.  $H \vdash h$ ,
3.  $H$  is minimal in the sense of set inclusion.

$H$  is called the support and  $h$  is called the conclusion of the argument.  $\mathcal{A}(\Sigma)$  denotes the set of all arguments which can be constructed from  $\Sigma$ .

<sup>1</sup> We will return to the relationship between the size of the sets and relations and the complexity of the BDD representation below.

This definition of an argument can be understood as a set of constraints on how information can be clustered as arguments. Condition (1) ensures that an argument is coherent. The coherence of an agent's information is defined in terms of the consistency of the language  $\mathcal{L}$  in which the information is written. Condition (2) can be understood as insisting that the conclusion of an argument should be supported by a set of information in the sense of inference in the language  $\mathcal{L}$ . Condition (3) can be understood as saying that no redundant information should appear in an argument.

**Definition 2.**  $(H', h')$  is a subargument of the argument  $(H, h)$  iff  $H' \subseteq H$ .

**Definition 3.** Let  $(H_1, h_1)$ ,  $(H_2, h_2)$  be two arguments of  $\mathcal{A}(\Sigma)$ .

1.  $(H_1, h_1)$  rebuts  $(H_2, h_2)$  iff  $h_1 \equiv \neg h_2$ .
2.  $(H_1, h_1)$  undercuts  $(H_2, h_2)$  iff  $\exists h \in H_2$  such that  $h_1 \equiv \neg h$ .
3.  $(H_1, h_1)$  contradicts  $(H_2, h_2)$  iff  $(H_1, h_1)$  rebuts a subargument of  $(H_2, h_2)$ .

The binary relations rebut, undercut, and contradict gather all pairs of arguments satisfying conditions (1), (2) and (3) respectively.

Definitions of rebut, undercut, and contradict will be given below and we will collectively refer to the relations as defeat if no distinction is necessary or we are describing them collectively. Following Dung's work [15], we have the following component definitions:

**Definition 4.** An argumentation framework is a pair,  $\text{Args} = \langle \mathcal{A}, \mathcal{R} \rangle$ , where  $\mathcal{A}$  is a set of arguments, and  $\mathcal{R}$  is the binary relation defeat over the arguments.

**Definition 5.** Let  $\langle \mathcal{A}, \mathcal{R} \rangle$  be an argumentation framework, and  $S \subseteq \mathcal{A}$ . An argument  $A$  is defended by  $S$  iff  $\forall B \in \mathcal{A}$  if  $(B, A) \in \mathcal{R}$  then  $\exists C \in S$  such that  $(C, B) \in \mathcal{R}$ .

**Definition 6.**  $S \subseteq \mathcal{A}$ .  $\mathcal{F}_{\mathcal{R}}(S) = \{A \in \mathcal{A} \mid A \text{ is defended by } S \text{ with respect to } \mathcal{R}\}$ .

Now, for a function  $F : D \rightarrow D$  where  $D$  is the domain and the range of the function, a fixed point of  $F$  is an  $x \in D$  such that  $x = F(x)$ . When the  $D$  is associated with an ordering  $P$  — for example,  $P$  can be set inclusion over the power set  $D$  of arguments —  $x$  is a least fixpoint of  $F$  if  $x$  is a least element of  $D$  with respect to  $P$  and  $x$  is a fixed point.

**Definition 7.** Let  $\langle \mathcal{A}, \mathcal{R} \rangle$  be an argumentation framework. The set of acceptable arguments, denoted by  $\text{Acc}_{\mathcal{R}}^F$ , is the least fixpoint of the function  $\mathcal{F}_{\mathcal{R}}$  with respect to set inclusion.

The least fixpoint semantics can be viewed as a mathematical translation of the principle that an argument survives if it can defend itself and be defended by a set of arguments which can also survive all the attacks made upon them.

## 4 Representing Arguments in QBFs and BDDs

We now turn our attention to using QBFs and BDDs to represent the components of an argumentation system, and then to perform the computations we need to carry out on that representation.

We can label each item  $f_i \in \Sigma$  with a proposition  $l_i$ . Namely, we will extend the language  $\mathcal{L}$  to contain both the information base  $\Sigma$  and the labels for these sentences. Formally, the proposition symbols can be extended to be  $\mathcal{P} = \mathcal{P}_D \cup \mathcal{P}_L$  where  $\mathcal{P}_D$  is the set of proposition symbols for the domain information, and  $\mathcal{P}_L$  is the set of system proposition symbols labeling the sentences in  $\Sigma$ . Given a finite information base  $\Sigma \subseteq \mathcal{L}$ ,  $|\mathcal{P}_L(\Sigma)| = |\Sigma|$ , namely each sentence  $f_i \in \Sigma$  has a corresponding label  $l_i$ .

For any formula  $\xi$  in  $\mathcal{L}$  based on  $\mathcal{P} = \mathcal{P}_D \cup \mathcal{P}_L$ ,  $\xi_D = \exists_{\mathcal{P}_L} \xi$  is the formula with only domain symbols left, and  $\xi_L = \exists_{\mathcal{P}_D} \xi$  is the formula with only the label symbols left.

### 4.1 Labeling

For representational convenience, we define

$$SEL(l_i) = l_i \wedge \bigwedge_{j \neq i} \neg l_j.$$

A sentence  $f_i$  of  $\Sigma$  corresponds to a pair  $\langle SEL(l_i), f_i \rangle$  which can be represented by  $SEL(l_i) \wedge f_i$ . Given a set of input information  $\Sigma = \{f_i\}$  for  $f_i \in \mathcal{L}_D$ , a labeling table  $A(\Sigma)$  can be expressed as follows

$$A(\Sigma) = \{\langle SEL(l_i), f_i \rangle\}$$

where  $f_i \in \Sigma$  and  $l_i \in \mathcal{P}_L$ , and the corresponding QBF representation

$$\xi(A(\Sigma)) = \bigvee_{f_i \in \Sigma} [SEL(l_i) \wedge f_i]$$

The above  $A(\Sigma)$  expression requires  $2 \times |\Sigma|$  QBF/BDD operations.<sup>2</sup> Given a subset  $\sigma \subseteq \Sigma$ ,

$$SEL(\sigma) = \bigwedge_{f_i \in \sigma} (l_i) \wedge \bigwedge_{f_j \notin \sigma} \neg l_j$$

### 4.2 Consistent Subsets

Since the support of an argument is a consistent set of propositions, a natural place to start thinking about argument computation is with the computation of consistent subsets. The set of all consistent subsets of  $\Sigma$  is

$$CONS(\Sigma) = \bigvee_{\sigma \subseteq \Sigma} [SEL(\sigma) \wedge \bigwedge_{f_i \in \sigma} f_i] \tag{1}$$

Computing the above expression directly requires an exponential number of QBF/BDD operations, so we want to find another way to compute it.

<sup>2</sup> The first condition of using QBF/BDD is to guarantee a way to express the information/specification that we need with only polynomial, linear, or even a logarithmic number of QBF/BDD operations; the second condition is to guarantee that the size of the initial, intermediate, and final BDDs corresponding to the information/specification is small enough to fit into memory. For the second condition, if the size of the BDD explodes, we may partition the expression into conjunctions or disjunctions, and modify the algorithms manipulating these BDDs correspondingly to try to avoid the explosion. If this still fails, then it means that the problem cannot be efficiently handled by BDDs. In this case, it usually also means that some aspect of the information required to solve the problem is simply too complex.

**Proposition 1.**  $CONS(\Sigma)$  can be constructed using  $2 \times |\Sigma| - 1$  operations as follows

$$CONS(\Sigma) = \bigwedge_{f_i \in \Sigma} [l_i \rightarrow f_i]. \quad (2)$$

*Proof.* The form of formula 2 follows from

$$\begin{aligned} CONS(\Sigma) &= \bigwedge_{f_i \in \Sigma} [l_i \rightarrow f_i] \\ &= \bigwedge_{f_i \in \Sigma} [l_i \rightarrow (l_i \wedge f_i)] \\ &= \bigwedge_{f_i \in \Sigma} [\neg l_i \vee (l_i \wedge f_i)] \\ &= \bigvee_{\sigma \subseteq \Sigma} \left[ \bigwedge_{f_j \notin \sigma} \neg l_j \wedge \bigwedge_{f_i \in \sigma} (l_i \wedge f_i) \right] \\ &= \bigvee_{\sigma \subseteq \Sigma} [SEL(\sigma) \wedge \bigwedge_{f_i \in \sigma} f_i] \end{aligned}$$

$(l_i \rightarrow f_i) \leftrightarrow (l_i \rightarrow (l_i \wedge f_i))$  follows from:

$$\begin{aligned} A \rightarrow B &\leftrightarrow \neg A \vee B \\ &\leftrightarrow (\neg A \vee A) \wedge (\neg A \vee B) \\ &\leftrightarrow \neg A \vee (A \wedge B) \\ &\leftrightarrow A \rightarrow (A \wedge B) \end{aligned}$$

With the above expression, we can exclude empty consistent subsets by □

$$CONS^+(\Sigma) = CONS(\Sigma) \wedge \left( \bigvee_{f_i \in \Sigma} l_k \right)$$

Because we are only interested in non-empty consistent subsets, from here on we will mean  $CONS^+(\Sigma)$  when we use  $CONS(\Sigma)$ . The set of subsets of selected sentences is

$$\begin{aligned} CONS(\Sigma)_L &= \exists_{\mathcal{P}_D} \left( \bigvee_{\sigma \in \Sigma} (SEL(\sigma) \wedge \bigwedge_{f_i \in \sigma} f_i) \right) \\ &= \bigvee_{\sigma \in \Sigma} (SEL(\sigma) \wedge \left[ \exists_{\mathcal{P}_D} \left( \bigwedge_{f_i \in \sigma} f_i \right) \right]) \\ &= \bigvee_{\sigma \subseteq \Sigma} SEL(\sigma) \end{aligned}$$

As we see, the complexity of  $CONS_L(\Sigma)$  is  $O(2 \times |\Sigma| - 1 + |\mathcal{P}_D|)$ . Let

$$CONJ(\sigma) = SEL(\sigma) \wedge \bigwedge_{f_i \in \sigma} f_i.$$

**Proposition 2.** *Given a sentence set selector  $\sigma \subseteq \Sigma$  represented by  $SEL(\sigma)$ , if the conjunction of the selected sentences in  $\sigma$  is consistent then it can be expressed as follows*

$$\begin{aligned} CONJ(\sigma) &= SEL(\sigma) \wedge CONS(\Sigma) \\ &= \bigvee_{\sigma \subseteq \Sigma} [SEL(\sigma) \wedge \bigwedge_{f_i \in \sigma} f_i] \end{aligned}$$

*Proof.*  $CONS(\Sigma)$  is a disjunction of conjunctions of all consistent subsets of  $\Sigma$ . Among these conjunctions,  $SEL(\sigma)$  only can make the one corresponding to the  $\sigma$  selection true, which is  $SEL(\sigma) \wedge \bigwedge_{f_i \in \sigma} f_i$ , and others false. Namely  $SEL(\sigma) \wedge CONS(\Sigma) = SEL(\sigma) \wedge \bigwedge_{f_i \in \sigma} f_i$   $\square$

Similarly, the set of conjunctions of a set of selected sentences can be expressed by:

$$CONJ(\{\sigma_i\}) = \bigvee_{\sigma_i} [SEL(\sigma_i)] \wedge CONS(\Sigma)$$

With this expression, we will be able to filter combinations of consistent and inconsistent sets of sentences into consistent sets.

### 4.3 QBF/BDD Representation of Arguments

We can extend the language  $\mathcal{P}$  further to contain  $\mathcal{P} = \mathcal{P}_L \cup \mathcal{P}_D \cup \mathcal{P}_{L,C} \cup \mathcal{P}_{D,C}$  where  $\mathcal{P}_{D,C}$  is a set of renaming symbols of  $\mathcal{P}_D$  to represent the conclusions of arguments;  $\mathcal{P}_{L,C}$  is an optional set of symbols to label an interesting sub-space of conclusions (the ones we want to compute arguments for). For example, if the sentences in  $\Sigma$  and their negations are of interest, then  $\mathcal{P}_{L,C} = 2\log|\Sigma|$  (we don't need to label a set of sentences, instead we just need to label individual sentences and their negations so that we need  $2\log|\Sigma|$  symbols). Similarly, we will denote  $\mathcal{P}_D$  by  $\mathcal{P}_{D,P}$  for premises when a distinction is needed.

An argument  $(H, h)$  in  $\mathcal{L}_D$  can then be represented by formula  $\xi(H, h)$  in  $\mathcal{L}$

$$\xi(H, h) = SEL(H) \wedge \bigwedge_{f_i \in H} f_i \wedge h[\mathcal{P}_{D,C}]$$

where  $h[\mathcal{P}_{D,C}]$  means the expression  $h$  is in terms of the symbols of  $\mathcal{P}_{D,C}$ .

The set of all arguments that can be constructed from  $\Sigma$  will be equivalent to

$$\mathcal{A}(\Sigma) = CONS(\Sigma)$$

for the moment by abstracting away the conclusions. Later we will reintroduce the conclusions to the representation during the query for conclusions and the defeat process.

### 4.4 Arguments for Conclusions

We can construct the set of arguments for a set of conclusions all at once as follows. Let us assume that, besides the input information base  $\Sigma$ , we also have a set of conclusions  $C$  that we wish to support.

$$C = \{h_k\}$$

with  $K = \log|C|$  and a set of labeling symbols

$$\mathcal{P}_{L,C} = \mathcal{L}_C = \{l_{1,C}, \dots, l_{K,C}\}$$

and let  $c_k$  be defined as  $\mathcal{L}_C = k$ , namely  $c_k$  is the encoding of integer  $k$  using the boolean symbols of  $\mathcal{P}_{L,C}$ .

The set of arguments for  $C$  based on  $\Sigma$  can be represented as

$$Args(\Sigma, C)_L = \forall_{\mathbf{x} \in \mathcal{P}_D \cup \mathcal{P}_{D,C}} \bigvee_{h_k \in C} \bigvee_{\sigma \subseteq \Sigma} \bigwedge_{f_i \in \Sigma} [( \bigwedge_{f_i \in \Sigma} f_i \rightarrow h_k ) \wedge SEL(\sigma) \wedge c_k]$$

and results in

$$Args(\Sigma, C) = Args(\Sigma, C)_L \wedge CONS(\Sigma) \wedge \bigvee_{h_k \in C} (c_k \wedge h_k)$$

**Proposition 3.**  $Args(\Sigma, C)_L$  can be expressed as

$$\forall_{\mathbf{x} \in \mathcal{P}_D \cup \mathcal{P}_{D,C}} CONS(\Sigma)_L \wedge \left[ \bigvee_{h_k \in C} (c_k) \right] \wedge \left[ \bigvee_{f_i \in \Sigma} (l_i \wedge \neg f_i) \vee \bigvee_{h_k \in C} (c_k \wedge h_k) \right]$$

using  $O(2 \times |C|) + O(|\Sigma|) + O(2 \times |\Sigma| + |\mathcal{P}_D|) + O(|\mathcal{P}_D \cup \mathcal{P}_{D,C}|)$  QBF/BDD operations.

*Proof.* Start with the first two items above,

$$CONS(\Sigma)_L \wedge \left[ \bigvee_{h_k \in C} (c_k) \right] \bigvee_{\sigma \subseteq \Sigma} \bigvee_{h_k \in H} [SEL(\sigma) \wedge c_k]$$

Conjoining with the remaining two items  $\left[ \bigvee_{f_i \in \Sigma} (l_i \wedge \neg f_i) \vee \bigvee_{h_k \in C} (c_k \wedge h_k) \right]$ , gives:

$$\begin{aligned} & \bigvee_{\sigma \subseteq \Sigma} \bigvee_{h_k \in C} \left[ SEL(\sigma) \wedge c_k \wedge \left( \bigvee_{f_i \in \Sigma} (l_i \wedge \neg f_i) \vee \bigvee_{h_k \in C} (c_k \wedge h_k) \right) \right] \\ &= \bigvee_{\sigma \subseteq \Sigma} \bigvee_{h_k \in \Sigma} \left[ \left( SEL(\sigma) \wedge \left( \bigvee_{f_i \in \sigma} \neg f_i \right) \right) \vee (c_k \wedge h_k) \right] \\ &= \bigvee_{\sigma \subseteq \Sigma} \bigvee_{h_k \in \Sigma} \left[ SEL(\sigma) \wedge c_k \wedge \left( \bigwedge_{f_i \in \sigma} (f_i) \rightarrow (h_k) \right) \right] \end{aligned}$$

The first line is derived using  $\bigvee_i A_i \wedge \bigvee_j B_j = \bigvee_i \left[ A_i \wedge \left( \bigvee_j B_j \right) \right]$ . The second line is derived using

$$SEL(\sigma) \wedge \left( \bigvee_{f_i \in \Sigma} (l_i \wedge \neg f_i) \right) = SEL(\sigma) \wedge \left( \bigvee_{f_i \in \sigma} \neg f_i \right)$$

since:

$$SEL(\sigma) \wedge (l_i \wedge \neg f_i) = FALSE$$

for any  $f_i \notin \sigma$ . The second line also employs  $c_k \wedge \bigvee_{h_k \in C} (c_k \wedge h_k) = c_k \wedge h_k$   $\square$

#### 4.5 Minimization of Consistent Sets with Respect to Conclusions

In the above,  $Args(\Sigma, C)$  may contain non-minimal arguments. To overcome this, we need to minimize the arguments in  $Args(\Sigma, C)$  with respect to their conclusions. Given a set of arguments  $Q \subseteq \mathcal{A}$  and a partial relation  $B \subseteq \mathcal{A} \times \mathcal{A}$  (e.g. the set-inclusion  $\subseteq$  relation on the supports of arguments) on  $\mathcal{A}$ , the set of minimal arguments in  $Q$  with respect to  $B$  is

$$Min(Q, B) = \{A \in Q \mid \text{for all } C \in Q, (C, A) \in B \text{ implies } (A, C) \in B\}$$

By encoding  $Q$  with a QBF formula  $Q[\mathcal{P}]$  based on a set  $\mathcal{P}$  of propositional symbols, and encoding the partial relation  $B$  with another QBF  $B[\mathcal{P}, \mathcal{P}']$  with the first component of  $B$  based on symbols in  $\mathcal{P}$  and the second component of  $B$  based on the symbols in  $\mathcal{P}'$ , we can compute  $Min(Q, B)$  as follows

$$Min(Q, B) = Q \wedge \forall_{\mathcal{Z}} [(Q[\mathcal{P}/\mathcal{Z}] \rightarrow (B[\mathcal{P}/\mathcal{Z}, \mathcal{P}'/\mathcal{P}] \rightarrow B[\mathcal{P}'/\mathcal{Z}]])]$$

where  $\mathcal{Z}$  is a temporary set of symbols renamed from  $\mathcal{P}$  to hold the intermediate results during the computation.

The set-inclusion relation between two sets of supports  $H_1[\mathcal{P}]$  and  $H_2[\mathcal{P}']$  can be implemented as:

$$\xi(\subseteq) = \bigwedge_{f_i \in \Sigma} [l_i \rightarrow l'_i].$$

This requires  $2 \times |\Sigma|$  QBF/BDD operations to construct. A linear BDD size implementation of  $\subseteq$  on the supports of  $\mathcal{A}$  is given in Algorithm 4.11.

The set of minimal supports which attack a sentence  $h_k \in C$  can be computed as

$$Args_{min}(\Sigma, h_k) = Min((Args(\Sigma, C) \wedge c_k)_L, \xi(\subseteq)).$$

The set of minimal supports with respect to each sentence in  $C$  can be computed as

$$Args_{min}(\Sigma, h_k) = \bigvee_{h_k \in C} Args_{min}(\Sigma, h_k).$$

For convenience of description, below we will use  $Args(\Sigma, C)$  for  $Args_{min}(\Sigma, C)$ .

#### 4.6 A QBF Representation of Defeat

A defeat relation  $\text{defeat}((H, h), (H', h'))$  can be represented by

$$\begin{aligned} \xi(H, h, H', h') = & CONJ(H) \wedge SEL(h)[\mathcal{P}_{L,C}] \wedge h[\mathcal{P}_{D,C}] \\ & \wedge CONJ(H')[\mathcal{P}'_D] \wedge SEL(h')[\mathcal{P}'_{L,C}] \wedge h'[\mathcal{P}'_{D,C}] \end{aligned}$$

by extending the language  $\mathcal{L}[\mathcal{P}]$  to be  $\mathcal{L}[\mathcal{P}] \cup \mathcal{L}[\mathcal{P}']$ . A defeat relation  $D = \{(A_i, A'_i)\}$  can be represented by a single QBF/BDD formula:

$$\xi(D) = \bigvee_{(A_i, A'_i) \in D} [\xi(A_i) \wedge \xi(A'_i)].$$

Now we need an expression with a polynomial number of operations to generate the set of all possible defeats from  $\Sigma$ . To do this, we need to inspect the specific types of defeats. We start with undercut:

<sup>3</sup> To the best of our knowledge only an exponential implementation exists in the literature [3].

**Algorithm 4.1.** Computing BDD for set-inclusion  $\subseteq$ 


---

```

1: Associate with each element in  $f_i \in \Sigma$  two BDD variables  $l_i$  and  $l'_i$ .
2: Take the variable order  $l_1, l'_1, l_2, l'_2, \dots, l_n, l'_n$  ( $n = |\Sigma|$ )
3: for each  $l_i$  do
4:   link  $l_i = 1$  to  $l'_i$ 
5:   link  $l_i = 0$  to  $l_{i+1}$ 
6: end for
7: for each  $l'_i \neq l'_n$  do
8:   link  $l'_i = 1$  to  $l_{i+1}$ 
9:   link  $l'_i = 0$  to terminal 0
10: end for
11: link  $l'_n = 0$ , to terminal 0
12: link  $l'_n = 1$ , to terminal 1

```

---

**Definition 8.** An argument  $(H_1, h_1)$  undercuts another argument  $(H_2, h_2)$  iff there exists an  $f \in H_2$  such that  $h_1 \equiv \neg f$ .

and this gives us:

**Proposition 4.** Let  $C = \Sigma \cup \{\neg f_i | f_i \in \Sigma\}$ , the set of all possible undercuts can be constructed as

$$\text{undercut}(\Sigma) = \text{Args}(\Sigma, C) \wedge \text{Args}(\Sigma', C') \wedge \left( \bigvee_{f'_i \in \Sigma'} (c_{\neg f_i} \wedge l_i) \right)$$

where  $c_{\neg f_i}$  denotes the encoding of the label that corresponds to  $\neg f_i$  in  $C$ .

*Proof.*  $\text{Args}(\Sigma, C)$  and  $\text{Args}(\Sigma', C')$  constructs the arguments for  $C$  based on  $\Sigma$  using two sets of symbols, and the corresponding selection of input sentences and conclusion sentences.  $\bigvee_{f'_i \in \Sigma'} (c_{\neg f_i} \wedge l_i)$  builds up the undercut relation between these two sets of arguments.  $\square$

Note that the setting of the conclusion points  $C = \Sigma \cup \{\neg f_i | f_i \in \Sigma\}$  can be changed according to any application-dependent argumentation process, for example  $CONS(\Sigma)$  and their negations as required.

Next we consider rebut:

**Definition 9.**  $(H_1, h_1)$  rebuts  $(H_2, h_2)$  iff  $h_1 \equiv \neg h_2$ .

We can construct the rebut relation in the same way as undercut, by assuming a set of interesting conclusions the status of which we want to discover. However, we can also construct the rebut relation in the following way thus leaving the conclusions unspecified, making the system more flexible.

**Definition 10.** Given a set  $H$  of sentences, let

$$\begin{aligned} S(H) &= \{s | s \models H\} \\ S(h) &= \{s | s \models h\} \end{aligned}$$

where  $s$  is an assignment to  $\mathcal{P}$ .  $H \vdash h$  iff  $S(H) \subseteq S(h)$ .

The definition of rebut is then:



**Definition 11.**  $H_1$  rebuts  $H_2$ , if there is some  $h$  such that  $H_1 \vdash h$  and  $H_2 \vdash \neg h$ .

and we have:

**Proposition 5.** Given two consistent sets of sentences  $H_1$  and  $H_2$ ,  $H_1$  rebuts  $H_2$  iff  $S(H_1) \cap S(H_2) = \emptyset$ , namely  $[CONJ(H_1) \wedge CONJ(H_2)] \leftrightarrow FALSE$ .

*Proof.* If  $H_1$  rebuts  $H_2$ , then there is a  $h$  such that  $H_1 \vdash h$  and  $H_2 \vdash \neg h$ . Since  $S(h) \cap S(\neg h) = \emptyset$ , and  $S(H_1) \subseteq S(h)$  and  $S(H_2) \subseteq S(\neg h)$ , therefore  $S(H_1) \cap S(H_2) = \emptyset$ .

If  $S(H_1) \cap S(H_2) = \emptyset$ , the rebutting point  $h$  can be constructed as follows. Let  $padding = \neg(H_1 \vee H_2)$ , and  $h = H_1 \vee padding$ . In this way,  $S(padding) = \mathcal{U} \setminus (S(H_1) \cup S(H_2))$ ,  $S(h) = S(H_1) \cup S(padding)$ ,  $S(\neg h) = \mathcal{U} \setminus (S(H_1) \cup S(padding)) = S(H_2)$ . Therefore  $S(H_1) \subseteq S(h)$  and  $S(H_2) \subseteq S(\neg h)$ , namely  $h$  is the rebutting point we are looking for such that  $H_1 \vdash h$  and  $H_2 \vdash \neg h$ .  $\square$

Actually  $h$  can be anything such that  $S(H_1) \subseteq S(h) \subseteq (S(H_1) \cup S(padding))$ , so we have the following corollary.

**Corollary 1.** Given two sets of sentences  $H_1$  and  $H_2$  which rebut each other, the rebut point  $h$  can be obtained by setting  $S(H_1) \subseteq S(h) \subseteq S(H_1) \cup S(padding)$  where  $padding = \neg(H_1 \vee H_2)$ . The choice of  $h = H_1 \vee \neg(H_1 \vee H_2)$  which makes  $H_1$  and  $H_2$  be the minimal sets of sentences such that  $H_1 \vdash h$  and  $H_2 \vdash \neg h$   $\square$

As a result, the set of all rebuts can be expressed as

$$\text{rebut}(\Sigma) = \bigvee_{\sigma \subseteq \Sigma, \sigma' \subseteq \Sigma} [CONJ(\sigma) \wedge CONJ'(\sigma') \wedge \neg(CONJ(\sigma)_D \wedge CONJ(\sigma')_D)]$$

and we have:

**Proposition 6.**  $\text{rebut}(\Sigma)$  can be expressed as

$$\text{rebut}(\Sigma) = CONS(\Sigma) \wedge CONS'(\Sigma) \wedge \left[ \bigvee_{f_i \in \Sigma} (l_i \wedge \neg f_i) \vee \bigvee_{f_j \in \Sigma} (l'_j \wedge \neg f_j) \right]$$

using  $2.O(CONS(\Sigma)) + 6.|\Sigma| + 3$  QBF/BDD operations.

*Proof.*

$$\begin{aligned} & \text{rebut}(\Sigma) \\ &= \bigvee_{\sigma \subseteq \Sigma, \sigma' \subseteq \Sigma} [CONJ(\sigma) \wedge CONJ'(\sigma') \wedge \neg(CONJ(\sigma)_D \wedge CONJ(\sigma')_D)] \\ &= \bigvee_{\sigma \subseteq \Sigma, \sigma' \subseteq \Sigma} [CONJ(\sigma) \wedge CONJ'(\sigma') \wedge (\neg CONJ(\sigma)_D \vee \neg CONJ(\sigma')_D)] \\ &= \bigvee_{\sigma \subseteq \Sigma, \sigma' \subseteq \Sigma} \left[ CONJ(\sigma) \wedge CONJ'(\sigma') \wedge \left( \bigvee_{f_i \in \sigma} \neg f_i \vee \bigvee_{f_j \in \sigma'} \neg f_j \right) \right] \end{aligned}$$

$$\begin{aligned}
&= \bigvee_{\sigma \subseteq \Sigma, \sigma' \subseteq \Sigma} \left[ \text{CONJ}(\sigma) \wedge \text{CONJ}'(\sigma') \wedge \left( \bigvee_{f_i \in \sigma} (l_i \wedge \neg f_i) \vee \bigvee_{f_j \in \sigma'} (l'_j \wedge \neg f_j) \right) \right] \\
&= \text{CONS}(\Sigma) \wedge \text{CONS}'(\Sigma) \wedge \left[ \bigvee_{f_i \in \Sigma} (l_i \wedge \neg f_i) \vee \bigvee_{f_j \in \Sigma} (l'_j \wedge \neg f_j) \right]
\end{aligned}$$

□

Finally we consider the computation of the contradict relation:

**Definition 12.**  $(H_1, h_1)$  *contradicts*  $(H_2, h_2)$  if and only if  $(H_1, h_1)$  *rebuts* a subargument of  $(H_2, h_2)$ .

The contradict relation can be computed by

$$\text{contradict}(\Sigma) = \exists_{\mathcal{Z}} (\text{rebut}(\Sigma)[\mathcal{P}'/\mathcal{Z}] \wedge \xi(\subseteq)[\mathcal{P}/\mathcal{Z}])$$

#### 4.7 Computing Fixed Points of Argumentation

The relations undercut, rebut and contradict give us the relationship between individual arguments, but, as is usual, we are more interested in computing properties of arguments such as whether arguments are *acceptable*, where such properties are defined as fixed-points.

**Definition 13.** An argument  $H$  *defends* another argument  $H'$  if there exists another argument  $H''$  such that  $H''$  *defeats*  $H'$  but  $H$  *defeats*  $H''$ .

The defend relation can be constructed from the defeat relation on the set of arguments as follows:

$$\text{defend}(\Sigma, \text{defeat}) = \exists_{\mathcal{Z}} (\text{defeat}(\Sigma)[\mathcal{P}'/\mathcal{Z}] \wedge \text{defeat}(\Sigma)[\mathcal{P}/\mathcal{Z}])$$

where  $\text{defeat}(\Sigma)$  is either  $\text{undercut}(\Sigma)$ ,  $\text{rebut}(\Sigma)$ ,  $\text{contradict}(\Sigma)$ , or any disjunction of the relations, such as  $\text{undercut}(\Sigma) \vee \text{rebut}(\Sigma)$ . The composition of two relations  $R_1$  and  $R_2$  on the set  $\mathcal{A}$  of arguments can be computed by

$$\text{Compose}R(R_1, R_2) = \exists_{\mathcal{Z}} R_1[\mathcal{P}'/\mathcal{Z}] \wedge R_2[\mathcal{P}/\mathcal{Z}].$$

With these constructs defined, the fixed point of argumentation can be computed using Algorithm 4.2. In Algorithm 4.2 the closure of a binary relation  $R$  on  $\mathcal{A}$ , is computed using a method called iterative squaring [6] which is guaranteed to terminate within  $O(\log|\mathcal{A}|)$  steps. In line 3:

$$\text{Old}R \leftarrow I_{\mathcal{P}_L} \cup \text{defend}_{\mathcal{P}_L}$$

the defend relation is first projected to sentence labeling symbols so that during the computation of the defending closure only the labels of arguments are considered without referring to their internal structure; the union with the identity relation

$$I_{\mathcal{P}_L} = \bigwedge_{f_i \in \Sigma} (l_i \leftrightarrow l'_i)$$

is to keep the defended arguments in the closure.

**Algorithm 4.2.** Computing Fixed Point of Argumentation

---

```

1: function ComputeFixedpoint( $\Sigma$ , defeat) {
  (1)  $\Sigma$ : The set of input information
  (2) defeat is binary relation on  $\mathcal{A}$  }
2: defend  $\leftarrow$  defend( $\Sigma$ , defeat)
3: OldR  $\leftarrow$   $I_{\mathcal{P}_L} \cup$  defend $_{\mathcal{P}_L}$ 
4: R  $\leftarrow$  FAIL
5: while (OldR  $\neq$  R) do
6:   tmpR  $\leftarrow$  R
7:   R  $\leftarrow$  ComposeR(OldR, OldR)
8:   OldR  $\leftarrow$  tmpR
9: end while
10: Undeferred  $\leftarrow$  CONS( $\Sigma$ )  $\wedge$   $\neg$  ( $\exists x \in \mathcal{P}$  defeat) [ $\mathcal{P}'/\mathcal{P}$ ]
11: Acc  $\leftarrow$   $\exists x \in \mathcal{P}$  (Undeferred  $\wedge$  R) [ $\mathcal{P}'/\mathcal{P}$ ]  $\vee$  Undeferred
12: return Acc end function

```

---

**Proposition 7.** Algorithm 4.2 computes the fixed point of the defend relation, namely the set of acceptable arguments constructed from  $\Sigma$ .

*Proof.* Let  $step(R)$  be the maximum length of paths between a pair  $(A, A') \in R$  in the induced graph of the defend relation defend. Let the starting  $R$  in line 3 denoted by  $R_0 = defend \cup I$ . In  $R_0$ , for every  $(A, A') \in R$ , either  $(A, A') \in defend$ , namely  $A$  defends  $A'$  using one step, or  $A$  is identical to  $A'$  namely  $A$  defends  $A'$  using 0 step, therefore  $step(R_0) = 1$ . Let the consequent content of  $R$  in each *while* iteration denoted by  $R_i$  where  $i$  is the number of the iteration. Each time, when  $R_{i+1} \leftarrow ComposeR(R_i, R_i)$  is applied in line 7,  $R_{i+1}$  will gather all the argument pairs of the form  $(A, A')$  such that  $A$  defends  $A'$  using defending steps less or equal than  $step(R_{i+1}) = 2 \times step(R_i)$  steps. Assume that  $i$  is the number such that  $R_{i+1} = R_i$ , if the iteration continues we will have:

$$R_{i+2} = ComposeR(R_{i+1}, R_{i+1}) = ComposeR(R_i, R_i) = R_{i+1} = R_i$$

namely for all  $j \geq i$ ,  $R_j = R_i$ . Therefore, after the *while* loop terminates  $R$  will gather all the argument pairs  $(A, A')$  via any number of defending steps. Since the number of arguments is finite, all the defending paths are of finite length, therefore the algorithm is guaranteed to terminate.  $\square$

**Proposition 8.** The complexity of Algorithm 4.2 is  $O(|\Sigma| \cdot K^2 \cdot |\mathcal{P}|)$  where  $K$  is the maximum size of the BDDs which appear during the fixed point computing process.

*Proof.* As analyzed in the proof of proposition 7, the  $step(R_i) = step^2(R_{i+1})$ . The maximum possible step of  $R_i$ s is the number of arguments which is  $2^{|\Sigma|}$ . Therefore the algorithm is guaranteed to terminate after  $m = \log_2 2^{|\Sigma|} = |\Sigma|$  iterations, therefore the number of iterations is bounded above by  $O(|\Sigma|)$ . In each iteration, *ComposeR* can be computed using  $O(1 + |\mathcal{P}|)$  number of BDD operations, each operation is of complexity  $O(K^2)$  where  $K$  is the maximum size of BDDs used. Therefore the whole algorithm is bounded above by  $O(|\Sigma|) \cdot O(K^2 \cdot |\mathcal{P}|)$ .  $\square$

## 5 Discussion

Proposition 8 shows that we can compute the fixed-point in a polynomial number of BDD operations. As we mentioned above, this is a long way from saying that we can do general logical inference in polynomial time, rather what we are saying is that while the complexity of algorithm 4.2 depends on the maximum size of the BDD ( $K$ ), this doesn't depend on the size of  $\Sigma$  but rather on the complexity of the information contained in  $\Sigma$ . In the worse case,  $K$  can still be exponential in  $|\mathcal{P}|$ , but in many practical applications  $K$  tends to be small.

Because of this feature of systems built using the QBF/BDD representation, there has been a lot of work on reducing the size of BDDs. Many successful approaches have been developed in literature, especially those developed for symbolic model checking in software and hardware verification [24], and in non-deterministic AI planning [10]. Examples of techniques for reducing the size of BDDs are early quantification [19], quantification scheduling [9], transition partitioning [7], iterative squaring [6,8], frontier simplification [12], input splitting [28,29], and state set  $A^*$  branching [22,21,23] (a BDD version of the  $A^*$  search heuristic [31]).

Another factor affecting the BDD size greatly is variable ordering. The problem of finding an optimal variable ordering is NP-complete [4]. Algorithms based on dynamic programming [14], heuristics [20], dynamic variable reordering [30] and machine learning approaches [18] have been proposed for finding a good variable ordering in reasonable time<sup>4</sup>.

We are currently working on an implementation of the reasoning mechanism proposed above with the aim of experimentally clarifying the nature of  $K$  for different argumentation problems.

## 6 Conclusions and Future Work

In this paper, we have proposed a symbolic model checking approach to compute argumentation. The computation only uses a polynomial number of BDD operations in terms of the number of sentences in the input and the number of symbols used in the input. A key idea in the approach is to construct the set of consistent arguments all together using a polynomial number of BDD operations. In the same way, the defeat relation among these arguments can also be computed all at once using a polynomial number of BDD operations. And with the iterative squaring technique, we are able to compute the fixed point of a set of arguments in polynomial number of BDD operations.

We are currently working on implementing the BDD-based argumentation system proposed in this paper, with the aim of conducting experiments to classify the nature of the BDDs constructed for argumentation. This will allow us to determine how effective this approach will be in general. This in turn may lead us to look for new heuristics for controlling the size of the BDDs we need to construct to compute arguments. Another direction that we are working on is to extend the current method to compute more sophisticated and controllable approaches argumentation, such as those based on argumentation schemes [27]. On the way, we will need to develop BDD techniques to efficiently specify application-dependent patterns of arguments (such as those captured by

<sup>4</sup> [18] is also a good source for other references on BDD variable (re-)ordering.

argument schemes), specify application-dependent patterns of defeats (defeat schemes), and extend the basic entailment-based reasoning modelled here to specify the necessary patterns of rule-based procedural reasoning. In combination with our continuing efforts to use BDD techniques in multiagent planning and dialogues [34,35,36,37,38] all these efforts are aimed at our ultimate goal of a practical argumentation-based dialogue model for multiagent planning.

## Acknowledgment

Research was sponsored by the U.S. Army Research Laboratory and the U.K. Ministry of Defence and was accomplished under Agreement Number W911NF-06-3-0001. The views and conclusions contained in this document are those of the author(s) and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Army Research Laboratory, the U.S. Government, the U.K. Ministry of Defence or the U.K. Government. The U.S. and U.K. Governments are authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation hereon.

## References

1. Amgoud, L., Cayrol, C.: Inferring from inconsistency in preference-based argumentation frameworks. *Journal of Automated Reasoning* 29(2), 125–169 (2002)
2. Amgoud, L., Cayrol, C.: A reasoning model based on the production of acceptable arguments. *Annals of Mathematics and Artificial Intelligence* 34(1-3), 197–215 (2002)
3. Berghammer, R., Fronk, A.: Exact computation of minimum feedback vertex sets with relational algebra. *Fundam. Inf.* 70(4), 301–316 (2005)
4. Bollig, B., Wegener, I.: Improving the variable ordering of OBDDs is NP-Complete. *IEEE Transactions on Computers* 45(9), 993–1002 (1996)
5. Bryant, R.E.: Symbolic boolean manipulation with ordered binary-decision diagrams. *ACM Computing Surveys* 24(3), 293–318 (1992)
6. Burch, J., Clarke, E., McMillan, K., Dill, D., Burch, L.H.J.R., Clarke, E.M., Mcmilla, K.L., Dill, D.L., Hwang, L.J.: Symbolic Model Checking:  $10^{20}$  States and Beyond. In: *Proceedings of the Fifth Annual IEEE Symposium on Logic in Computer Science*, pp. 1–33. IEEE Computer Society Press, Washington, D.C (1990)
7. Burch, J.R., Clarke, E.M., Long, D.E.: Symbolic model checking with partitioned transition relations. In: *Proceedings of International Conference on Very Large Scale Integration*, pp. 49–58. North-Holland, Amsterdam (1991)
8. Cabodi, G., Camurati, P., Lavagno, L., Quer, S.: Disjunctive partitioning and partial iterative squaring: an effective approach for symbolic traversal of large circuits. In: *DAC 1997: Proceedings of the 34th Annual Conference on Design Automation*, pp. 728–733. ACM Press, New York (1997)
9. Chauhan, P., Clarke, E.M., Jha, S., Kukula, J., Shiple, T., Veith, H., Wang, D.: Non-linear quantification scheduling in image computation. In: *ICCAD 2001: Proceedings of the 2001 IEEE/ACM International Conference on Computer-aided Design*, pp. 293–298. IEEE Computer Society Press, Piscataway (2001)
10. Cimatti, A., Pistore, M., Roveri, M., Traverso, P.: Weak, strong, and strong cyclic planning via symbolic model checking. *Artificial Intelligence* 147(1-2), 35–84 (2003)

11. Coudert, O., Berthet, C., Madre, J.C.: Verification of synchronous sequential machines based on symbolic execution. In: *Automatic Verification Methods for Finite State Systems*, pp. 365–373 (1989)
12. Coudert, O., Madre, J.C.: The implicit set paradigm: a new approach to finite state system verification. *Formal Methods in System Design* 6(2), 133–145 (1995)
13. Dimopoulos, Y., Nebel, B., Toni, F.: On the computational complexity of assumption-based argumentation for default reasoning. *Artificial Intelligence* 141, 57–78 (2002)
14. Drechsler, R., Drechsler, N., Günther, W.: Fast exact minimization of BDDs. In: *DAC 1998: Proceedings of the 35th Annual Conference on Design Automation*, pp. 200–205. ACM Press, New York (1998)
15. Dung, P.M.: On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence* 77(2), 321–358 (1995)
16. Dunne, P.E., Bench-capon, T.J.M.: Two party immediate response disputes: properties and efficiency. *Artificial Intelligence* 149, 2003 (2001)
17. Gordon, T., Karacapilidis, N.: The zeno argumentation framework. In: *Proceedings of the Sixth International Conference on AI and Law*, pp. 10–18. ACM Press, New York (1997)
18. Grumberg, O., Livne, S., Markovitch, S.: Learning to order BDD variables in verification. *Journal of Artificial Intelligence Research (JAIR)* 18, 83–116 (2003)
19. Hojati, R., Krishnan, S.C., Brayton, R.K.: Early quantification and partitioned transition relations. In: *ICCD 1996: Proceedings of the 1996 International Conference on Computer Design, VLSI in Computers and Processors*, pp. 12–19. IEEE Computer Society Press, Washington, DC, USA (1996)
20. Jain, J., Adams, W., Fujita, M.: Sampling schemes for computing OBDD variable orderings. In: *ICCAD 1998: Proceedings of the 1998 IEEE/ACM International Conference on Computer-aided Design*, pp. 631–638. ACM Press, New York (1998)
21. Jensen, R.M., Bryant, R.E., Veloso, M.M.: An efficient BDD-based A\* algorithm. In: *Proceedings of AIPS 2002 Workshop on Planning via Model Checking* (2002)
22. Jensen, R.M., Bryant, R.E., Veloso, M.M.: SetA\*: An efficient BDD-based heuristic search algorithm. In: *Proceedings of 18th National Conference on Artificial Intelligence (AAAI 2002)*, pp. 668–673 (2002)
23. Jensen, R.M., Veloso, M.M., Bryant, R.E.: State-set branching: Leveraging BDDs for heuristic search. *Artificial Intelligence* 172(2-3), 103–139 (2008)
24. Jhala, R., Majumdar, R.: Software model checking. *ACM Comput. Surv.* 41(4), 1–54 (2009)
25. Kakas, A.C., Toni, F.: Computing argumentation in logic programming. *Journal of Logic and Computation* 9, 515–562 (1999)
26. Karacapilidis, N., Papadias, D.: Computer supported argumentation and collaborative decision making: The hermes system. *Information Systems* 26, 259–277 (2001)
27. Katzav, J., Reed, C.: On argumentation schemes and the natural classification of arguments. *Argumentation* 18(2) (2004)
28. Meinel, C., Theobald, T.: *Algorithms and Data Structures in VLSI Design*. Springer-Verlag New York, Inc, Secaucus (1998)
29. Moon, I.-H., Kukula, J.H., Ravi, K., Somenzi, F.: To split or to conjoin: the question in image computation. In: *DAC 2000: Proceedings of the 37th Conference on Design Automation*, pp. 23–28. ACM Press, New York (2000)
30. Panda, S., Somenzi, F., Plessier, B.F.: Symmetry detection and dynamic variable ordering of decision diagrams. In: *ICCAD 1994: Proceedings of the 1994 IEEE/ACM International Conference on Computer-aided Design*, pp. 628–631. IEEE Computer Society Press, Los Alamitos (1994)
31. Pearl, J.: *Heuristics: intelligent search strategies for computer problem solving*. Addison-Wesley Longman Publishing Co., Inc., Boston (1984)

32. Prakken, H.: Coherence and flexibility in dialogue games for argumentation. *Journal of Logic and Computation* 15(6), 1009–1040 (2005)
33. Rahwan, I., Ramchurn, S.D., Jennings, N.R., Mcburney, P., Parsons, S., Sonenberg, L.: Argumentation-based negotiation. *The Knowledge Engineering Review* 18(4), 343–375 (2003)
34. Tang, Y., Norman, T.J., Parsons, S.: A model for integrating dialogue and the execution of joint plans. In: *Proceedings of the Eighth International Joint Conference on Autonomous Agents and Multiagent Systems*, Budapest, Hungary, May 10-15 (2009)
35. Tang, Y., Norman, T.J., Parsons, S.: Towards the implementation of a system for planning team activities. In: *Proceedings of the Second Annual Conference of the ITA*. University of Maryland University College, Maryland (2009)
36. Tang, Y., Parsons, S.: Argumentation-based dialogues for deliberation. In: *Proceedings of the Fourth International Joint Conference on Autonomous Agents and Multiagent Systems*, pp. 552–559. ACM Press, New York (2005)
37. Tang, Y., Parsons, S.: Using argumentation-based dialogues for distributed plan management. In: *Proceedings of the AAAI Spring Symposium on Distributed Plan and Schedule Management*, Stanford (2006)
38. Tang, Y., Parsons, S.: A dialogue mechanism for public argumentation using conversation policies. In: *Proceedings of the Seventh International Joint Conference on Autonomous Agents and Multiagent Systems*, Estoril, Portugal, May 12-16, pp. 445–452 (2008)

# On Strategic Argument Selection in Structured Argumentation Systems

Matthias Thimm<sup>1</sup> and Alejandro J. García<sup>2</sup>

<sup>1</sup> Technische Universität Dortmund, Germany

<sup>2</sup> Universidad Nacional del Sur, Bahía Blanca, Argentina

**Abstract.** This paper deals with strategical issues of arguing agents in a multi-agent setting. We investigate different scenarios of such argumentation games that differ in the protocol used for argumentation, i. e. direct, synchronous, and dialectical argumentation protocols, the awareness that agents have on other agents beliefs, and different settings for the preferences of agents. We give a thorough investigation and classification of these scenarios employing structured argumentation frameworks which are an extension to Dung's abstract argumentation frameworks that give a simple inner structure to arguments. We also provide some game theoretical results that characterize a specific argumentation game as strategy-proof and develop some argumentation selection strategies that turn out to be the dominant strategies for other specific argumentation games.

## 1 Introduction

The study of computational models of argumentation [4] is a relatively novel research area in the field of artificial intelligence and non-monotonic reasoning with logic-based formalisms for knowledge representation. There are a lot of approaches to model argumentation in different kinds of logics, e. g. classical logic [5] or defeasible logic [19,12] and also abstract formalizations of argumentation [11] are widely used to talk about computational argumentation in general. In abstract argumentation, arguments are represented as atomic entities and the interrelationships between different arguments are modeled using an attack relation. Abstract argumentation has been thoroughly investigated in the past ten years and there is quite a lot of work on, e. g. semantical issues [3] and extensions of abstract argumentation frameworks [16,2].

In the context of agent and multi-agent systems, there are mainly two applications of formal argumentation. First, using argumentation techniques as a non-monotonic reasoning process within a single agent and second, using argumentation in dialogues between different agents in order to realize persuasion, cooperation, planning, or general conflict solving. Here, we focus on the second application where reasoning is performed involving the whole system of agents, see e. g. [13,11,7,25] for formalizations. In a dialogue, agents take turns in bringing up arguments for some given claim and depending on the interrelationships of



the arguments the claim is accepted or rejected by the agents (either individually or jointly). Up until recently, strategic issues in argumentation dialogues have been mostly ignored with few exceptions, e. g. [22]. By considering game theoretical aspects in argumentation dialogues [20] the interest in strategies for the selection of arguments and the general connection of game theory and argumentation has grown. From the point of view of game theory, an argumentation dialogue can be represented as a strategic game involving a set of self-interested agents and in choosing the “right” arguments agents can influence the outcome of the argumentation and reach a more desirable result according to their own preferences. In [20,21,17] Rahwan et al. investigate *direct argumentation mechanisms* in which agents have to state all arguments they wish at once. Under specific circumstances of the underlying argumentation framework they were able to prove strategy-proofness, i. e. the *dominant strategy* of each agent is to truthfully report all their arguments. Besides this scenario of direct argumentation there are other formalizations of specific argumentation games, e. g. [22,7]. But up till now, to our knowledge there has been no comprehensive overview on the different argumentation settings and the different scenarios where agents can argue with each other.

The contribution of this paper is twofold. The main contribution lies in a classification of the different argumentation games agents can play within a multi-agent setting. We make a first attempt to characterize argumentation games by means of the game protocol, the awareness of the agents on other agents’ beliefs, and the structure of the preferences of the agents. We use structured argumentation frameworks, a novel approach which generalizes abstract argumentation frameworks, to model argumentation between different agents. The second contribution lies in generalizing the strategy-proofness result of [20] and investigating several other settings for argumentation games in terms of the strategical issues involving argument selection.

This paper is a slightly extended version of a previously published paper [24] and is organized as follows. In Section 2 we give a brief overview on abstract argumentation and introduce the novel approach of structured argumentation. We continue in Section 3 with applying structured argumentation onto a multi-agent setting. Section 4 develops a classification of argumentation games in the multi-agent setting in terms of game protocol, awareness, and agent types. We investigate several strategical issues in some instances of argumentation games in Section 5 and conclude in Section 6.

## 2 Preliminaries

We first give a brief overview on *abstract argumentation frameworks* [11] and continue by introducing *structured argumentation frameworks* which extend abstract argumentation frameworks and are the means to model argumentation games in this paper.

## 2.1 Abstract Argumentation

*Abstract argumentation frameworks* [11] take a very simple view on argumentation as they do not presuppose any internal structure of an argument. Abstract argumentation frameworks only consider the interactions of arguments by means of an attack relation between arguments.

**Definition 1 (Abstract Argumentation Framework).** *An abstract argumentation framework AF is a tuple  $AF = (\text{Arg}, \text{attacks})$  where  $\text{Arg}$  is a set of arguments and attacks is a relation  $\text{attacks} \subseteq \text{Arg} \times \text{Arg}$ .*

For two arguments  $\mathcal{A}, \mathcal{B} \in \text{Arg}$  the relation  $(\mathcal{A}, \mathcal{B}) \in \text{attacks}$  means that argument  $\mathcal{A}$  attacks argument  $\mathcal{B}$ . Abstract argumentation frameworks can be concisely represented as directed graphs, where arguments are represented as nodes and edges model the attack relation.

*Example 1.* Consider the abstract argumentation framework  $AF = (\text{Arg}, \text{attacks})$  depicted in Figure 1. Here it is  $\text{Arg} = \{\mathcal{A}_1, \mathcal{A}_2, \mathcal{A}_3, \mathcal{A}_4\}$  and  $\text{attacks} = \{(\mathcal{A}_1, \mathcal{A}_2), (\mathcal{A}_2, \mathcal{A}_3), (\mathcal{A}_2, \mathcal{A}_4), (\mathcal{A}_3, \mathcal{A}_2), (\mathcal{A}_3, \mathcal{A}_4)\}$ .

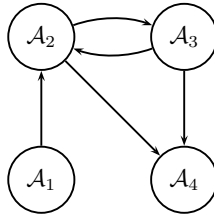


Fig. 1. A simple argumentation framework

Semantics are given to abstract argumentation frameworks by means of extensions. An *extension*  $E$  of an  $AF = (\text{Arg}, \text{attacks})$  is a set of arguments  $E \subseteq \text{Arg}$  that gives some coherent view on the argumentation underlying  $AF$ . In the literature [11,8] a wide variety of different types of extensions has been proposed. All these different types of extensions require some basic properties as *conflict-freeness* and *admissibility*. A set  $S \subseteq \text{Arg}$  is *conflict-free* if and only if there are no two arguments  $\mathcal{A}, \mathcal{B} \in \text{Arg}$  with  $(\mathcal{A}, \mathcal{B}) \in \text{attacks}$ . An argument  $\mathcal{A} \in \text{Arg}$  is *acceptable* with respect to a set of arguments  $S \subseteq \text{Arg}$  if and only if for every argument  $\mathcal{B} \in \text{Arg}$  with  $(\mathcal{B}, \mathcal{A}) \in \text{attacks}$  there is an argument  $\mathcal{C} \in S$  with  $(\mathcal{C}, \mathcal{B}) \in \text{attacks}$ . A set  $S \subseteq \text{Arg}$  is *admissible* if and only if it is conflict-free and every argument  $a \in S$  is acceptable with respect to  $S$ .

Extensions of an abstract argumentation framework can be described using the characteristic function  $F_{AF}(S) = \{\mathcal{A} \in \text{Arg} \mid \mathcal{A} \text{ is acceptable wrt. } S\}$  defined for sets  $S \subseteq \text{Arg}$ .

**Definition 2 (Extensions).** *Let  $AF = (\text{Arg}, \text{attacks})$  be an abstract argumentation framework and  $S \subseteq \text{Arg}$  an admissible set.*

- $S$  is a complete extension if and only if  $S = F_{AF}(S)$ .
- $S$  is a grounded extension if and only if it is a minimal complete extension (with respect to set inclusion).
- $S$  is a preferred extension if and only if it is a maximal complete extension (with respect to set inclusion).
- $S$  is a stable extension if and only if it is a complete extension and attacks each  $\mathcal{A} \in \text{Arg} \setminus S$ .

*Example 2.* We continue Example 1. As  $F_{AF}(\{\mathcal{A}_1, \mathcal{A}_3\}) = \{\mathcal{A}_1, \mathcal{A}_3\}$  the set  $\{\mathcal{A}_1, \mathcal{A}_3\}$  is a complete extension. Furthermore it is the only complete extension and also grounded, preferred, and stable.

Note that the grounded extension is uniquely determined and always exists [11].

## 2.2 Structured Argumentation

In the following, we introduce *structured argumentation frameworks* which extend Dung’s abstract argumentation frameworks and are a slightly modified variant of dynamic argumentation frameworks [23]. In structured argumentation frameworks arguments are built using a very simple propositional language, so let  $\text{Prop}$  denote a finite and fixed set of propositions. The basic structure for structured argumentation frameworks are *basic arguments* which represent atomic inference rules by connecting some set of propositions (the support) to another proposition (the claim).

**Definition 3 (Basic Argument).** *Let  $\alpha_1, \dots, \alpha_n, \beta \in \text{Prop}$  be some propositions with  $\beta \notin \{\alpha_1, \dots, \alpha_n\}$ . Then a basic argument  $\mathcal{A}$  is a tuple  $\mathcal{A} = (\{\alpha_1, \dots, \alpha_n\}, \beta)$ . We abbreviate  $\text{supp}(\mathcal{A}) = \{\alpha_1, \dots, \alpha_n\}$  (the support of  $\mathcal{A}$ ) and  $\text{cl}(\mathcal{A}) = \beta$  (the claim of  $\mathcal{A}$ ).*

For the rest of this paper, let  $U$  be some fixed and finite set of basic arguments, called the *universal set of basic arguments*. As such,  $U$  represents all possible basic arguments under consideration. To keep things simple, we assume that  $U$  does not contain any *cyclic dependencies*, i.e. there is no infinite sequence  $\mathcal{A}_1, \mathcal{A}_2, \dots \in U$  with  $\text{cl}(\mathcal{A}_i) \in \text{supp}(\mathcal{A}_{i+1})$  for all  $i > 0$ . Together with an attack relation  $\rightarrow \subseteq U \times U$  the set of basic arguments form a *structured argumentation framework* (SAF)  $\mathfrak{F} = (U, \rightarrow)$  [1].

*Example 3.* Consider the SAF  $\mathfrak{F}_1 = (U, \rightarrow)$  given by

$$\begin{aligned}
 U = \{ & \mathcal{A}_1 = (\emptyset, a), & \mathcal{A}_2 = (\{a\}, b), & \mathcal{A}_3 = (\emptyset, c) \\
 & \mathcal{A}_4 = (\emptyset, d), & \mathcal{A}_5 = (\{d\}, e), & \mathcal{A}_6 = (\{b\}, f) \\
 & \mathcal{A}_7 = (\emptyset, g) & \}
 \end{aligned}$$

and

$$\begin{aligned}
 \rightarrow = \{ & (\mathcal{A}_3, \mathcal{A}_2), (\mathcal{A}_2, \mathcal{A}_4), (\mathcal{A}_5, \mathcal{A}_6), \\
 & (\mathcal{A}_5, \mathcal{A}_7), (\mathcal{A}_6, \mathcal{A}_7), (\mathcal{A}_7, \mathcal{A}_5) & \} .
 \end{aligned}$$

---

<sup>1</sup> Although SAFs have the same structure as abstract argumentation frameworks, we deliberately use different notations to avoid ambiguity.

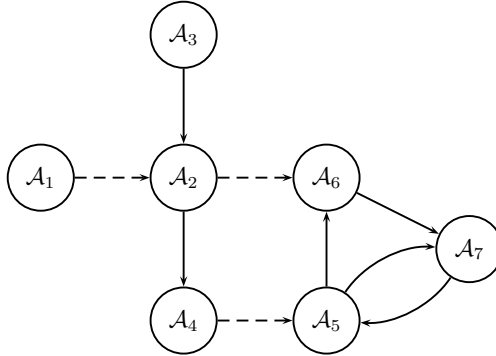


Fig. 2. The SAF  $\mathfrak{F}_1$

The rough structure of  $\mathfrak{F}_1$  is depicted in Figure 2, where the attack relation is represented by solid arrows and “support” by dashed arrows. Notice that Figure 2 does not contain all the information represented by  $\mathfrak{F}_1$  as the propositions the arguments relate to have been omitted.

A set  $S \subseteq U$  is *conflict-free* if and only if there are no two basic arguments  $\mathcal{A}, \mathcal{B} \in S$  with  $\mathcal{A} \rightarrow \mathcal{B}$ . A finite sequence  $[\mathcal{A}_1, \dots, \mathcal{A}_n]$  of basic arguments is conflict-free if and only if  $\{\mathcal{A}_1, \dots, \mathcal{A}_n\}$  is conflict-free. Basic arguments are used to form inference chains called *argument structures*.

**Definition 4 (Argument Structure).** Let  $S \subseteq U$  be a set of basic arguments and  $\mathcal{A} \in S$  a basic argument. An argument structure  $AS$  for  $\mathcal{A}$  with respect to  $S$  is a minimal (with respect to set inclusion) conflict-free sequence of basic arguments  $AS = [\mathcal{A} = \mathcal{A}_1, \dots, \mathcal{A}_n]$  with  $\{\mathcal{A}_2, \dots, \mathcal{A}_n\} \subseteq S$  such that for any  $\mathcal{A}_i \in AS$  and for any  $\alpha \in \text{supp}(\mathcal{A}_i)$  there is an  $\mathcal{A}_j \in AS$  with  $j > i$  and  $\text{cl}(\mathcal{A}_j) = \alpha$  (for  $1 \leq i, j \leq n$ ). Let  $\text{ArgStruct}_S(\mathcal{A})$  denote the set of argument structures for  $\mathcal{A}$  with respect to  $S$  and let  $\text{ArgStruct}_S = \bigcup_{\mathcal{A} \in S} \text{ArgStruct}_S(\mathcal{A})$  be the set of all argument structures with respect to  $S$ .

For an argument structure  $AS = [\mathcal{A}_1, \dots, \mathcal{A}_n]$  let  $\text{top}(AS) = \mathcal{A}_1$  denote the first basic argument in  $AS$ . The attack relation  $\rightarrow$  on basic arguments can be extended on argument structures by defining  $AS_1 \rightarrow AS_2$  if and only if there is an  $\mathcal{A} \in AS_2$  with  $\text{top}(AS_1) \rightarrow \mathcal{A}$  for two argument structures  $AS_1$  and  $AS_2$ . An argument structure  $AS_1$  *indirectly attacks* an argument structure  $AS_2$ , denoted by  $AS_1 \hookrightarrow AS_2$  if  $AS_1 \rightarrow \dots \rightarrow AS_2$  with an odd number of attacks.

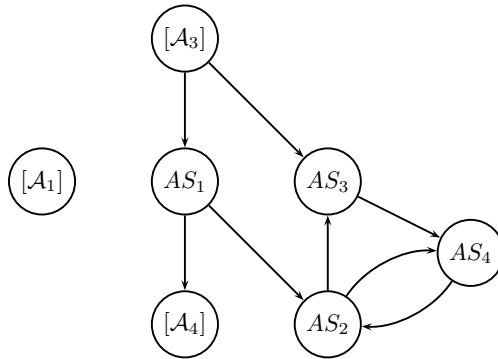
*Example 4.* We continue Example 3. In  $\mathfrak{F}_1$  the following sequences are argument structures

$$\begin{aligned} AS_1 &= [\mathcal{A}_2, \mathcal{A}_1] & AS_2 &= [\mathcal{A}_5, \mathcal{A}_4] \\ AS_3 &= [\mathcal{A}_6, \mathcal{A}_2, \mathcal{A}_1] & AS_4 &= [\mathcal{A}_7] \end{aligned}$$

Due to  $\mathcal{A}_2 \rightarrow \mathcal{A}_4$  it holds  $AS_1 \rightarrow AS_2$ . Similarly, it holds  $AS_2 \rightarrow AS_3$ ,  $AS_3 \rightarrow AS_4$ ,  $AS_2 \rightarrow AS_4$ ,  $AS_4 \rightarrow AS_2$ , and especially  $AS_1 \hookrightarrow AS_4$ .

Using the extended attack relation, a structured argumentation framework  $\mathfrak{F}$  induces an abstract argumentation framework  $AF_{\mathfrak{F}} = (\text{Arg}_{\mathfrak{F}}, \text{attacks}_{\mathfrak{F}})$  with  $\text{Arg}_{\mathfrak{F}} = \text{ArgStruct}_U$  and  $\text{attacks}_{\mathfrak{F}} = \{(AS_1, AS_2) \mid AS_1 \rightarrow AS_2\}$ . Let  $Sem$  denote one of the Dung-style semantics, cf. Subsection 2.1. Given a structured argumentation framework  $\mathfrak{F}$  and a semantics  $Sem$  the *output* of  $\mathfrak{F}$  denotes the set of all conclusions acceptable with the semantics  $Sem$  in the induced abstract argumentation framework  $AF_{\mathfrak{F}}$ , cf. [9]. More precisely, if  $E_1, \dots, E_n$  are the extensions of  $AF_{\mathfrak{F}}$  under  $Sem$ , then  $\text{Output}_{Sem}(\mathfrak{F}) = \{\alpha \in \text{Prop} \mid \forall i : \exists AS \in E_i : \text{cl}(\text{top}(AS)) = \alpha\}$ .

*Example 5.* A graphical representation of the induced abstract argumentation framework  $AF_{\mathfrak{F}_1}$  of  $\mathfrak{F}_1$  from Example 3 is depicted in Figure 3. Note that we abbreviated some argument structures by their names introduced in Example 4. The grounded extension  $E_G$  of  $AF_{\mathfrak{F}_1}$  computes to  $E_G = \{[A_1], [A_3], [A_4], AS_2\}$  and therefore  $\text{Output}_{\text{grounded}}(\mathfrak{F}_1) = \{a, c, d, e\}$ .



**Fig. 3.** The induced abstract argumentation framework of  $\mathfrak{F}_1$  from Example 3

Structured argumentation frameworks are a clear generalization of abstract argumentation frameworks as every abstract argumentation framework can be cast into a structured argumentation framework while retaining semantics.

**Definition 5 (Equivalent Structured Argument Framework).** *Let  $AF = (\text{Arg}, \text{attacks})$  be an abstract argumentation framework. For every argument  $A \in \text{Arg}$  introduce a new proposition  $\mathcal{A} \in \text{Prop}$ . The equivalent structured argumentation framework  $\mathfrak{F}_{AF} = (U, \rightarrow)$  to  $AF$  is defined as*

$$U = \{(\emptyset, \mathcal{A}) \mid \mathcal{A} \in \text{Arg}\}$$

$$\rightarrow = \{((\emptyset, \mathcal{A}), (\emptyset, \mathcal{B})) \mid (\mathcal{A}, \mathcal{B}) \in \text{attacks}\}$$

The following theorem states that structured argumentation frameworks are a clear generalization of abstract argumentation frameworks and can easily be verified.

**Theorem 1.** *Let  $AF$  be an abstract argumentation framework with extensions  $E_1, \dots, E_n$  under some semantics  $Sem$  and let  $E'_1, \dots, E'_m$  be the extensions*

of  $\text{AF}_{\mathfrak{F}_{\text{AF}}}$  under Sem. Then there is bijective function  $T : \{E_1, \dots, E_n\} \rightarrow \{E'_1, \dots, E'_m\}$  such that  $T(\{\mathcal{A}_1, \dots, \mathcal{A}_k\}) = \{(\emptyset, \mathcal{A}_1), \dots, (\emptyset, \mathcal{A}_k)\}$  for every  $E_i = \{\mathcal{A}_1, \dots, \mathcal{A}_k\}$ ,  $1 \leq i \leq n$ . In particular, it is  $n = m$ .

So far we have motivated the use of structured argumentation frameworks as a computational model for argumentation. We now turn to the setting of argumentation in dialogs. Usually in a multi-agent setting, the universal set of basic arguments  $U$  is unknown to all agents because of lack of expertise or just lack of knowledge. When considering a multi-agent setting, every agent may only have a partial view on  $U$  and the attack relation.

**Definition 6 (View).** A view  $V_{\mathfrak{F}}$  on a structured argumentation framework  $\mathfrak{F} = (U, \rightarrow)$  is a structured argumentation framework  $V_{\mathfrak{F}} = (U', \rightarrow')$  with  $U' \subseteq U$  and  $\rightarrow' = \{(\mathcal{A}_1, \mathcal{A}_2) \in \rightarrow \mid \mathcal{A}_1, \mathcal{A}_2 \in U'\}$ .

We will omit the subscript of  $V_{\mathfrak{F}}$  when the SAF  $\mathfrak{F}$  is clear from context. Definition 6 implies that, in general, games played on some structured argumentation framework are *incomplete* as not every possible move of an agent might be known by other agents. Nonetheless, when a move is played (i. e. an argument has been put forward) all agents agree on the attack relation. So with respect to the attack relation the information distributed among the agents is complete.

### 3 The Multi-agent Setting

The scenario we consider can be intuitively described as follows. At the beginning every agent has some view on the underlying SAF  $\mathfrak{F}$  and some preferences over the output of the argumentation. The common view considered by all agents as starting point is empty and the agents take turn by bringing up some basic arguments from their own view and incorporating them into the common view. When no agent can bring up more arguments the argumentation ends and an abstract argumentation framework is computed with respect to the final common view. Lastly, this abstract argumentation framework is used to compute the output of the argumentation given some predefined semantics. In the following, we formalize this intuition.

The multi-agent setting is divided into two parts, one describing the basic contents of the scenario, namely the underlying argumentation framework and the agents, and one describing the dynamic part of an evolving argumentation.

**Definition 7 (Structured Argumentation System).** A structured argumentation system (SAS)  $\Pi$  is a tuple  $\Pi = (\mathfrak{F}, \text{Ag})$  with a structured argumentation framework  $\mathfrak{F}$  and a set of agent identifiers  $\text{Ag}$ .

As a simplification we assume that the universal set of basic arguments  $U$  of  $\mathfrak{F}$  contains exactly the union of the basic arguments appearing in the views of the agents. Hence, any basic argument in  $U$  is known by at least one agent. This is not a restriction as an argument not appearing in any view cannot be used at all. In particular, we do not allow agents to “make up” arguments as in [21].

A SAS  $\Pi$  describes the functionality and the underlying language of an argumentation game. Dynamism is introduced by considering evolving *states* of  $\Pi$ . At any time the state of  $\Pi$  is determined by a current common view  $V^0$ , the views of each agent  $V^i$ , and the outcome of the argumentation.

**Definition 8 (State).** A state  $\Gamma^\Pi$  of  $\Pi = (\mathfrak{F}, Ag)$  with  $Ag = \{A_1, \dots, A_n\}$  is a tuple  $\Gamma^\Pi = (V^0, \{V^1, \dots, V^n\}, O)$  with views  $V^0, \dots, V^n$  on  $\mathfrak{F}$ , and a set  $O \subseteq \text{Prop}$ . Let  $\Delta^\Pi$  denote the set of all states of  $\Pi$ .

We will omit the superscripts  $\Pi$  when  $\Pi$  is clear from context. The final component  $O$  of a state  $\Gamma$  denotes the output of the argumentation if  $\Gamma$  is the final state. If the final state has not been reached yet, we set  $O = \text{nil}$ , where  $\text{nil}$  is a special identifier denoting no output. For a state  $\Gamma = (V^0, \{V^1, \dots, V^n\}, O)$  we denote  $V^i(\Gamma) = V^i$ , and  $O(\Gamma) = O$ . The initial state of a SAS  $\Pi$  is denoted by  $\Gamma_0^\Pi$  with  $O(\Gamma_0^\Pi) = \text{nil}$ . The state of a SAS  $\Pi$  evolves over time when agents bring up new basic arguments from their own views. The protocol of an argumentation game might restrict an agent to only bring up one basic argument at a time or all basic arguments he wants at once. We will elaborate some of these possible protocols in the next section. In the general case, if an agent has to take turn in an argumentation he does so by using its *selection function*. Given a common view of a SAF and an agent's own view a selection function selects a set of basic arguments of the agent's view to come up with. Let  $\mathfrak{P}(S)$  denote the power set of a set  $S$ .

**Definition 9 (Selection Function).** Let  $A_k$  be an agent identifier. A selection function  $\text{sel}^{A_k}$  for  $A_k$  is a function  $\text{sel}^{A_k} : \Delta \rightarrow \mathfrak{P}(U)$  such that  $\text{sel}^{A_k}(\Gamma) \subseteq (U^k \setminus U^0)$  for any  $\Gamma \in \Delta$  with  $V^k(\Gamma) = (U^k, \rightarrow^k)$  and  $V^0(\Gamma) = (U^0, \rightarrow^0)$ .

The condition  $\text{sel}^{A_k}(\Gamma) \subseteq (U^k \setminus U^0)$  ensures that the agent brings up new basic arguments that are not already part of the common view. Notice also that an agent may bring up no basic arguments at all via  $\text{sel}^{A_k}(\Gamma) = \emptyset$ . Intuitively spoken, a selection function implements the strategy of an agent in an argumentation in a game theoretical sense. As said before, in our framework the strategy of an agent only allows for hiding arguments but not for making up new arguments, cf. [17]. To our understanding this is not a drawback as arguments that could be made up by an agent could also be integrated in the agent's view from the beginning. From the perspective of knowledge representation this is a more adequate formalization as making up arguments requires the agent to have an understanding of rational inference chains (i.e. atomic arguments) to support their claim.

In game theory, the *performance* of an agent's strategy is evaluated by using the agent's preferences on the outcomes of a game. As in our framework the outcome of the argumentation game is determined by the output of the final common view of the underlying  $\mathfrak{F}$  the agent's utility is determined by its utility function which maps sets of propositions, i.e. possible outcomes, to natural numbers, thus describing a ranking on the output.

**Definition 10 (Utility Function).** An utility function  $\text{util}^A$  for an agent identifier  $A$  is a function  $\text{util}^A : \mathfrak{P}(\text{Prop}) \rightarrow \mathbb{N}$ .

An agent  $A$  with a utility function  $\text{util}^A$  prefers the outcome (i. e. the output)  $L_1 \subseteq \text{Prop}$  over  $L_2 \subseteq \text{Prop}$  if  $\text{util}^A(L_1) > \text{util}^A(L_2)$ . By taking a selection function and an utility function together we obtain the basic characteristics of an agent.

Other agents observe new basic arguments and integrate these in their own views respectively. As a convenience we abbreviate this operation as follows.

**Definition 11 (View Update).** *Let  $V = (U', \rightarrow')$  be a view on  $\mathfrak{F} = (U, \rightarrow)$  and  $\mathcal{A} \in U$  a basic argument. The view update of  $V$  with  $\mathcal{A}$  is a view  $V' = V \otimes \mathcal{A}$  on  $\mathfrak{F}$  with  $V' = (U'', \rightarrow'')$  defined as  $U'' = U' \cup \{\mathcal{A}\}$  and*

$$\rightarrow'' = \rightarrow' \cup \{(\mathcal{A}, \mathcal{B}) \in \rightarrow' \mid \mathcal{B} \in U'\} \cup \{(\mathcal{B}, \mathcal{A}) \in \rightarrow' \mid \mathcal{B} \in U'\}$$

Definition 11 suggests that agents are fully aware of attacks between known arguments. This means that when agents incorporate new basic arguments into their view, all attacks between this argument and arguments already known are incorporated as well. This assumption corresponds to the assumption of *perfect information* in e. g. [22]. For a set of basic arguments  $\mathfrak{A} = \{\mathcal{A}_1, \dots, \mathcal{A}_n\} \subseteq U$  we define  $V \otimes \mathfrak{A} = (\dots((V \otimes \mathcal{A}_1) \otimes \mathcal{A}_2) \otimes \dots) \otimes \mathcal{A}_n$ .

## 4 Argumentation Games

The type of argumentation game that agents play directly influences the strategies agents should use in order to obtain the best outcomes. In [20,21,17] Rahwan et. al. investigate mechanism design techniques [14] in order to determine suitable mechanisms, i. e. types of games, for abstract argumentation. For a special case of mechanism and a special type of agents they were able to identify this scenario as a strategy-proof game. As such, the best strategy for the agents is to be truthful about their views and bring up all arguments they know of. In their works, Rahwan et. al. focus on *direct mechanisms*, i. e., mechanisms where every agent reports his arguments at once without having the possibility to react on other agents' arguments. Restricting the attention on these simple games is not as limiting as it seems. Due to the *revelation principle*—a well-known result in mechanism design—if some *social choice function* can be implemented with some equilibrium by some mechanism it can also be implemented by a direct and truthful mechanism [20]. Roughly, this means that when designing a game one does not lose expressive power by just considering direct mechanisms. Instead one gains the additional advantage that agents have to be truthful. Nonetheless, there is some criticism on the revelation principle, especially when it comes to natural representation of games or computational issues. Implementing a game in a direct fashion might put a computational intractable task onto the evaluator of the game or create an exponential overhead in communication [10,2]. Furthermore, as a direct mechanism expects an agent to (truthfully) report its type, e. g. in our framework his arguments, confidentiality issues might be considered as well [6]. Hence, besides direct mechanisms we also investigate more natural

<sup>2</sup> Thanks to Iyad Rahwan for pointing that out to us.



forms of argumentation dialogs in the following. We obtain a similar result as in [20] of strategy-proofness for a special scenario of a SAS but we have a look on strategies for non-strategy-proof games as well, cf. Section 5.

In this section we give an overview on different settings for argumentation games. To this end we identify three key parameters as follows.

1. *Game protocol*: How do agents take turn and when does the game terminates?
2. *Awareness*: Does an agent have knowledge on the views of other agents?
3. *Agent types*: How are the preferences of an agent organized?

As discussed above, we assume for all scenarios that every action undertaken by any agent is recorded by all other agents and the agents agree on the structure of the attack relation.

### 4.1 Game Protocol

A protocol describes the extensional rules of an argumentation game and prescribes how agents take turns and which actions can be undertaken. More formally, we describe argumentation game protocols by means of state transition rules as in operational semantics [18] that transform one state of a SAS  $\Pi$  into a new one. Given a SAS  $\Pi$  and some initial state  $\Gamma_0^\Pi$  of  $\Pi$  the rules of a protocol  $P$  are applied to  $\Gamma_0^\Pi$  and its successor state until a *final state*  $\text{final}_P(\Gamma_0^\Pi)$  with  $O(\text{final}_P(\Gamma_0^\Pi)) \neq \text{nil}$  is reached. In this paper, we do not allow for probabilistic decision in the agents' strategies, so  $\text{final}_P(\Gamma_0^\Pi)$  is uniquely determined. An investigation on indeterministic strategies is part of future work. For an agent  $A$  its *gain* for  $\Gamma_0^\Pi$  and  $P$  is defined as  $\text{gain}_A^P(\Gamma_0^\Pi) = \text{util}_A(O(\text{final}_P(\Gamma_0^\Pi)))$ , i. e. the agent's utility for the outcome of the argumentation. In the following, let  $\Pi$  be a SAS with  $\Pi = (\mathfrak{F}, \{A_1, \dots, A_n\})$  and  $\Gamma$  a state.

### Direct Argumentation Mechanism

A direct argumentation mechanism [20] allows only one single step in the argumentation game. Every agent may put forward any set of basic arguments at once. After this, the mechanism terminates. This can be realized with the single state transition rule  $T_1^d$  defined as follows.

$$[T_1^d] \frac{\mathfrak{A} = \text{sel}^{A_1}(\Gamma) \cup \dots \cup \text{sel}^{A_n}(\Gamma)}{\Gamma \longrightarrow (V^{0'}, \{V^{1'}, \dots, V^{n'}\}, \text{Output}_{Sem}(V^{0'}))}$$

$$\text{with: } V^{i'} = V^i(\Gamma) \otimes \mathfrak{A} \quad (0 \leq i \leq n)$$

Obviously, the direct argumentation protocol  $P^d = \{T_1^d\}$  always terminates after one execution step.

### Synchronous Argumentation Mechanism

A generalization of the direct argumentation mechanism is the *synchronous argumentation mechanism*. There, every agent may bring up a set of basic arguments at the same time but the process is repeated until no agent wants to bring up any more basic arguments. There are two variants of this mechanism, one where agents are allowed to bring up new basic arguments even if they have not done so in a previous step, and one where agents cannot bring up any new basic arguments if they previously decided not to do so. We call the second variant a *rigid protocol*. When using a rigid protocol, agents have to carefully deliberate whether they choose to not bring forward any arguments, because they do not get any other chance to do so. In the following, we only consider the non-rigid variant. The non-rigid variant is realized with the following transition rules.

$$[T_1^s] \frac{\mathfrak{A} = \text{sel}^{A_1}(\Gamma) \cup \dots \cup \text{sel}^{A_n}(\Gamma) \text{ and } \mathfrak{A} \neq \emptyset}{\Gamma \longrightarrow (V^{0'}, \{V^{1'}, \dots, V^{n'}\}, \text{nil})}$$

with:  $V^{i'} = V^i(\Gamma) \otimes \mathfrak{A} \quad (0 \leq i \leq n)$

$$[T_2^s] \frac{\text{sel}^{A_1}(\Gamma) \cup \dots \cup \text{sel}^{A_n}(\Gamma) = \emptyset}{\Gamma \longrightarrow (\cdot, \cdot, \text{Output}_{Sem}(V^0(\Gamma)))}$$

The *synchronous argumentation protocol*  $P^s = \{T_1^s, T_2^s\}$  also clearly terminates after a finite number of steps, because the number of basic arguments is finite. Note, that in the synchronous and the direct argumentation mechanism the assumption of *perfect information* is restrained due to the simultaneous moves of the agents. Therefore, the selection of arguments to put forward can only depend on the moves of other agents from the previous steps but not on those in the current step.

### Dialectical Argumentation Mechanism

In natural dialogues agents usually alternately take turns when bringing up arguments. In general, this can be realized by a *dialectical argumentation mechanism* where we assume some order of the agents and basic arguments can be brought up with respect to this order. As for the synchronous argumentation mechanism two variants are possible with respect to rigidity of the protocol. Anyway, the protocol needs some extra meta information for the states to select the next agent appropriately and we have to ensure that the protocol terminates if no agent wants to bring up new arguments. To this end we introduce some meta information  $M = (k_1, k_2) \in \mathbb{N}^2$  such that  $k_1$  is the index of the agent that last took turn and  $k_2$  counts the number of agents that skipped bringing up new basic arguments since the last one that did. For an initial state  $\Gamma_0^{\text{II}}$  we set  $M = (0, 0)$ . Then this protocol is realized by the following transition rules.

$$[T_1^t] \frac{k_2 < n \text{ and } \mathfrak{A} = \text{sel}^{A_{k_1}'}(\Gamma)}{\Gamma \longrightarrow (V^{0'}, \{V^{1'}, \dots, V^{n'}\}, \text{nil})} \\ M = (k_1, k_2) \longrightarrow M' = (k_1', k_2')$$

$$\text{with: } V^{i'} = V^i(\Gamma) \otimes \mathfrak{A} \quad (0 \leq i \leq n) \\ k_1' = (k_1 \bmod n) + 1 \\ k_2' = \begin{cases} 0 & \text{if } \mathfrak{A} \neq \emptyset \\ k_2 + 1 & \text{otherwise} \end{cases}$$

$$[T_2^t] \frac{k_2 = n}{\Gamma \longrightarrow (\cdot, \cdot, \text{Output}_{Sem}(V^0))} \\ M = (k_1, k_2) \longrightarrow M$$

As for the synchronous argumentation protocol the termination of the *dialectical argumentation protocol*  $P^t = \{T_1^t, T_2^t\}$  is ensured due to the finiteness of the universal set of basic arguments  $U$ .

Notice that a variant of the rigid version of the dialectical argumentation mechanism has been previously employed for an argumentation game in [22].

The general protocols described above allow an agent to bring forward an arbitrary number of arguments at any step. For the synchronous and dialectical mechanisms a restricted variant would be allow an agent to bring forward only a single argument at any step. We call such a protocol an *atomic-step* protocol. More formally, an atomic-step protocol  $P$  can only be applied to a SAS  $(\mathfrak{F}, Ag)$  if for all  $A \in Ag$  it is  $|\text{sel}_A(S, \mathfrak{F})| \leq 1$  for any  $S \in \mathfrak{P}(U)$  and every  $\mathfrak{F}$ . Together with the option of rigidness we obtain each four variants of the synchronous and dialectical mechanisms. Notice also, that we do not restrict the agents to follow some dialectical structure such as always replying to the last argument brought forward. The above protocols can be refined in order to implement such restrictions but this is outside the scope of this paper. Assuming a fair implementation of the protocols they fulfill most of the desiderata expected for argumentation protocols proposed in [15] such as *separation of syntax and semantics* and *discouragement of disruption*.

## 4.2 Awareness

Our definition of selection functions (Definition 9) is quite general as it takes the whole state of the system into account when determining the basic argument that should be brought forward. In particular, a selection function might be heavily influenced by the views of other agents. Usually, an agent does not have complete and accurate knowledge on the subjective views of other agents. One extreme is that an agent has *no awareness* of other agents views. More formally, a selection function  $\text{sel}^{A_k}$  of an agent  $A_k \in Ag$  is *ignorant* if for all  $\Gamma_1, \Gamma_2 \in \Delta$  it holds: If  $V_0(\Gamma_1) = V_0(\Gamma_2)$  and  $V^k(\Gamma_1) = V^k(\Gamma_2)$ , then it is  $\text{sel}^{A_k}(\Gamma_1) = \text{sel}^{A_k}(\Gamma_2)$ .

This means that the decision of agent  $A_k$  is at any time only dependent on the agent's own view and the common view.

Usually, an agent has some subjective beliefs about the views of other agents. Let  $\text{Bel}_{A_k}(A_j, \Gamma)$  the subjective belief of agent  $A_k$  on the view of agent  $A_j$  in state  $\Gamma$ , i.e.  $\text{Bel}_{A_k}(A_j, \Gamma)$  is itself a view. Then, a selection function  $\text{sel}^{A_k}$  of  $A_k$  is *belief-based* if for all  $\Gamma_1, \Gamma_2 \in \Delta$  it holds: If  $V^0(\Gamma_1) = V^0(\Gamma_2)$  and  $V^k(\Gamma_1) = V^k(\Gamma_2)$  and for all  $j \neq k$  it is  $\text{Bel}_{A_k}(A_j, \Gamma_1) = \text{Bel}_{A_k}(A_j, \Gamma_2)$ , then it is  $\text{sel}^{A_k}(\Gamma_1) = \text{sel}^{A_k}(\Gamma_2)$ . An agent  $A_k$  has *full awareness* if his selection function  $\text{sel}^{A_k}$  is belief-based and  $\text{Bel}_{A_k}(A_j, \Gamma) = V^j(\Gamma)$  for every state  $\Gamma \in \Delta$  and  $j \neq k$ .

In between no awareness and full awareness there is a wide range of incomplete and uncertain awareness of other agents' views, but we will not discuss this topic in the current paper.

### 4.3 Agent Types

Under the term *agent type* we understand in this paper the way the preferences of the agent are organized. The main reason for arguing with other agents is to persuade other agents or to prove some statement. This goal is represented by the agent's utility function which ranks the possible outcomes of the argumentation. In the following we identify some simple utility functions.

The most simple attitude of an agent towards the outcome of an argumentation is the desire to prove a single proposition, no matter what else is proven.

**Definition 12 (Indicator Utility Function).** *Let  $\alpha \in \text{Prop}$ . The utility function  $\text{util}_\alpha$  is called an indicator utility function for  $\alpha$  if for any  $L \subseteq \text{Prop}$  it is  $\text{util}_\alpha(L) = 1$  if  $\alpha \in L$  and  $\text{util}_\alpha(L) = 0$  otherwise.*

The choice of 0 and 1 as the only values for the indicator utility function is arbitrary. Any utility function  $\text{util}$  with  $\text{util}(L) = k$  and  $\text{util}(L') = l$  for any  $L, L' \subseteq \text{Prop}$  with  $\alpha \in L$  and  $\alpha \notin L'$  for some  $\alpha$  can be normalized to an indicator utility function if  $k > l$ . Note that the definition of indicator utility functions resembles the rationale behind *focal arguments* in [20]. Because of this, if  $\text{util}_\alpha$  is the utility function of an agent  $A$  we call  $\alpha$  the *focal element* of  $A$ .

The definition of an indicator function can be extended to comprehend for multiple focal elements as follows.

**Definition 13 (Multiple Indicator Utility Function).** *The utility function  $\text{util}_S$  is called a multiple indicator utility function for  $S \subseteq \text{Prop}$  if for any  $L \subseteq \text{Prop}$  it is  $\text{util}_S(L) = 0$  if  $S \not\subseteq L$  and  $\text{util}_S(L) = 1$  if  $S \subseteq L$ .*

Notice that it holds  $\text{util}_{\{\alpha\}} = \text{util}_\alpha$ . This general definition does not demand that  $S$  has to be "consistent", i.e. there may be argument structures  $AS_1$  resp.  $AS_2$  for some  $\alpha \in \text{Prop}$  resp.  $\alpha' \in \text{Prop}$  such that  $AS_1 \rightarrow AS_2$ . Another variant of an agent's preferences can be characterized by a counting utility function which is similar in spirit to the notion of *acceptability maximising preferences* in [20].

**Definition 14 (Counting Utility Function).** *Let  $S \subseteq \text{Prop}$ . The utility function  $\text{util}_S^\#$  is called a counting utility function for  $S$  if for any  $L \subseteq \text{Prop}$  it is  $\text{util}_S^\#(L) = |L \cap S|$ .*

Notice that it holds  $\text{util}_{\{\alpha\}}^{\#} = \text{util}_{\alpha}$ . The difference between a counting utility function and a multiple indicator utility function is that for a multiple indicator utility function all focal elements have to be in the output of an argumentation in order to yield a better utility than zero. An agent with a counting utility function tries to prove as many of his focal elements as possible.

In general, there has to be no direct relationship between an agent’s view and his utility function. For example, an agent with an indicator utility function  $\text{util}_{\alpha}$  may have no basic argument for  $\alpha$  in his own view or, more drastically, his view can give reasons to not believe in  $\alpha$ . A special form of views are *subjective views* in which an agent’s utility function is consistent with its own view.

**Definition 15 (Subjective View).** *Let  $V$  be a view on  $\mathfrak{F}$ .  $V$  is a subjective view on  $\mathfrak{F}$  with respect to a utility function  $\text{util}$  if and only if  $\text{util}(\text{Output}_{\text{Sem}}(V))$  is a maximum of  $\text{util}$ .*

Furthermore, a view  $V = (U', \rightarrow')$  is *globally consistent with respect to a SAF  $\mathfrak{F}$*  if there are no two argument structures  $AS_1, AS_2$  in  $\mathfrak{F}$  such that  $AS_1 \leftrightarrow AS_2$  and  $AS_1 \cap U' \neq \emptyset$  and  $AS_2 \cap U' \neq \emptyset$ . This means that no two basic arguments in  $V$  can be used to construct argument structures that are, in any way, inconsistent to one another.

Figure 4 summarizes the different game parameters we investigate in this paper, ordered by their “complexity”. Distance from the origin indicates a more demanding setting with respect to the complexity of the strategy for argument selection.

## 5 Strategies for Selecting Arguments

In the following, we investigate some strategies for argument selection in different argumentation games as defined in the previous section. The most simple selection function one can think of is the one that just reports all basic arguments of the agent’s view. Let  $A_k \in Ag$  be an agent identifier and  $\Gamma$  a state. Then the *truthful selection function*  $\text{sel}_{\top}^{A_k}$  is defined as  $\text{sel}_{\top}^{A_k}(\Gamma) = U^k \setminus U^0$  with  $V^k(\Gamma) = (U^k, \rightarrow^k)$  and  $V^0(\Gamma) = (U^0, \rightarrow^0)$ . In other words, the selection function  $\text{sel}_{\top}^{A_k}$  always returns all basic arguments of an agent’s view that aren’t already present in the common view of the SAS. As being truthful does not demand for strategic decisions the function is the same for direct, synchronous, and dialectical argumentation protocols. For an atomic-step protocol the truthful selection function can be serialized, i. e., a serialized variant would select an arbitrary basic argument each turn until all arguments have been brought forward.

In general, we are interested in finding selection functions that maximize an agent’s gain in an argumentation game. Here, an *argumentation game*  $AG$  is defined as a tuple  $AG = (\Pi, P)$  with a SAS  $\Pi$  and a protocol  $P$ . The strongest concept of a selection function that maximizes utility is that of a *dominant strategy*. Let  $\Pi$  be a SAS and let  $\Pi'$  be the same as  $\Pi$  except possibly different selection functions of the agents. Then the selection function  $\text{sel}^{A_k}$  of agent  $A_k$  is a *dominant selection function* if for any such  $\Pi'$  it is  $\text{gain}_{A_k}^P(\Gamma_0^{\Pi}) \geq \text{gain}_{A_k}^P(\Gamma_0^{\Pi'})$ .

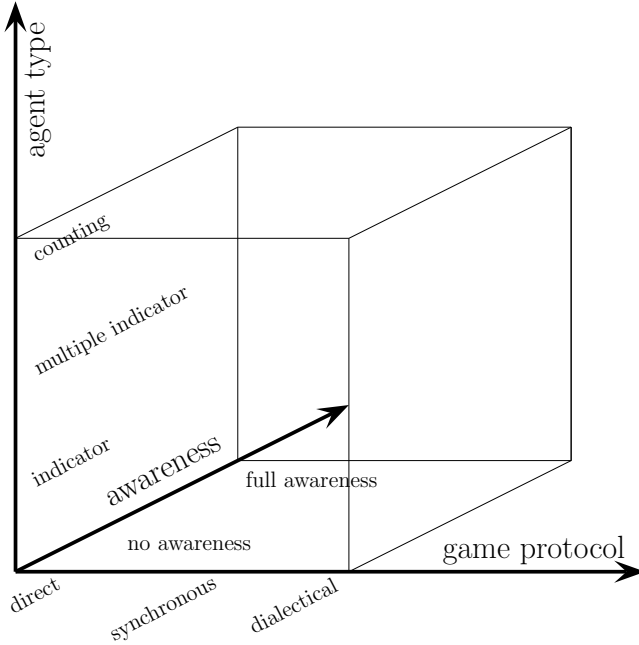


Fig. 4. Complexity of game parameters

This means, regardless of how the other agents select their arguments, the selection function  $\text{sel}^{A_k}$  maximizes the gain of agent  $A_k$ .<sup>3</sup> The truthful strategy is of special interest in game theory as it is the dominant strategy for *strategy-proof* games. Therefore, given a strategy-proof argumentation game it is the best choice for each agent to truthfully report all their basic arguments. In [20] Rahwan et. al. identified a special type of direct argumentation game as strategy-proof. We can restate and extend their result in our framework as follows.

**Theorem 2.** *Let  $\Pi = (\mathfrak{F}, Ag)$  be a SAS. If the initial view  $V^i(\Gamma_0^\Pi)$  of each agent  $A_i \in Ag$  is subjective and globally consistent with respect to  $\mathfrak{F}$  and the utility function  $\text{util}^{A_i}$  of each agent  $A_i$  is a counting utility function, then  $(\Pi, P^d)$  is strategy-proof.*

Observe that the above statement is independent of the actual chosen semantics due to the skeptical definition of **Output**. Theorem 2 states that the dominant strategy for subjective and globally consistent views is to use the truthful selection function  $\text{sel}_\top$ . It is a clear extension of Theorem 32 stated in [20] as our underlying argumentation framework is a structured argumentation framework. The statement of Theorem 2 easily extends to indicator utility functions, multiple indicator utility functions as well as synchronous and dialectical argumentation protocols (the latter because of the *revelation principle*, see above). However, the condition of a globally consistent view is hard to check for an agent

<sup>3</sup> Notice that agent  $A_k$  may have the same selection function  $\text{sel}^{A_k}$  in  $\Pi$  and  $\Pi'$ .

who has no idea of the structure of the underlying framework  $\mathfrak{F}$ . Given a basic argument  $\mathcal{A}$  in his view he may not know if  $\mathcal{A}$  can be used to construct an argument structure against one of his “own” arguments. Due to this observation Theorem 2 is only applicable for an agent if the global consistency is assured by a trustworthy third party or if the agents have full awareness of the other agent’s views and thus can verify the global consistency by themselves. Otherwise an agent cannot know if the best strategy is to be truthful.

In general, full awareness is not a realistic assumption in argumentation. When agents cannot verify the global consistency of their view, some strategic deliberations are mandatory as the following example shows.

*Example 6.* Consider the following SAF  $\mathfrak{F}_2 = (U, \rightarrow)$ .

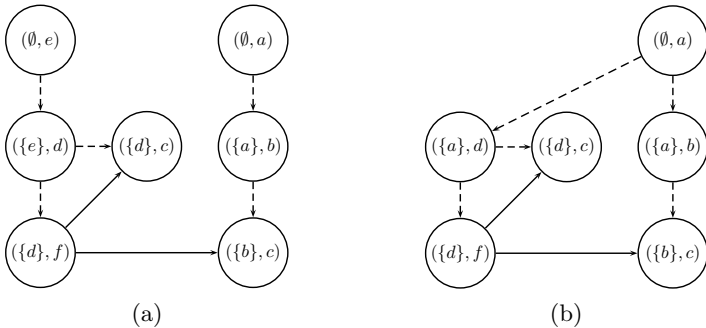
$$\begin{aligned}
 U = & \{ (\emptyset, a), (\{a\}, b), (\{b\}, c), (\emptyset, e), \\
 & \{ \{e\}, d \}, \{ \{d\}, f \}, \{ \{d\}, c \} \} \\
 \rightarrow = & \{ (\{ \{d\}, f \}, \{ \{d\}, c \}), (\{ \{d\}, f \}, \{ \{b\}, c \}) \}
 \end{aligned}$$

An overview of  $\mathfrak{F}_2$  is given in Figure 5(a). Let  $\Pi = (\mathfrak{F}_2, \{A_1, A_2\})$  be a SAS and the initial state  $\Gamma_0^\Pi = (\emptyset, \{V^1, V^2\}, \text{nil})$  of  $\Pi$  be given as follows.

$$V^1 = (U \setminus \{ \{ \{d\}, f \} \}, \cdot) \qquad V^2 = (\{ \{ \{d\}, f \} \}, \emptyset)$$

The attack relation of  $V^1$  is omitted but can be determined via Definition 6. Note that view  $V^1$  is subjective but not globally consistent. Imagine  $A_1$  wants to prove  $c$ , i. e., the utility function of  $A_k$  is  $\text{util}_c$ . Note that there are two argument structures in  $\mathfrak{F}_2$  to prove  $c$  while one of them ( $(\{ \{d\}, c \}, \{ \{e\}, d \}, (\emptyset, e))$ ) enables  $A_2$  to bring up an attacker, namely  $(\{ \{d\}, f \}, \{ \{e\}, d \}, (\emptyset, e))$ . From a self-interested point of view  $A_1$  should only bring forward the arguments that do not allow  $A_2$  to counterargue.

In the following, we develop some simple strategies for argument selection that generalize the truthful strategy in scenarios where the agent may not have a



**Fig. 5.** The structured argumentation frameworks (a)  $\mathfrak{F}_2$  from Example 6 and (b)  $\mathfrak{F}_3$  from Example 8

globally consistent view and that are more cautious in bringing forward arguments. In order to ensure that an agent brings forward only the arguments that are not harmful for proving his focal elements, we define the *attack set* as follows.

**Definition 16 (Attack Set).** *Let  $\mathfrak{F} = (U, \rightarrow)$  be a SAF and  $\alpha \in \text{Prop}$ . The attack set  $\text{AttackSet}_{\mathfrak{F}}(\alpha)$  of  $\alpha$  in  $\mathfrak{F}$  is defined as*

$$\text{AttackSet}_{\mathfrak{F}}(\alpha) = \{ \mathcal{A} \in U \mid \exists AS_1, AS_2 \in \text{ArgStruct}_U : \mathcal{A} \in AS_1 \wedge \text{cl}(\text{top}(AS_2)) = \alpha \wedge AS_1 \leftrightarrow AS_2 \}$$

Intuitively, the set  $\text{AttackSet}_{\mathfrak{F}}(\alpha)$  contains all arguments that can be harmful to  $\alpha$  in any way. For example, for any argument  $\mathcal{A}$  with claim  $\alpha$ , the set  $\text{AttackSet}_{\mathfrak{F}}(\alpha)$  contains all attackers on  $\mathcal{A}$ . More generally,  $\text{AttackSet}_{\mathfrak{F}}(\alpha)$  contains every argument that belongs to an argument structure that indirectly attacks an argument structure for  $\alpha$ . Using attack sets we can define a simple strategy that brings only forward arguments that cannot be harmful in any way.

**Definition 17 (Overcautious Selection Function).** *Let  $\alpha \in \text{Prop}$  and  $A_k$  an agent identifier. Let  $s_{\alpha, A_k}^{\text{oc}}$  be the selection function defined as*

$$s_{\alpha, A_k}^{\text{oc}}(\Gamma) = \text{sel}_{\top}^{A_k}(\Gamma) \setminus \text{AttackSet}_{V^k(\Gamma)}(\alpha)$$

for every state  $\Gamma$ . The function  $s_{\alpha}^{\text{oc}}$  is called the overcautious selection function for  $\alpha$ .

Although the overcautious strategy is more careful in bringing forward arguments one should note that the determination of  $\text{AttackSet}_{V^k(\Gamma)}(\alpha)$  depends on the current view of the agent and might not be complete. The overcautious selection function can be extended to a belief-based selection function by incorporating the beliefs of  $A_k$  on the views of the other agents, into the determination of  $\text{AttackSet}_{V^k(\Gamma)}(\alpha)$ . However, we will not formalize this in the current paper.

*Example 7.* We continue Example 6 but suppose  $\text{sel}^{A_1} = s_{c, A_1}^{\text{oc}}$ . Here,  $A_1$  will not bring forward arguments  $(\emptyset, e)$  and  $(\{e\}, d)$  as they all belong to  $\text{AttackSet}_{V_1}(c)$ . Note that this strategy is independent of the strategy of any other agent.

Although the overcautious strategy is a very simple strategy for argument selection it is the dominant strategy in a simple class of argumentation games. If an agent has a complete view, i. e., he knows of every argument in the system, but has no awareness on the other agents beliefs, then its best choice is to avoid bringing forward possibly harmful arguments.

**Theorem 3.** *Let  $\Pi = (\mathfrak{F}, Ag)$  be a SAS. For an agent  $A_i \in Ag$ , if  $V_i(\Gamma_0^{\Pi}) = \mathfrak{F}$  and  $A_i$  has no awareness then the overcautious selection function is a dominant strategy for  $A_i$  in  $(\Pi, P^d)$ .*

The limitations of this simple strategy are reached very quickly as the following small modification of Example 6 shows.



*Example 8.* Consider the following SAF  $\mathfrak{F}_3 = (U, \rightarrow)$ , cf. Figure 5(b).

$$\begin{aligned}
 U &= \{ (\emptyset, a), (\{a\}, b), (\{b\}, c), (\{a\}, d), (\{d\}, f), (\{d\}, c) \} \\
 \rightarrow &= \{ ((\{d\}, f), (\{d\}, c)), ((\{d\}, c), (\{d\}, f)) \}
 \end{aligned}$$

Let  $\Pi = (\mathfrak{F}_3, \{A_1, A_2\})$  be a SAS and  $\Gamma_0^\Pi = (\emptyset, \{V_1, V_2\}, \text{nil})$  the initial state of  $\Pi$  with  $V^1 = \mathfrak{F}_3$  and  $V^2 = (U \setminus \{(\{a\}, d)\}, \cdot)$ . Suppose  $\text{util}_{A_1} = \text{util}_c$  and  $\text{sel}_{A_1} = \text{s}_{c, A_k}^{\text{oc}}$ . Here,  $A_1$  will never bring forward argument  $(\emptyset, a)$  as  $(\emptyset, a) \in \text{AttackSet}_{V_1}(c)$ . As a consequence,  $A_1$  will never be able to proof any argument for  $c$ .

As Example 8 showed it is advisable to bring forward arguments that on the one side may be harmful to one own’s desires but on the other side necessary to actually reach the desires. So we refine the overcautious strategy by allowing the agent to bring forward arguments that are inherently necessary for constructing an argument structure for his focal element.

**Definition 18 (Necessary Arguments).** Let  $\mathfrak{F} = (U, \rightarrow)$  be a SAF and  $\alpha \in \text{Prop}$ . The set of necessary arguments  $\text{NecArg}_{\mathfrak{F}}(\alpha)$  for  $\alpha$  in  $\mathfrak{F}$  is defined as

$$\text{NecArg}_{\mathfrak{F}}(\alpha) = \bigcap_{A \in U, \text{cl}(A) = \alpha, AS \in \text{ArgStruct}_V(A)} AS$$

**Definition 19 (Cautious Selection Function).** Let  $\alpha \in \text{Prop}$ ,  $A_k$  and agent with a view  $V$  and  $\text{s}_{\alpha, A_k}^c$  be the selection function defined as

$$\text{s}_{\alpha, A_k}^c(\Gamma) = \text{sel}_{\Gamma}^{A_k}(\Gamma) \setminus (\text{AttackSet}_V(\alpha) \setminus \text{NecArg}_V(\alpha))$$

$\text{s}_{\alpha, A_k}^c$  is called the cautious selection function for  $\alpha$ .

*Example 9.* We continue Example 8 but suppose  $\text{util}_{A_1} = \text{util}_c$  and  $\text{sel}_{A_1} = \text{s}_{c, A_k}^c$ . Here,  $A_1$  will bring forward argument  $(\emptyset, a)$  because it is inherently necessary to construct any argument structure for  $c$ .

The cautious strategy performs well in the above example and can be seen as a lower bound for direct argumentation protocols, i.e. the cautious strategy returns as few arguments as necessary.

## 6 Summary and Future Work

In this work we have introduced structured argumentation frameworks, a formalism that extends Dung’s abstract argumentation frameworks [11] and are a slightly modified variant of dynamic argumentation frameworks [23]. We have used structured argumentation frameworks for defining a multi-agent setting that contains two elements: one describing the basic contents of the scenario, i.e. the underlying argumentation framework and the set of agents; and a second element that describes the dynamic part of an evolving argumentation and

determines how the state of the multi-agent system evolves in time. In our framework every agent has its own view on the underlying argumentation framework and its own preferences over the output of the argumentation process. We proposed a first attempt to characterize argumentation games by means of the used game protocol, the awareness of the agents on other agents beliefs, and the structure of the preferences of the agents. We used structured argumentation systems to model argumentation among a group of agents. We have also presented some properties for the proposed framework and protocols.

For future work we plan to investigate the concept of strategies based on (uncertain) beliefs of other agents' views. In natural dialogues strategic argumentation is all about what an agents expects of his opponents beliefs and attitudes as even weak arguments can win an argumentation if the opponent has no counter-argument available. Especially when considering dialectical argumentation the possibility to learn from an agent's previous moves and thus building up beliefs on the other agent's view incrementally might bring advantage in the ongoing argumentation.

## References

1. Atkinson, K., Bench-capon, T., Mcburney, P.: A Dialogue Game protocol for Multi-Agent Argument over Proposals for Action. *Journal of Autonomous Agents and Multi-Agent Systems*, 149–161 (2004)
2. Baroni, P., Cerutti, F., Giacomin, M., Guida, G.: Encompassing Attacks to Attacks in Abstract Argumentation Frameworks. In: *Proceedings of the 10th European Conference on Symbolic and Quantitative Approaches to Reasoning with Uncertainty*, pp. 83–94 (2009)
3. Baroni, P., Giacomin, M.: Skepticism relations for comparing argumentation semantics. *International Journal of Approximate Reasoning* 50(6), 854–866 (2009)
4. Bench-Capon, T.J.M., Dunne, P.E.: *Argumentation in Artificial Intelligence*. *Artificial Intelligence* 171, 619–641 (2007)
5. Besnard, P., Hunter, A.: *Elements of Argumentation*. The MIT Press, Cambridge (2008)
6. Biskup, J., Kern-Isberner, G., Thimm, M.: Towards enforcement of confidentiality in agent interactions. In: Pagnucco, M., Thielscher, M. (eds.) *Proceedings of the 12th International Workshop on Non-Monotonic Reasoning (NMR 2008)*, pp. 104–112. University of New South Wales, Sydney (2008) Technical Report No. UNSW-CSE-TR-0819 (September 2008)
7. Black, E., Hunter, A.: *An Inquiry Dialogue System*. *Autonomous Agents and Multi-Agent Systems* (2009)
8. Caminada, M.: Semi-stable semantics. In: *Proceedings of the First Conference on Computational Models of Argument (COMMA 2006)*, pp. 121–130 (2006)
9. Caminada, M., Amgoud, L.: An Axiomatic Account of Formal Argumentation. In: *20th National Conference on Artificial Intelligence (AAAI 2005)*, Pittsburgh, pp. 608–613 (July 2005)
10. Conitzer, V., Sandholm, T.: Computational criticisms of the revelation principle. In: *Proceedings of the 5th ACM Conference on Electronic Commerce (EC 2004)*, New York, NY, USA, pp. 262–263 (2004)

11. Dung, P.M.: On the Acceptability of Arguments and its Fundamental Role in Nonmonotonic Reasoning, Logic Programming and n-Person Games. *Artificial Intelligence* 77(2), 321–358 (1995)
12. Garía, A., Simari, G.R.: Defeasible Logic Programming: An Argumentative Approach. *Theory and Practice of Logic Progr.* 4(1-2), 95–138 (2004)
13. Karunatillake, N.C., Jennings, N.R., Rahwan, I., McBurney, P.: Dialogue Games that Agents Play within a Society. *Artificial Intelligence* 173, 935–981 (2009)
14. Mas-Colell, A., Whinston, M.D., Green, J.R.: *Microeconomic Theory*. Oxford University Press, Oxford (1995)
15. McBurney, P., Parsons, S., Wooldridge, M.: Desiderata for Agent Argumentation Protocols. In: *Proceedings of the First International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2002)*, Bologna, Italy (July 2002)
16. Modgil, S.: Reasoning about Preferences in Argumentation Frameworks. *Artificial Intelligence* 173(9-10), 901–934 (2009)
17. Pan, S., Larson, K., Rahwan, I.: Argumentation mechanism design for preferred semantics. In: *In Proceedings of the 3rd International Conference on Computational Models of Argument (COMMA 2010)* (September 2010)
18. Plotkin, G.D.: *A Structural Approach to Operational Semantics*. Technical report, Department of Computer Science, Aarhus University, Aarhus, Denmark (1981)
19. Prakken, H., Vreeswijk, G.: Logical Systems for Defeasible Argumentation. In: Gabbay, D., Guenther, F. (eds.) *Handbook of Philosophical Logic*, vol. 4, pp. 219–318. Kluwer, Dordrecht (2002)
20. Rahwan, I., Larson, K.: Mechanism Design for Abstract Argumentation. In: *Proceedings of 7th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, Estoril, Portugal, pp. 1031–1038 (2008)
21. Rahwan, I., Larson, K., Tohmé, F.: A Characterisation of Strategy-Proofness for Grounded Argumentation Semantics. In: *Proceedings of the 21st International Joint Conference on Artificial Intelligence (IJCAI)*, Pasadena, California, USA (2009)
22. Riveret, R., Prakken, H., Rotolo, A., Sartor, G.: Heuristics in Argumentation: A Game-Theoretical Investigation. In: *Proceedings of COMMA 2008*, pp. 324–335 (2008)
23. Rotstein, N.D., Maguillansky, M.O., García, A.J., Simari, G.R.: An Abstract Argumentation Framework for Handling Dynamics. In: *Proceedings of NMR 2008*, pp. 131–139 (2008)
24. Thimm, M., García, A.J.: Classification and Strategical Issues of Argumentation Games on Structured Argumentation Frameworks. In: *Proceedings of the Ninth International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS 2010)* (May 2010)
25. Thimm, M., Garcia, A.J., Kern-Isberner, G., Simari, G.R.: Using Collaborations for Distributed Argumentation with Defeasible Logic Programming. In: Pagnucco, M., Thielscher, M. (eds.) *Proceedings of the 12th International Workshop on Non-Monotonic Reasoning (NMR 2008)*, Sydney, Australia, pp. 179–188. University of New South Wales (September 2008) Technical Report No. UNSW-CSE-TR-0819

# Preference-Based Argumentation Capturing Prioritized Logic Programming

Toshiko Wakaki

Shibaura Institute of Technology  
307 Fukasaku, Minuma-ku, Saitama-city, Saitama, 337-8570 Japan  
twakaki@sic.shibaura-it.ac.jp

**Abstract.** First, we present a novel approach to an abstract preference-based argumentation framework (an abstract *PAF*), which generalizes Dung's abstract argumentation framework (*AF*) to deal with additional preferences over a set of arguments. In our formalism, the semantics of such a *PAF* is given as  $\mathcal{P}$ -extensions that are selected from extensions of acceptability semantics by taking into account such preferences. Second, using a prioritized logic program (PLP) capable of representing priority information along with integrity constraints, the proposed method defines the non-abstract preference-based argumentation framework (the non-abstract *PAF*) translated from a PLP, whose semantics is also given by  $\mathcal{P}$ -extensions instantiating those of an abstract one. Finally we show the interesting result that,  $\mathcal{P}$ -extensions of such a non-abstract *PAF* under stable semantics capture *preferred* answer sets of a PLP, which ensures the advantages as well as the correctness of our approach.

## 1 Introduction

In the research field of argumentation, Dung's frameworks of abstract argumentation [11] have gained wide acceptance and are the basis for the implementation of concrete formalisms. In his paper [11], Dung showed that argumentation can be viewed as a special form of logic programming with negation as failure and gives a series of theorems that relate semantics of logic programs and semantics of argumentation frameworks. And as its application, there have been a number of proposals for negotiation between multiagents that make use of argumentation.

Recently, several approaches to generalize Dung's theory have been proposed in order to handle additional information such as preferences as well as constraints which a negotiating agent generally has knowledge of, because preferences are useful to solve conflicts between arguments and constraints are needed to eliminate extensions not satisfying the required conditions. With respect to handling preferences, quite recently, Amgoud and Vesic [2] pointed out that, there is the critical problem such that extensions are not conflict-free w.r.t. the attack relation for existing preference-based argumentation frameworks such as Amgoud and Cayrol's approach [1]. Conflict-freeness for extensions is important since it ensures sound results. Hence they [2] proposed a new abstract

preference-based argumentation framework whose semantics ensures requirements of *conflict-freeness* along with *generalization* to recover Dung's acceptability semantics in the case where preferences are not available.

On the other hand, with respect to formalisms for integrating logic programming and argumentation, Dung [10] showed that answer set semantics [13] of an extended logic program (an ELP, for short) is captured by stable semantics of Dung's argumentation framework, whereas Prakken and Sartor [15] introduced an argument-based formalism for extended logic programming with defeasible properties, which instantiated Dung's grounded semantics if it is restricted to static priorities. With respect especially to logic programming based on answer set semantics [12,13], a significant amount of studies have been done such as Brewka and Eiter's preferred answer sets for extended logic programs [5], Sakama and Inoue's prioritized logic programming [19], Delgrande and Schaub's ordered logic programs [9] and so on. However, as far as we know, few works have been achieved respecting the semantic relation between such logic programming capable of handling preferences and preference-based argumentation which generalizes Dung's argumentation framework [11].

Under such circumstances, first, we present a new approach of an abstract preference-based argumentation framework (an abstract *PAF*, for short), which generalizes Dung's abstract argumentation framework to deal with additional preferences with meeting requirements of conflict-freeness and generalization. In our formalism, the semantics of such a *PAF* is given as  $\mathcal{P}$ -extensions that are selected from extensions of Dung's acceptability semantics by taking into account such preferences.

Second, since Sakama and Inoue's formalism of a prioritized logic program (PLP) is capable of representing priority information along with integrity constraints in a nonmonotonic logic program, we use such PLP as the underlying logic to construct a non-abstract *PAF* instantiating an abstract *PAF*. That is, the proposed method defines a non-abstract preference-based argumentation framework (a non-abstract *PAF*, for short) translated from a PLP, whose semantics is also given by  $\mathcal{P}$ -extensions instantiating those of the corresponding abstract *PAF*. As a result, we can show the interesting result that,  $\mathcal{P}$ -extensions of such a non-abstract *PAF* under stable semantics capture *preferred* answer sets of a PLP, which generalizes Dung's theorem about relation between answer sets of an ELP  $P$  and stable extensions of the argumentation framework associated with  $P$ . Thus this property ensures the advantages as well as the correctness of our approach.

Finally under an inconsistent knowledge base, the PLP system [21,20] is unable to reason at all, whereas the non-abstract *PAF* translated from a PLP is able to reason the intended results based on preferred semantics, i.e. preferred  $\mathcal{P}$ -extensions. Therefore we can regard such a non-abstract *PAF* as the enhanced PLP so that it can also reason paraconsistently from inconsistent knowledge bases.

This paper is organized as follows: Section 2 gives the preliminaries. Section 3 presents a new abstract *PAF*. Section 4 presents the non-abstract *PAF* translated

from a PLP and the semantics. Section 5 discusses the related work and Section 6 concludes the paper.

## 2 Preliminaries

We briefly review the basic notions used throughout this paper.

### 2.1 Extended Logic Programs and Answer Set Semantics

The logic programs we consider in this paper are extended logic programs (ELPs), which have two kinds of negation, i.e. classical negation ( $\neg$ ) along with negation as failure (*not*) defined as follows.

**Definition 1.** An extended logic program (ELP) [13,12] is a set of rules of the form:

$$L \leftarrow L_1, \dots, L_m, \text{not}L_{m+1}, \dots, \text{not}L_n, \quad (1)$$

or of the form:

$$\leftarrow L_1, \dots, L_m, \text{not}L_{m+1}, \dots, \text{not}L_n, \quad (2)$$

where  $L$  and  $L_i$ 's are literals, i.e. either atoms or atoms preceded by the classical negation sign  $\neg$  and  $n \geq m \geq 0$ . The symbol “not” denotes negation as failure. We call a literal preceded by “not” a NAF-literal. For a rule  $r$  of the form (1), we call  $L$  the head of the rule,  $\text{head}(r)$ , and  $\{L_1, \dots, L_m, \text{not}L_{m+1}, \dots, \text{not}L_n\}$  the body of the rule,  $\text{body}(r)$ . Especially,  $\text{body}(r)^+$  and  $\text{body}(r)^-$  denote  $\{L_1, \dots, L_m\}$  and  $\{L_{m+1}, \dots, L_n\}$  respectively. We often write  $L \leftarrow \text{body}(r)^+, \text{not } \text{body}(r)^-$  instead of (1) by using sets,  $\text{body}(r)^+$  and  $\text{body}(r)^-$ . Each rule of the form (2) is called an integrity constraint. For a rule with an empty body, we may write  $L$  instead of  $L \leftarrow$ . As usual, a rule with variables stands for the set of its ground instances.

The semantics of an ELP is given by the answer sets [13,12] as follows.

**Definition 2.** Let  $\text{Lit}_P$  be the set of all ground literals in the language of  $P$ . First, let  $P$  be a not-free ELP (i.e., for each rule  $m = n$ ). Then,  $S \subseteq \text{Lit}_P$  is an answer set of  $P$  if  $S$  is a minimal set satisfying the conditions:

1. For each ground instance of a rule  $L \leftarrow L_1, \dots, L_m$  in  $P$ , if  $\{L_1, \dots, L_m\} \subseteq S$ , then  $L \in S$ . In particular, for each integrity constraint  $\leftarrow L_1, \dots, L_m$  in  $P$ ,  $\{L_1, \dots, L_m\} \not\subseteq S$  holds;
2. If  $S$  contains a pair of complementary literals, then  $S = \text{Lit}_P$ .

Second, let  $P$  be any ELP and  $S \subseteq \text{Lit}_P$ . The reduct of  $P$  by  $S$  is a not-free ELP  $P^S$  whose form is  $L \leftarrow L_1, \dots, L_m$ , or  $\leftarrow L_1, \dots, L_m$ , iff there is a ground rule of the form (1), (2) in  $P$  s.t.  $\{L_{m+1}, \dots, L_n\} \cap S = \emptyset$ . Then,  $S$  is an answer set of  $P$  if  $S$  is an answer set of  $P^S$ .

An answer set is consistent if it is not  $\text{Lit}_P$ . A program  $P$  is consistent if it has a consistent answer set; otherwise,  $P$  is inconsistent. We write  $P \models L$  if a literal  $L$  is included in every answer set of  $P$ .

## 2.2 Prioritized Logic Programs and Preferred Answer Sets

A prioritized logic program (PLP) [19] is defined as follows.

**Definition 3 (Priorities).** *Given an ELP  $P$  and the set of ground literals  $Lit_P$ , a reflexive and transitive relation  $\preceq$  is defined on  $Lit_P$ . For any element  $e_1$  and  $e_2$  from  $Lit_P$ ,  $e_1 \preceq e_2$  is called a priority, and we say  $e_2$  has a higher priority than  $e_1$ . We write  $e_1 \prec e_2$  if  $e_1 \preceq e_2$  and  $e_2 \not\preceq e_1$ , and say  $e_2$  has a strictly higher priority than  $e_1$ .*

**Definition 4 (Prioritized Logic Programs, PLPs).** *A prioritized logic program (PLP, for short) is defined as a pair  $(P, \Phi)$ , where  $P$  is an ELP [1] and  $\Phi$  is a set of priorities on  $Lit_P$ .*

The declarative semantics of a PLP  $(P, \Phi)$  is given by *preferred answer sets* which are selected from answer sets of  $P$  based on the preference relation  $\sqsubseteq_{as}$  derived from priorities in  $\Phi$ . In what follows, the closure  $\Phi^*$  is defined as the set of priorities which are reflexively or transitively derived using priorities in  $\Phi$ .

**Definition 5 (Preferences between answer sets).** *Given a PLP  $(P, \Phi)$ , the preference relation  $\sqsubseteq_{as}$  over answer sets of  $P$  is defined as follows: For any answer sets  $S_1, S_2$  and  $S_3$  of  $P$ ,*

1.  $S_1 \sqsubseteq_{as} S_1$ ,
2.  $S_1 \sqsubseteq_{as} S_2$  if for some literal  $e_2 \in S_2 \setminus S_1$ ,
  - (i) there is a literal  $e_1 \in S_1 \setminus S_2$  such that  $e_1 \preceq e_2 \in \Phi^*$ , and
  - (ii) there is no literal  $e_3 \in S_1 \setminus S_2$  such that  $e_2 \prec e_3 \in \Phi^*$ ,
3. if  $S_1 \sqsubseteq_{as} S_2$  and  $S_2 \sqsubseteq_{as} S_3$ , then  $S_1 \sqsubseteq_{as} S_3$ .

We say that  $S_2$  is preferable to  $S_1$  with respect to  $\Phi$  if  $S_1 \sqsubseteq_{as} S_2$  holds. We write  $S_1 \sqsubset_{as} S_2$  if  $S_1 \sqsubseteq_{as} S_2$  and  $S_2 \not\sqsubseteq_{as} S_1$ . Hereafter, each  $S_1 \sqsubseteq_{as} S_2$  is called a preference between answer sets.

**Definition 6 (Preferred answer sets).** *Let  $(P, \Phi)$  be a PLP. Then, an answer set  $S$  of  $P$  is called a preferred answer set (or p-answer set, for short) of  $(P, \Phi)$  if  $S \sqsubseteq_{as} S'$  implies  $S' \sqsubseteq_{as} S$  (with respect to  $\Phi$ ) for any answer set  $S'$  of  $P$ .*

## 2.3 Abstract/Non-abstract Argumentation Frameworks and Acceptability Semantics

Dung presented an abstract argumentation framework and acceptability semantics [11] defined as follows.

**Definition 7 (Abstract Argumentation Frameworks).** *An abstract argumentation framework is a pair  $AF=(\mathcal{A}, \mathcal{R})$  where  $\mathcal{A}$  is a set of arguments and  $\mathcal{R}$  is a binary relation over  $\mathcal{A}$ , i.e.  $\mathcal{R} \subseteq \mathcal{A} \times \mathcal{A}$ .  $(a, b) \in \mathcal{R}$ , or equivalently  $a \mathcal{R} b$ , means that  $a$  attacks  $b$ . A set  $S$  of arguments attacks an argument  $a$  if  $a$  is attacked by an argument of  $S$ .*

<sup>1</sup> In this paper, for a PLP  $(P, \Phi)$ ,  $P$  is restrictedly given as an ELP though such  $P$  is originally allowed to be a GEDP, i.e. a member of the superclass of an ELP [19].

**Definition 8 (Acceptable sets / Conflict-free sets).** Let  $AF=(A, \mathcal{R})$  be an argumentation framework. A set  $S \subseteq A$  is conflict-free iff there are no arguments  $a$  and  $b$  in  $S$  such that  $a$  attacks  $b$ . An argument  $a \in A$  is acceptable w.r.t. a set  $S \subseteq A$  iff for any  $b \in A$  such that  $(b, a) \in \mathcal{R}$ , there exists  $c \in S$  such that  $(c, b) \in \mathcal{R}$ .

**Definition 9 (Acceptability Semantics).** Let  $AF=(A, \mathcal{R})$  be an argumentation framework and  $E \subseteq A$  be a conflict-free set of arguments. Let  $F : 2^A \rightarrow 2^A$  be a function with  $F(E) = \{a \mid a \text{ is acceptable w.r.t. } E\}$ .

Acceptability Semantics such as complete (resp. stable, preferred, grounded) semantics is given by the respective extensions defined as follows.  $E$  is admissible iff  $E \subseteq F(E)$ .  $E$  is a complete extension iff  $E = F(E)$ .  $E$  is a grounded extension iff  $E$  is a minimal (w.r.t. set-inclusion) complete extension.  $E$  is a preferred extension iff  $E$  is a maximal (w.r.t. set-inclusion) complete extension.  $E$  is a stable extension iff  $E$  is a preferred extension that attacks every argument in  $A \setminus E$ .

**Definition 10 (Credulous Justification vs Skeptical Justification).**

Let  $AF=(A, \mathcal{R})$  be an argumentation framework. and  $S_{name}$  be one of complete, stable, preferred, and grounded. Then for an argument  $a \in A$ ,

- $a$  is credulously justified (w.r.t.  $(A, \mathcal{R})$ ) under a  $S_{name}$  semantics iff  $a$  is contained in at least one  $S_{name}$  extension of  $(A, \mathcal{R})$ ;
- $a$  is skeptically justified (w.r.t.  $(A, \mathcal{R})$ ) under a  $S_{name}$  semantics iff  $a$  is contained in every  $S_{name}$  extension of  $(A, \mathcal{R})$ .

Non-abstract argumentation formalisms for ELPs [15,17] are defined as follows.

**Definition 11 (Arguments).** [17] Let  $P$  be an extended logic program whose rules have the form (1). An argument associated with  $P$  is a finite sequence  $Ag = [r_1; \dots; r_n]$  of ground instances of rules  $r_i \in P$  such that for every  $1 \leq i \leq n$ , for every literal  $L_j$  in the body of  $r_i$  there is a  $k > i$  such that  $\text{head}(r_k) = L_j$ . The head of a rule in  $Ag$ , i.e.  $\text{head}(r_i)$  is called a conclusion of  $Ag$ , whereas a NAF-literal not  $L$  in the body of a rule of  $Ag$  is called an assumption of  $Ag$ . We write  $\text{assm}(Ag)$  for the set of assumptions and  $\text{conc}(Ag)$  for the set of conclusions of an argument  $Ag$ . Especially we call the head of the first rule  $r_1$  the claim of an argument  $Ag$  as written  $\text{claim}(Ag)$ .

A subargument of  $Ag$  is a subsequence of  $Ag$  which is an argument. An argument  $Ag$  with a conclusion  $L$  is a minimal argument for  $L$  if there is no subargument of  $Ag$  with conclusion  $L$ . An argument  $Ag$  is minimal if it is minimal for its claim, i.e.  $\text{claim}(Ag)$ . Given an extended logic program  $P$ , the set of minimal arguments associated with  $P$  is denoted by  $\text{Args}_P$ .

As usual, the notions of attack such as “rebut”, “undercut”, “attack”, “defeat” abbreviated to  $\mathbf{r}$ ,  $\mathbf{u}$ ,  $\mathbf{a}$ ,  $\mathbf{d}$  are defined as a binary relation over  $\text{Args}_P$  as follows.

**Definition 12 (Rebut, Undercut, Attack, Defeat).** For two arguments,  $Ag_1$  and  $Ag_2$ , the notions of attack such as rebut, undercut, attack, defeat ( $\mathbf{r}$ ,  $\mathbf{u}$ ,  $\mathbf{a}$ ,  $\mathbf{d}$  for short) are defined as follows:



- $Ag_1$  rebuts  $Ag_2$ , i.e.  $(Ag_1, Ag_2) \in \mathbf{r}$  if there exists a literal  $L$  such that  $L \in \text{conc}(Ag_1)$  and  $\neg L \in \text{conc}(Ag_2)$ ;
- $Ag_1$  undercuts  $Ag_2$ , i.e.  $(Ag_1, Ag_2) \in \mathbf{u}$  if there exists a literal  $L$  such that  $L \in \text{conc}(Ag_1)$  and not  $L \in \text{assm}(Ag_2)$ ;
- $Ag_1$  attacks  $Ag_2$ , i.e.  $(Ag_1, Ag_2) \in \mathbf{a}$  if  $Ag_1$  rebuts or undercuts  $Ag_2$ ;
- $Ag_1$  defeats  $Ag_2$ , i.e.  $(Ag_1, Ag_2) \in \mathbf{d}$  if  $Ag_1$  undercuts  $Ag_2$ , or  $Ag_1$  rebuts  $Ag_2$  and  $Ag_2$  does not undercut  $Ag_1$ .

**Definition 13 (Abstract vs non-Abstract Argumentation Frameworks).**

Let  $P$  be an ELP,  $Args_P$  be the set of minimal arguments associated with  $P$  and  $attacks_P$  be the binary relation over  $Args_P$  defined according to some notion of attack (e.g.  $\mathbf{r}$ ,  $\mathbf{u}$ ,  $\mathbf{a}$ ,  $\mathbf{d}$ ). Then we call  $AF_P \stackrel{\text{def}}{=} (Args_P, attacks_P)$  the “non-abstract argumentation framework” associated with  $P$ .

Although Dung’s acceptability semantics is defined as the set of extensions under the specified argumentation semantics w.r.t. an abstract  $AF = (\mathcal{A}, \mathcal{R})$ , it is also given as the set of extensions w.r.t. the non-abstract  $AF_P = (Args_P, attacks_P)$  instantiating  $AF$  using an ELP  $P$ .

**2.4 Answer Set Programming as Argumentation**

Dung [10] showed that stable extensions of the argumentation framework  $AF_P$  associated with an ELP  $P$  without integrity constraints capture answer set semantics of  $P$  as follows.

**Theorem 1.** *Let  $P$  be an ELP having no integrity constraints, and  $AF_P = (Args_P, attack_P)$  be the concrete argumentation framework associated with  $P$ , where  $attacks_P$  is the binary relation over  $Args_P$  defined according to undercut (i.e.  $\mathbf{u}$ ) as the notion of attack. Then  $S$  is an answer set of  $P$  iff there is a stable extension  $E$  of  $AF_P$  such that*

$$S = \{ L \mid L \text{ is a literal s.t. } L = \text{claim}(Ag) \text{ for an argument } Ag \in E \} \quad \square$$

**3 A New Approach of an Abstract Preference-Based Argumentation Framework**

We present a new approach of an abstract preference-based argumentation framework (an abstract  $PAF$  for short), where Dung’s acceptability semantics is extended in a natural way so as to take into account additional preferences.

An abstract  $PAF$  takes as input three elements: a set  $\mathcal{A}$  of arguments, an attack relation  $\mathcal{R}$  on  $\mathcal{A}$ , and a preorder  $\leq$  on  $\mathcal{A}$ , where a pair  $AF = (\mathcal{A}, \mathcal{R})$  coincides with Dung’s argumentation framework. It returns *extensions* that are subsets of  $\mathcal{A}$  satisfying two basic requirements as addressed by Amgoud [2] as follows.

**Conflict-freeness:** If  $E$  is an extension (i.e.  $\mathcal{P}$ -extension in our approach) of

---

<sup>2</sup> In [10], it is expressed that  $S = \{ L \mid L \text{ is supported by an argument from } E \}$ .

$PAF=(\mathcal{A}, \mathcal{R}, \leq)$ , then  $E$  is conflict free w.r.t.  $\mathcal{R}$ .

**Generalization:** Dung's acceptability semantics of  $AF=(\mathcal{A}, \mathcal{R})$  is captured as the special case of the semantics of  $PAF=(\mathcal{A}, \mathcal{R}, \leq)$ .

Our formalism of  $PAF$  satisfies these basic requirements. It is based on the idea that an arguing agent wants to filter out extensions of the traditional acceptability semantics according to his/her preferences. Our approach is defined as follows.

**Definition 14 (Priorities between arguments).** A reflexive and transitive relation  $\leq$  is defined over  $A$ . For any element  $a_1$  and  $a_2$  from  $\mathcal{A}$ ,  $a_1 \leq a_2$ , or equivalently  $(a_1, a_2) \in \leq$ , is called a *priority*, and we say  $a_2$  has a higher priority than  $a_1$ . We write  $a_1 < a_2$  if  $a_1 \leq a_2$  and  $a_2 \not\leq a_1$ , and say  $a_2$  strictly has a higher priority than  $a_1$ .

**Definition 15 (Preference-based Argumentation Frameworks).**

A preference-based argumentation framework ( $PAF$ ) is a tuple  $PAF=(\mathcal{A}, \mathcal{R}, \leq)$ , where  $\mathcal{A}$  is a set of arguments,  $\mathcal{R}$  is an attack relation on  $\mathcal{A}$ , and  $\leq$  is a preorder (i.e., a reflexive and transitive relation) on  $\mathcal{A}$ , called a priority relation.

The semantics of an abstract  $PAF=(\mathcal{A}, \mathcal{R}, \leq)$  is given as *preferable extensions* (or  $\mathcal{P}$ -extensions) which are maximal arguments in  $\mathcal{A}$  w.r.t.  $\sqsubseteq_{ex}$  defined as follows.

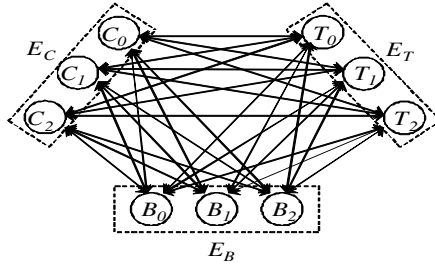
**Definition 16 (Preferences between extensions).** Given  $PAF=(\mathcal{A}, \mathcal{R}, \leq)$  and  $Sname \in \{\text{complete, stable, preferred, grounded}\}$ , let  $\mathcal{E}$  be the set of extensions for  $AF = (\mathcal{A}, \mathcal{R})$  under  $Sname$  semantics. Then the preference relation  $\sqsubseteq_{ex}$  over  $\mathcal{E}$  (i.e.,  $\sqsubseteq_{ex} \subseteq \mathcal{E} \times \mathcal{E}$ ) is defined as follows. For any  $Sname$  extensions  $E_1, E_2$  and  $E_3$  from  $\mathcal{E}$ ,

1.  $E_1 \sqsubseteq_{ex} E_1$ ,
2.  $E_1 \sqsubseteq_{ex} E_2$  if for some argument  $a_2 \in E_2 \setminus E_1$ ,
  - (i) there is an argument  $a_1 \in E_1 \setminus E_2$  such that  $a_1 \leq a_2$  w.r.t.  $\leq$ , and
  - (ii) there is no argument  $a_3 \in E_1 \setminus E_2$  such that  $a_2 < a_3$  w.r.t.  $\leq$ ,
3. if  $E_1 \sqsubseteq_{ex} E_2$  and  $E_2 \sqsubseteq_{ex} E_3$ , then  $E_1 \sqsubseteq_{ex} E_3$ .

Note that  $\sqsubseteq_{ex}$  is reflexive and transitive according to the items no.1 and no.3. We say that  $E_2$  is preferable to  $E_1$  with respect to  $\leq$  if  $E_1 \sqsubseteq_{ex} E_2$  holds. We write  $E_1 \sqsubseteq_{ex} E_2$  if  $E_1 \sqsubseteq_{ex} E_2$  and  $E_2 \not\sqsubseteq_{ex} E_1$ . Hereafter, each  $E_1 \sqsubseteq_{ex} E_2$  is called a *preference between extensions*.

*Example 1* Consider  $PAF=(\mathcal{A}, \mathcal{R}, \leq)$ , where  $\mathcal{A} = \{a, b, c, d\}$ ,  $\mathcal{R} = \{(a, b), (b, a), (c, d), (d, c), (c, a), (b, d)\}$  and  $\leq$  is the reflexive and transitive closure of  $\{(a, b), (a, c), (b, d), (c, d)\}$ . Then both of  $\{a, d\}$  and  $\{b, c\}$  are preferred extensions as well as stable extensions of  $AF = (\mathcal{A}, \mathcal{R})$ , and  $\{b, c\} \sqsubseteq_{ex} \{a, d\}$ . Note that  $\{a, d\} \not\sqsubseteq_{ex} \{b, c\}$  by the presence of  $b \leq d$  and  $c \leq d$  in  $\leq$ .

**Definition 17 ( $\mathcal{P}$ -extensions).** Let  $\mathcal{E}$  be the set of  $Sname$  extensions (e.g. a set of preferred extensions) for  $AF = (\mathcal{A}, \mathcal{R})$ . Given  $PAF = (\mathcal{A}, \mathcal{R}, \leq)$ , a  $Sname$  extension  $E \in \mathcal{E}$  (e.g. a preferred extension) is called a  $Sname$   $\mathcal{P}$ -extension (e.g.



**Fig. 1.** The Argumentation Framework (AF) of Example 3

a preferred  $\mathcal{P}$ -extension) of PAF if  $E \sqsubseteq_{ex} E'$  implies  $E' \sqsubseteq_{ex} E$  (with respect to  $\leq$ ) for any Sname extension  $E' \in \mathcal{E}$ . In other words, a Sname extension  $E$  is a Sname  $\mathcal{P}$ -extension of PAF if  $E \not\sqsubseteq_{ex} E'$  (with respect to  $\leq$ ) for any Sname extension  $E' \in \mathcal{E}$ .

*Example 2 (Ex. 1, Cont.).*  $\{a, d\}$  is the preferred  $\mathcal{P}$ -extension as well as the stable  $\mathcal{P}$ -extension of PAF w.r.t.  $\leq$ , but  $\{b, c\}$  is neither of them.<sup>3</sup>

**Proposition 1 (Generalization).**  $E$  is a Sname extension of  $AF=(\mathcal{A}, \mathcal{R})$  iff  $E$  is a Sname  $\mathcal{P}$ -extension of  $PAF = (\mathcal{A}, \mathcal{R}, \leq)$  when  $\leq$  is empty.

*Example 3* We can illustrate our PAF by using the example of travel arrangements. Let us suppose that there are three alternative vehicles for traveling by car, train or bicycle. Then we get nine arguments as follows.

- $C_0$  : We travel by car if there is no evidence that we travel by bicycle or train.
- $C_1$  : We will be able to carry more baggage in case of traveling by car.
- $C_2$  : We will be less tired in case of traveling by car.
- $T_0$  : We travel by train if there is no evidence that we travel by car or bicycle.
- $T_1$  : We will not encounter traffic jams in case of traveling by a train.
- $T_2$  : We will be less tired in case of traveling by train.
- $B_0$  : We travel by bicycle if there is no evidence that we travel by train or car.
- $B_1$  : Our health will be promoted in case of traveling by bicycle.
- $B_2$  : We will not encounter traffic jams in case of traveling by bicycle.

W.r.t. these arguments, we can represent the argumentation framework AF:

$$AF = (\mathcal{A}, \mathcal{R}), \text{ where } \mathcal{A} = \{C_i, T_i, B_i \mid 0 \leq i \leq 2\} \text{ and } \mathcal{R} = \{(C_i, T_j), (T_j, C_i), (B_i, T_j), (T_j, B_i), (C_i, B_j), (B_j, C_i) \mid 0 \leq i, j \leq 2\},$$

which is shown as the directed graph in Fig. 1. Then there are three preferred extensions of this AF as follows.

$$E_C = \{C_0, C_1, C_2\}, E_T = \{T_0, T_1, T_2\}, E_B = \{B_0, B_1, B_2\}$$

<sup>3</sup> Based on Amgoud and Cayrol’s approach [1], the unique preferred (resp. stable) extension,  $\{b, d\}$  is obtained for this PAF. But it is not conflict free w.r.t.  $\mathcal{R}$ .

Now suppose that an agent has two preferences such as,

- (1) being less tired is preferable compared to promoting health, and
- (2) not encountering traffic jams is preferable compared to carrying more baggage, which are expressed as follows:

$$B_1 \leq T_2, B_1 \leq C_2, C_1 \leq T_1, C_1 \leq B_2$$

Taking account of these priorities, this example is represented by  $PAF = (\mathcal{A}, \mathcal{R}, \leq)$ , where  $\leq = \{(B_1, T_2), (B_1, C_2), (C_1, T_1), (C_1, B_2)\} \cup \{(x, x) | x \in \mathcal{A}\}$ .

According to Definition 16, preferences between extensions are derived as follows.

$$E_B \sqsubseteq_{ex} E_T, E_C \sqsubseteq_{ex} E_T, E_C \sqsubseteq_{ex} E_B, E_B \sqsubseteq_{ex} E_C$$

i.e.,  $E_T \not\sqsubseteq_{ex} E_B, E_T \not\sqsubseteq_{ex} E_C$ , Therefore we obtain  $E_T$  as the unique  $\mathcal{P}$ -extension of this  $PAF$  under preferred or stable semantics, which is the expected result.

**Remark.** It should be noted that according to Amgoud and Vesic's approach [3] (resp. Amgoud and Cayrol's approach [1]), the  $PAF$  of Example 3 has the same extensions as those of the basic  $AF$ , i.e.  $E_C, E_T, E_B$  under preferred or stable semantics respectively. This means that preferences do not filter the extensions of its basic framework based on their respective methods. Moreover Amgoud and Vesic [3] proposed rich  $PAFs$  to refine  $AFs$  by preferences. However, extensions of the rich  $PAF$  for Example 3 are obtained as  $\text{Max}(\{E_C, E_T, E_B\}, \succeq_d) = \{E_C, E_T, E_B\}$ , where  $\succeq_d$  is the *democratic relation* [4] given by them. This also indicates that preferences do not work well to filter extensions based on their rich  $PAF$  [3], whereas preferences effectively work well to select the intended extension  $E_T$  based on the proposed method, as shown in Example 3.

## 4 Preference-Based and Constrained Argumentation Capturing Prioritized LP

In Section 3, an abstract  $PAF$  is presented. In this section, a non-abstract preference-based argumentation framework (a non-abstract  $PAF$ , for short) compiled from a PLP expressing domain knowledge is proposed as follows:

In the following, let  $P$  (resp.  $IC$ ) be a set of rules of form (1) (resp. (2)). Then in answer set programming (ASP), the semantics of an ELP  $P \cup IC$  is given by answer sets which are selected from answer sets of  $P$  by taking into account the set  $IC$  of integrity constraints. On the other hand, the semantics of a prioritized logic program, i.e. a PLP  $(P \cup IC, \Phi)$ , is given by *preferred answer sets* ( $p$ -*answer sets*, for short) which are selected from *answer sets* of  $P \cup IC$  by taking into account the set  $\Phi$  of priorities between literals, where integrity constraints from  $IC$  and priorities from  $\Phi$  are regarded as hard and soft constraints respectively.

A similar idea is also applied to our formalisms of argumentation for handling preferences and constraints. That is, our basic idea is that, given an ELP  $P \cup IC$  as the underlying logic and specified a particular Dung's argumentation semantics, the semantics of the constrained argumentation framework,

<sup>4</sup> The *democratic relation* [3] is the variant of the Hoare ordering. (See [6].)

i.e.  $CAF(P, IC) = (Args_P, attacks_P, IC)$  is given by  $\mathcal{C}$ -extensions which are selected from extensions of the non-abstract argumentation framework  $AF_P = (Args_P, attacks_P)$  by taking into account integrity constraints from  $IC$ , whereas given a PLP  $(P \cup IC, \Phi)$  as the underlying logic and specified a particular Dung's argumentation semantics, the semantics of the non-abstract preference-based argumentation framework, i.e.  $PAF(P \cup IC, \Phi) = (Args_P, attacks_P, IC, \leq)$  translated from the PLP is given by  $\mathcal{P}$ -extensions which are selected from  $\mathcal{C}$ -extensions of the  $CAF$  by taking into account the set  $\leq$  of priorities between arguments as constructed via priorities between literals from  $\Phi$ .

#### 4.1 Constrained AFs Built on ELPs with Integrity Constraints

First, we show a constrained argumentation framework whose underlying logic is an ELP with integrity constraints as follows.

**Definition 18 (From ELPs with constraints to constrained AFs).**

The constrained argumentation framework  $CAF(P, IC)$  associated with an ELP  $P \cup IC$  is defined as follows:

$$CAF(P, IC) \stackrel{def}{=} (Args_P, attacks_P, IC),$$

where  $P$  and  $IC$  are sets of rules of form (1) and (2) respectively.

After defining the claims of a set of arguments, we show the definition of satisfiability of an extension w.r.t. constraints as follows.

**Definition 19 (The claims of a set of arguments).** Let  $E$  be a set of arguments. Then  $claims(E)$  which we call the claims of  $E$  is defined as follows:

$$claims(E) \stackrel{def}{=} \{L \mid L \text{ is a literal s.t. } L = claim(Ag) \text{ for an argument } Ag \in E\}.$$

**Definition 20 (Satisfiability).** Let  $CAF(P, IC) \stackrel{def}{=} (Args_P, attacks_P, IC)$  be a constrained argumentation framework associated with  $P \cup IC$ . Note that for a rule,  $r_{ic}$  from  $IC$  whose form is (2) as follows:

$$r_{ic} : \quad \leftarrow L_1, \dots, L_m, not L_{m+1}, \dots, not L_n, \quad (2)$$

$body(r_{ic})^+ = \{L_1, \dots, L_m\}$  and  $body(r_{ic})^- = \{L_{m+1}, \dots, L_n\}$ .

Then for  $E \subseteq Args_P$ , whether  $E$  satisfies  $IC$  is defined as follows.

- $E$  violates  $IC$  iff  $claims(E) \cup IC$  is inconsistent  
iff  $\exists r_{ic} \in IC$  s.t.  $body(r_{ic})^+ \subseteq claims(E)$  and  $body(r_{ic})^- \cap claims(E) = \emptyset$ .
- $E$  satisfies  $IC$  iff  $E$  does not violate  $IC$  iff  $claims(E) \cup IC$  is consistent  
iff  $\forall r_{ic} \in IC$  if  $body(r_{ic})^- \cap claims(E) = \emptyset$ , then  $body(r_{ic})^+ \not\subseteq claims(E)$ .

The semantics of a constrained argumentation framework is defined as follows.

**Definition 21 ( $\mathcal{C}$ -extensions).** Let  $CAF(P, IC) = (Args_P, attacks_P, IC)$  be the constrained argumentation framework associated with an ELP  $P \cup IC$  and

$AF_P = (Args_P, attacks_P)$  be the argumentation framework associated with  $P$ . Then the semantics of  $CAF(P, IC)$  is given by  $\mathcal{C}$ -extensions defined as follows. For  $E \subseteq Args_P$  and  $Sname \in \{\text{complete, preferred, stable, grounded}\}$ ,

- $E$  is  $\mathcal{C}$ -admissible iff  $E$  is admissible for  $AF_P$  and satisfies  $IC$ .
- $E$  is a  $Sname$   $\mathcal{C}$ -extension of  $CAF(P, IC)$  under  $Sname$  semantics iff  $E$  is an extension of  $AF_P$  under  $Sname$  semantics and satisfies  $IC$ .

The following Theorem extends Theorem 1 to handle integrity constraints.

**Theorem 2.** Let  $CAF(P, IC) \stackrel{def}{=} (Args_P, attacks_P, IC)$  be the constrained argumentation framework associated with an ELP  $P \cup IC$ , where  $P$  and  $IC$  are the sets of rules of the form (1) and (2) respectively, and  $attacks_P$  is the binary relation over  $Args_P$  which is defined according to undercut (i.e.  $\mathbf{u}$ ) as the notion of attack. Then  $S$  is an answer set of  $P \cup IC$  iff there is a stable  $\mathcal{C}$ -extension  $E$  of  $CAF(P, IC)$  such that  $S = claims(E)$ .

*Proof:* See appendix.

*Example 4.* Let us consider the following ELP  $P \cup IC$ :

$$\begin{aligned}
 P: \quad & p \leftarrow not\ q, \quad q \leftarrow not\ p, \\
 & q \leftarrow not\ r, \quad r \leftarrow not\ q. \\
 IC: \quad & \leftarrow p, r.
 \end{aligned}$$

$P$  has two answer sets,  $S_1$  and  $S_2$  such that  $S_1 = \{p, r\}$  and  $S_2 = \{q\}$ , whereas  $P \cup IC$  has the unique answer set,  $S_2$ . On the other hand, the set  $Args_P$  of minimal arguments associated with  $P$  is  $\{A, B, C, D\}$  such that

$$\begin{aligned}
 A &= [p \leftarrow not\ q], & B &= [q \leftarrow not\ p] \\
 C &= [q \leftarrow not\ r], & D &= [r \leftarrow not\ q],
 \end{aligned}$$

and the attack relation,  $attacks_P$  is derived according to undercut as follows,

$$attacks_P = \{(A, B), (B, A), (C, D), (D, C), (C, A), (B, D)\}.$$

Thus, w.r.t.  $AF_P = (Args_P, attacks_P)$  whose graph is shown on the left of Fig. 2, there are two preferred as well as stable extensions,  $E_1$  and  $E_2$  as follows:

$$E_1 = \{A, D\}, \quad E_2 = \{B, C\}$$

with  $claims(E_1) = \{p, r\}$  and  $claims(E_2) = \{q\}$ .



**Fig. 2.** Argumentation Frameworks ( $AF_P$ s) of Example 4 and Example 5

Instead,  $E_2$  is the preferred  $\mathcal{C}$ -extension as well as the stable  $\mathcal{C}$ -extension of  $CAF(P, IC) = (Args_P, attacks_P, IC)$ , but  $E_1$  is neither of them since  $claims(E_1) \cup IC$  is inconsistent, but  $claims(E_2) \cup IC$  is consistent. Note that  $claims(E_2)$  coincides with the answer set  $S_2 = \{q\}$  of  $P \cup IC$  as addressed by Theorem 2.

**Remark.** It is well-known in answer set programming (ASP) that we can express an integrity constraint by means of a rule of form (1) instead of (2), where its head is expressed by a newly introduced propositional symbol, say  $\alpha$ , as follows:

$$\alpha \leftarrow L_1, \dots, L_m, not L_{m+1}, \dots, not L_n, not \alpha. \tag{3}$$

It is obvious that integrity constraints of the form (3) are effective in not deriving stable extensions which violate such integrity constraints based on Theorem 1. However, they are ineffective for extensions under the other argumentation semantics like preferred semantics.

### 4.2 Preference-Based AFs Translated from PLPs

Here, we show the non-abstract preference-based argumentation framework translated from a PLP  $(P \cup IC, \Phi)$ .

**Definition 22 (From PLPs to Preference-based AFs).**

Given a PLP  $(P \cup IC, \Phi)$ , the non-abstract preference-based argumentation framework  $PAF(P, IC, \Phi)$  associated with the PLP is defined as follows:

$$PAF(P, IC, \Phi) = (Args_P, attacks_P, IC, \leq)$$

where  $\leq$  is a priority relation on  $Args_P$  such that,  $Ag_1 \leq Ag_2$  iff  $e_1 \preceq e_2 \in \Phi^*$  for  $claim(Ag_1) = e_1$  and  $claim(Ag_2) = e_2$ . For any argument  $Ag_1$  and  $Ag_2$  from  $Args_P$ ,  $Ag_1 \leq Ag_2$  or  $(Ag_1, Ag_2) \in \leq$  is called “a priority between arguments”, and we say  $Ag_2$  has a higher priority than  $Ag_1$ . We write  $Ag_1 < Ag_2$  if  $Ag_1 \leq Ag_2$  and  $Ag_2 \not\leq Ag_1$ , and say “ $Ag_2$  has a strictly higher priority than  $Ag_1$ ”. Note that  $\leq$  is a preorder, i.e. a reflexive and transitive relation. When  $IC$  is empty, we may write

$$PAF(P, \Phi) = (Args_P, attacks_P, \leq)$$

instead of  $PAF(P, \emptyset, \Phi) = (Args_P, attacks_P, \emptyset, \leq)$ .

In our approach, given a PLP  $(P \cup IC, \Phi)$ , preferences between  $\mathcal{C}$ -extensions are defined w.r.t.  $PAF(P, IC, \Phi)$  as follows.

**Definition 23 (Preferences between  $\mathcal{C}$ -extensions).**

For a PLP  $(P \cup IC, \Phi)$  and  $Sname \in \{\text{complete, preferred, stable, grounded}\}$ , let  $PAF(P, IC, \Phi) = (Args_P, attacks_P, IC, \leq)$  be the non-abstract preference-based argumentation framework associated with the PLP, and  $\mathcal{E}$  be the set of  $Sname$   $\mathcal{C}$ -extensions for  $CAF(P, IC)$  associated with  $P \cup IC$  under  $Sname$  semantics. Then the preference relation  $\sqsubseteq_{ex}$  over  $\mathcal{E}$  (i.e.,  $\sqsubseteq_{ex} \subseteq \mathcal{E} \times \mathcal{E}$ ) is defined as follows. For any  $\mathcal{C}$ -extensions,  $E_1, E_2$  and  $E_3$  from  $\mathcal{E}$ ,

1.  $E_1 \sqsubseteq_{ex} E_1$ ,
2.  $E_1 \sqsubseteq_{ex} E_2$  if for some argument  $Ag_2 \in E_2 \setminus E_1$ ,
  - (i) there is an argument  $Ag_1 \in E_1 \setminus E_2$  s.t.  $Ag_1 \leq Ag_2$  w.r.t.  $\leq$ , and
  - (ii) there is no argument  $Ag_3 \in E_1 \setminus E_2$  s.t.  $Ag_2 < Ag_3$  w.r.t.  $\leq$ ,
3. if  $E_1 \sqsubseteq_{ex} E_2$  and  $E_2 \sqsubseteq_{ex} E_3$ , then  $E_1 \sqsubseteq_{ex} E_3$ .

Note that  $\sqsubseteq_{ex}$  is reflexive and transitive according to the items no.1 and no.3. We say that  $E_2$  is preferable to  $E_1$  with respect to  $\leq$  if  $E_1 \sqsubseteq_{ex} E_2$  holds. We write  $E_1 \sqsubset_{ex} E_2$  if  $E_1 \sqsubseteq_{ex} E_2$  and  $E_2 \not\sqsubseteq_{ex} E_1$ . Hereafter, each  $E_1 \sqsubseteq_{ex} E_2$  is called “a preference between  $\mathcal{C}$ -extensions”.

The semantics of  $PAF(P, IC, \Phi)$  is given by  $\mathcal{P}$ -extensions as follows.

**Definition 24** ( *$\mathcal{P}$ -extensions*). For a PLP  $(P \cup IC, \Phi)$  and  $Sname \in \{\text{complete, preferred, stable, grounded}\}$ , let  $\mathcal{E}$  be the set of the  $Sname$   $\mathcal{C}$ -extensions for  $CAF(P, IC)$  associated with  $P \cup IC$ . Then a  $\mathcal{C}$ -extension  $E \in \mathcal{E}$  is called a  $Sname$   $\mathcal{P}$ -extension of  $PAF(P, IC, \Phi)$  associated with the PLP under  $Sname$  semantics if  $E \sqsubseteq_{ex} E'$  implies  $E' \sqsubseteq_{ex} E$  (with respect to  $\leq$ ) for any  $E' \in \mathcal{E}$ . In other words,  $E \in \mathcal{E}$  is called a  $Sname$   $\mathcal{P}$ -extension of  $PAF(P, IC, \Phi)$  iff  $E \not\sqsubset_{ex} E'$  with respect to  $\leq$  for any  $E' \in \mathcal{E}$ .

The following theorem shows that stable  $\mathcal{P}$ -extensions of  $PAF(P, IC, \Phi)$  capture preferred answer sets of a PLP  $(P \cup IC, \Phi)$ , which extends Theorem 2.

**Theorem 3.** Let  $PAF(P, IC, \Phi) = (Args_P, attacks_P, IC, \leq)$  be the non-abstract preference-based argumentation framework associated with a PLP  $(P \cup IC, \Phi)$ , where  $attacks_P$  is the binary relation over  $Args_P$  defined according to undercut (i.e.  $\mathbf{u}$ ). Then  $S$  is a preferred answer set (i.e.  $p$ -answer set) of a PLP  $(P \cup IC, \Phi)$  iff there is a stable  $\mathcal{P}$ -extension  $E$  of  $PAF(P, IC, \Phi)$  such that  $S = claims(E)$ .

*Proof:* See appendix.

The skeptical (resp. credulous) query-answering problem is uniformly handled for our preference-based argumentation framework as follows.

**Definition 25** (*Credulous / Skeptical query-answering*).

Let  $PAF(P, IC, \Phi) = (Args_P, attacks_P, IC, \leq)$  be the preference-based argumentation framework associated with a PLP  $(P \cup IC, \Phi)$ . Then for an argument  $Ag \in Args_P$  and  $Sname \in \{\text{complete, preferred, stable, grounded}\}$ ,

- $Ag$  is credulously justified w.r.t.  $PAF(P, IC, \Phi)$  under  $Sname$  semantics iff  $Ag$  is contained in at least one  $Sname$   $\mathcal{P}$ -extension of  $PAF(P, IC, \Phi)$ ;
- $Ag$  is skeptically justified w.r.t.  $PAF(P, IC, \Phi)$  under  $Sname$  semantics iff  $Ag$  is contained in every  $Sname$   $\mathcal{P}$ -extension of  $PAF(P, IC, \Phi)$ .

The following proposition denotes that Dung’s acceptability semantics is the special case of our preference-based argumentation semantics.

**Proposition 2.** For a PLP  $(P \cup IC, \Phi)$  whose  $IC$  and  $\Phi$  are empty,  $E$  is a  $Sname$  extension of an argumentation framework  $AF_P$  associated with  $P$  iff  $E$  is a  $Sname$   $\mathcal{P}$ -extension of  $PAF(P, IC, \Phi)$  associated with the PLP.



In the following examples, each  $attacks_P$  is constructed based on *undercut* as the notion of attack in order to illustrate Theorem 3.

*Example 5.* Let us consider the PLP  $(P, \Phi)$  of Example 4.2 in [19] as follows:

$$\begin{aligned}
 P: \quad & p \leftarrow not\ q, not\ r, \\
 & q \leftarrow not\ p, not\ r, \\
 & r \leftarrow not\ p, not\ q, \\
 & s \leftarrow p. \\
 \Phi: \quad & p \preceq q, r \preceq s.
 \end{aligned}$$

$P$  has three answer sets  $S_1 = \{p, s\}$ ,  $S_2 = \{q\}$ ,  $S_3 = \{r\}$ , whereas the PLP  $(P, \Phi)$  has the unique p-answer set,  $S_2 = \{q\}$  since  $S_3 \sqsubseteq_{as} S_1$ ,  $S_1 \sqsubseteq_{as} S_2$  and  $S_3 \sqsubseteq_{as} S_2$  due to  $p \preceq q$ ,  $r \preceq s$  from  $\Phi$ .

On the other hand, we obtain  $PAF(P, \Phi) = (Args_P, attacks_P, \leq)$  compiled from this  $(P, \Phi)$  according to Definition 22, where  $Args_P$  is  $\{A, B, C, D\}$  s.t.

$$\begin{aligned}
 A &= [p \leftarrow not\ q, not\ r], & B &= [q \leftarrow not\ p, not\ r] \\
 C &= [r \leftarrow not\ p, not\ q], & D &= [s \leftarrow p; p \leftarrow not\ q, not\ r]
 \end{aligned}$$

with  $claim(A) = \{p\}$ ,  $claim(B) = \{q\}$ ,  $claim(C) = \{r\}$  and  $claim(D) = \{s\}$ ,  $attacks_P$  is the binary relation derived according to *undercut* as follows,

$$\begin{aligned}
 & \{(A, B), (B, A), (C, A), (A, C), (B, C), (C, B), (B, D), (D, B), (C, D), (D, C)\} \\
 & \text{and } \leq = \{(A, B), (C, D)\} \cup \{(x, x) \mid x \in Args_P\}, \text{ since } p \preceq q \in \Phi \text{ for } claim(A) = \{p\}, \\
 & claim(B) = \{q\} \text{ and } r \preceq s \in \Phi \text{ for } claim(C) = \{r\}, claim(D) = \{s\}.
 \end{aligned}$$

Now  $AF_P = (Args_P, attacks_P)$  associated with this  $P$  whose graph is on the right of Fig 2 has three preferred as well as stable extensions as follows:

$$E_1 = \{A, D\}, \quad E_2 = \{B\}, \quad E_3 = \{C\}$$

where  $claims(E_1) = \{p, s\}$ ,  $claims(E_2) = \{q\}$  and  $claims(E_3) = \{r\}$ . Therefore  $E_2$  is the unique preferred (resp. stable)  $\mathcal{P}$ -extension of  $PAF(P, \Phi)$  since  $E_3 \sqsubseteq_{ex} E_1$ ,  $E_1 \sqsubseteq_{ex} E_2$  and  $E_3 \sqsubseteq_{ex} E_2$  due to  $(A, B) \in \leq$ ,  $(C, D) \in \leq$  and transitive law of  $\sqsubseteq_{ex}$ . Noted that the unique p-answer set,  $S_2$  of the PLP coincides with  $claims(E_2)$  for the stable  $\mathcal{P}$ -extension,  $E_2$  of  $PAF(P, \Phi)$ .

*Example 6.* Consider the PLP  $(P \cup IC, \Phi)$ , where  $P$  and  $\Phi$  are given in Example 5 and  $IC$  has the integrity constraint as follows:

$$IC: \quad \leftarrow q.$$

Then  $P$  has two answer sets  $S_1 = \{p, s\}$  and  $S_3 = \{r\}$ , whereas the PLP  $(P \cup IC, \Phi)$  has the unique p-answer set,  $S_1$  since  $S_3 \sqsubseteq_{as} S_1$ .

On the other hand, we obtain  $PAF(P, IC, \Phi) = (Args_P, attacks_P, IC, \leq)$  compiled from this PLP  $(P \cup IC, \Phi)$ , where  $AF_P = (Args_P, attacks_P)$  and  $\leq$  are the same as the ones shown in Example 5.

Though there are three preferred as well as stable extensions,  $E_1, E_2$  and  $E_3$  for this  $AF_P$ , both  $E_1$  and  $E_3$  are preferred as well as stable  $\mathcal{C}$ -extensions of this

$CAF(P, IC)$  but  $E_2$  is not because both  $claims(E_1) \cup IC$  and  $claims(E_3) \cup IC$  are consistent, but  $claims(E_2) \cup IC$  is inconsistent. As a result, according to Definition 23,  $E_1 = \{A, D\}$  is not only the unique preferred  $\mathcal{P}$ -extension but also the unique stable  $\mathcal{P}$ -extension of  $PAF(P, IC, \Phi)$ , but  $E_3$  is not. Note that, the unique p-answer set,  $S_1 = \{p, s\}$  of this PLP  $(P \cup IC, \Phi)$  coincides with  $claims(E_1)$  for  $E_1 = \{A, D\}$  of  $PAF(P, IC, \Phi)$ .

*Example 7* (Gordon’s Perfected Shipping Problem).

Let us consider the famous legal reasoning example from Gordon [14]. The problem is described as follows:

*“A person wants to find out if her security interest in a certain ship is perfected. According to the Uniform Commercial Code (UCC) which is a state law, a security interest in goods may be perfected by taking possession of the collateral. However, the federal Ship Mortgage Act (SMA) states that a security interest in a ship may only be perfected by filing a financing statement. She currently has possession of the ship, but a statement has not been filed. Both UCC and SMA are applicable: the question is which takes precedence here.”*

The situation is presented by the ELP  $P_1$  as follows.

$$\begin{aligned}
 P_1: \quad &perfected \leftarrow posses, ucc, && (UCC) \\
 &\neg perfected \leftarrow ship, \neg file, sma, && (SMA) \\
 &posses \leftarrow, \quad ship \leftarrow, \quad \neg file \leftarrow, \\
 &ucc \leftarrow not \neg perfected, \quad sma \leftarrow not perfected.
 \end{aligned}$$

Since the two laws are in conflict with one another, they lead to two answer sets  $S_1$  and  $S_2$  of  $P_1$  as follows.

$$\begin{aligned}
 S_1 &= \{perfected, posses, ship, \neg file, ucc\}. \\
 S_2 &= \{\neg perfected, posses, ship, \neg file, sma\}.
 \end{aligned}$$

Now, there are two well-known legal principles for resolving such conflict between laws as follows.

*“The principle of Lex Posterior gives precedence to newer laws, and the principle of Lex Superior gives precedence to laws supported by the higher authority. In our case, UCC is newer than the SMA, and the SMA has higher authority since it is a federal law.”* Such knowledge may be described as the following sets:

$$\Phi_1 = \{sma \preceq ucc\}, \quad \Phi_2 = \{ucc \preceq sma\}, \quad \Phi_3 = \{sma \preceq ucc, ucc \preceq sma\},$$

where  $\Phi_1$  takes account of only the principle of Lex Posterior,  $\Phi_2$  only Lex Superior, and  $\Phi_3$  both. Then  $S_1$  (resp.  $S_2$ ) is the unique p-answer set of  $(P_1, \Phi_1)$  (resp.  $(P_1, \Phi_2)$ ), but both of  $S_1$  and  $S_2$  become tie p-answer sets of  $(P_1, \Phi_3)$  since  $S_1 \sqsubseteq_{as} S_2$  and  $S_2 \sqsubseteq_{as} S_1$  due to the conflict between these principles.

On the other hand, this problem can be represented by the preference-based argumentation framework  $PAF(P_1, \Phi_i) = (Args_P, attacks_P, \leq_i)$  compiled from  $(P_1, \Phi_i)$  ( $1 \leq i \leq 3$ ), where  $Args_{P_1}$  is  $\{A, B, C, D, F, G, H\}$  such that,

$$\begin{aligned}
A &= [\textit{perfected} \leftarrow \textit{posses}, \textit{ucc}; \textit{posses}; \textit{ucc} \leftarrow \textit{not} \neg\textit{perfected}], \\
B &= [\neg\textit{perfected} \leftarrow \textit{ship}, \neg\textit{file}, \textit{sma}; \textit{ship}; \neg\textit{file}; \textit{sma} \leftarrow \textit{not} \textit{perfected}], \\
C &= [\textit{ucc} \leftarrow \textit{not} \neg\textit{perfected}], \\
D &= [\textit{sma} \leftarrow \textit{not} \textit{perfected}], \\
F &= [\textit{posses} \leftarrow], \quad G = [\textit{ship} \leftarrow], \quad H = [\neg\textit{file} \leftarrow]
\end{aligned}$$

with  $\textit{claim}(A) = \{\textit{perfected}\}$ ,  $\textit{claim}(B) = \{\neg\textit{perfected}\}$ ,  $\textit{claim}(C) = \{\textit{ucc}\}$ ,  $\textit{claim}(D) = \{\textit{sma}\}$ ,  $\textit{claim}(F) = \{\textit{posses}\}$ ,  $\textit{claim}(G) = \{\textit{ship}\}$ ,  $\textit{claim}(H) = \{\neg\textit{file}\}$ ,  $\textit{attacks}_{P_1}$  is  $\{(A, B), (B, A), (A, D), (B, C)\}$  derived according to *undercut*, and each  $\leq_i$  is the binary relation over  $\textit{Args}_{P_1}$  such that  $\leq_1 = \{(D, C)\} \cup \Psi$ ,  $\leq_2 = \{(C, D)\} \cup \Psi$ ,  $\leq_3 = \{(C, D), (D, C)\} \cup \Psi$  where  $\Psi = \{(x, x) | x \in \textit{Args}_{P_1}\}$  due to the respective  $\Phi_i$ . In this case,  $AF_{P_1} = (\textit{Args}_{P_1}, \textit{attacks}_{P_1})$  has two preferred as well as stable extensions,  $E_1 = \{A, C, F, G, H\}$  and  $E_2 = \{B, D, F, G, H\}$  with  $\textit{claims}(E_1) = \{\textit{perfected}, \textit{ucc}, \textit{posses}, \textit{ship}, \neg\textit{file}\}$  and  $\textit{claims}(E_2) = \{\neg\textit{perfected}, \textit{sma}, \textit{posses}, \textit{ship}, \neg\textit{file}\}$ .

According to  $\leq_1$  (resp.  $\leq_2$ ),  $E_1$  (resp.  $E_2$ ) is the unique preferred as well as stable  $\mathcal{P}$ -extension of  $PAF(P_1, \Phi_1)$  (resp.  $PAF(P_1, \Phi_2)$ ), but both  $E_1$  and  $E_2$  are the preferred as well as stable  $\mathcal{P}$ -extensions of  $PAF(P_1, \Phi_3)$  since  $E_1 \sqsubseteq_{ex} E_2$  and  $E_2 \sqsubseteq_{ex} E_1$  due to  $\leq_3$ .

The following example shows that even for a PLP  $(P, \Phi)$  whose  $P$  is inconsistent, intended results of argumentation are derived based on the PAF.

*Example 8 (Ex. 7 Cont.)*. Consider the PLP  $(P_2, \Phi_i)$  ( $1 \leq i \leq 3$ ) such that  $P_2 = P_1 \cup \{ab \leftarrow \textit{not} ab\}$ , Due to the added rule to  $P_1$ ,  $P_2$  is inconsistent since it has no answer sets. Hence the PLP  $(P_2, \Phi_i)$  with any  $\Phi_i$  has no  $p$ -answer sets. This reveals the limitation of answer set programming which is only applicable to consistent knowledge bases. Instead, for the PLP  $(P_2, \Phi_i)$ , we have

$$PAF(P_2, \Phi_i) = (\textit{Args}_{P_2}, \textit{attacks}_{P_2}, \leq_i) \quad (\textit{for } 1 \leq i \leq 3),$$

where  $\textit{Args}_{P_2} = \textit{Args}_{P_1} \cup \{I\}$  such that  $I = [ab \leftarrow \textit{not} ab]$  and  $\textit{attacks}_{P_2} = \textit{attacks}_{P_1} \cup \{(I, I)\}$  as derived according to *undercut*. In this case,  $AF_{P_2} = (\textit{Args}_{P_2}, \textit{attacks}_{P_2})$  has no stable extensions but has the same two preferred extensions,  $E_1$  and  $E_2$  that  $AF_{P_1}$  has. Similarly, each  $PAF(P_2, \Phi_i)$  ( $1 \leq i \leq 3$ ) has no stable  $\mathcal{P}$ -extensions but has the same preferred  $\mathcal{P}$ -extensions that  $PAF(P_1, \Phi_i)$  has.

## 5 Related Work

Amgoud and Vesic [2] proposed only a new abstract PAF, whereas we present not only a new approach of an abstract PAF but also propose a non-abstract PAF constructed from a prioritized logic program. In our approach, we can show Theorem 3 for such a non-abstract PAF as is the generalization of Theorem 1 presented by Dung [10]. This property ensures the advantages as well as the correctness of our approach.

Coste-Marquis *et al.* [8] proposed an abstract *CAF* where constraints are expressed by a propositional formula defined over the set of abstract arguments, whereas in our approach, a non-abstract *CAF* is defined where constraints are given as nonmonotonic rules embedded in an extended logic program expressing an agent's domain knowledge. From the computational point of view, Besnard and Doutre's approach [4] for encoding acceptable semantics can be applied to their *CAF*, whereas a non-abstract *CAF* presented in this paper can be easily encoded in ASP setting by extending our previous work [22] to compute argumentation semantics in ASP based on Caminada's reinstatement labellings [7].

Šeřrnek [18] presented the semantics, i.e. preferred answer sets of a prioritized logic program  $(P, \prec, \mathcal{N})$  based on argumentation, where  $P$  is an ELP,  $\prec$  is a strict partial order on rules of  $P$  and  $\mathcal{N}$  is a function assigning names to rules of  $P$ . He proposed an argumentation framework translated from such a prioritized logic program, and defined preferred answer sets in his framework. However, not only argumentation framework proposed in [18] is inapplicable to an inconsistent  $P$  but also it is not the generalization of Dung's argumentation framework for handling additional preferences.

## 6 Conclusion

To handle preferences along with constraints, we presented a new abstract preference-based argumentation framework as well as a non-abstract one translated from a prioritized logic program. In our approach, we can show Theorem 3 such that, stable  $\mathcal{P}$ -extensions of the preference-based argumentation framework associated with a PLP  $(P \cup IC, \Phi)$  capture p-answer sets of the PLP. Hence the advantages and the correctness of our approach are ensured.

On the other hand, when agent's knowledge expressed by an ELP  $P$  is inconsistent, we cannot reason anything from the PLP  $(P, \Phi)$  as well as from our  $PAF(P, \Phi)$  under stable semantics, since there are no p-answer sets of the PLP as well as no stable  $\mathcal{P}$ -extensions of  $PAF(P, \Phi)$ . However, with such inconsistent  $P$ , we can infer the intended results from a non-abstract  $PAF(P, \Phi)$  under preferred semantics because there exists a preferred  $\mathcal{P}$ -extension for  $PAF(P, \Phi)$ . Thus in some sense, a non-abstract  $PAF$  presented in the paper can be regarded as the extended PLP.

Applying the techniques used in our previous work [20,21,22], the encoding to compute  $\mathcal{P}$ -extension of a non-abstract  $PAF$  can easily be established in an ASP setting. Thus such a system which encodes our  $PAF$  presented in the paper will behave as the enhanced PLP system such that not only it can compute p-answer sets of a PLP via the stable  $\mathcal{P}$ -extensions, but also it can infer intended results via the preferred  $\mathcal{P}$ -extensions even if  $P$  is inconsistent.

Our future work will not only investigate computational complexity of the proposed method but also will implement the  $PAF$  system in an ASP setting so that it may be used in practical multiagent systems of negotiation based on the proposed preference-based argumentation.

**Acknowledgments.** This research is partially supported by Grant-in-Aid for Scientific Research from JSPS, No. 20500141. The author thanks Nobuo Nagayoshi for his support to this research.

## References

1. Amgoud, L., Cayrol, C.: A reasoning model based on the production of acceptable arguments. *Annals of Mathematics and Artificial Intelligence* 34(1-3), 197–215 (2002)
2. Amgoud, L., Vesic, S.: Repairing preference-based argumentation frameworks. *Proceedings of IJCAI 2009*, 665–670 (2009)
3. Amgoud, L., Vesic, S.: Handling inconsistency with preference-based argumentation. In: Deshpande, A., Hunter, A. (eds.) *SUM 2010*. LNCS, vol. 6379, pp. 56–69. Springer, Heidelberg (2010)
4. Besnard, P., Doutre, S.: Checking the acceptability of a set of arguments. In: *Proceedings of 10th International Workshop on Non-Monotonic Reasoning (NMR-2004)*, pp. 59–64 (2004)
5. Brewka, G., Eiter, T.: Preferred answer sets for extended logic programs. *Artificial Intelligence* 109, 297–356 (1999)
6. Brewka, G., Truszczynski, M., Woltran, S.: Representing preferences among sets. In: *Proceedings of AAAI 2010*, pp. 273–278 (2010)
7. Caminada, M.: On the issue of reinstatement in argumentation. In: Fisher, M., van der Hoek, W., Konev, B., Lisitsa, A. (eds.) *JELIA 2006*. LNCS (LNAI), vol. 4160, pp. 111–123. Springer, Heidelberg (2006)
8. Coste-Marquis, S., Devred, C., Marquis, P.: Constrained argumentation frameworks. In: *Proceedings of KR 2006*, pp. 112–122 (2006)
9. Delgrande, J.P., Schaub, T., Tompits, H.: A framework for compiling preferences in logic programs. *Theory and Practice of Logic Programming* 3(2), 129–187 (2003)
10. Dung, P.M.: An argumentation semantics for logic programming with explicit negation. In: *Proceedings of ICLP 1993*, pp. 616–630. MIT Press, Cambridge (1993)
11. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning. logic programming, and n-person games. *Artificial Intelligence* 77, 321–357 (1995)
12. Gelfond, M., Lifschitz, V.: The stable model semantics for logic programming. In: *Proceedings of the Fifth International Conference and Symposium on Logic Programming (ICLP/SLP 1988)*, pp. 1070–1080. MIT Press, Cambridge (1988)
13. Gelfond, M., Lifschitz, V.: Classical negation in logic programs and disjunctive databases. *New Generation Computing* 9, 365–385 (1991)
14. Gordon, T.F.: *The pleadings game: An Artificial Intelligence Model of Procedural Justice*. Dissertation, TU Darmstadt (1993)
15. Prakken, H., Sartor, G.: Argument-based extended logic programming with defeasible priorities. *Journal of Applied Non-Classical Logics* 7(1), 25–75 (1997)
16. Prakken, H., Vreeswijk, G.A.W.: Logics for defeasible argumentation. In: Gabbay, D.M., Guenther, F. (eds.) *Handbook of Philosophical Logic*, vol. 4, pp. 218–319. Kluwer, Dordrecht (2001)
17. Schweimeier, R., Schroeder, M.: A Parameterized hierarchy of argumentation semantics for extended logic programming and its application to the well-founded semantics. *Theory and Practice of Logic Programming* 5(1,2), 207–242 (2005)

18. Šefracuteánek, J.: Preferred answer sets supported by arguments. In: Proceedings of 12th International Workshop on Non-Monotonic Reasoning (NMR 2008), pp.232-240 (2008)
19. Sakama, C., Inoue, K.: Prioritized logic programming and its application to commonsense reasoning. *Artificial Intelligence* 123, 185–222 (2000)
20. Wakaki, T., Inoue, K., Sakama, C., Nitta, K.: Computing preferred answer sets in answer set programming. In: Vardi, M.Y., Voronkov, A. (eds.) LPAR 2003. LNCS, vol. 2850, pp. 259–273. Springer, Heidelberg (2003)
21. Wakaki, T., Inoue, K., Sakama, C., Nitta, K.: The PLP system. In: Alferes, J.J., Leite, J. (eds.) JELIA 2004. LNCS (LNAI), vol. 3229, pp. 706–709. Springer, Heidelberg (2004)
22. Wakaki, T., Nitta, K.: Computing argumentation semantics in answer set programming. In: van Hoesve, W.-J., Hooker, J.N. (eds.) CPAIOR 2009. LNCS, vol. 5547, pp. 254–269. Springer, Heidelberg (2009)

## Appendix: Proofs of Theorems

### Proof of Theorem 2

*Proof.* ( $\Leftarrow$ ) Suppose  $E$  is a stable  $\mathcal{C}$ -extension of  $CAF(P, IC) = (Args_P, attacks_P, IC)$ . According to Definition 21,  $E$  is a stable extension of  $AF_P = (Args_P, attacks_P)$  and satisfies  $IC$ . Therefore, there is the answer set  $S$  of  $P$  such that  $S = claims(E)$  due to Theorem 1. Thus according to Definition 2,  $S$  is also an answer set of the *not*-free  $P^S$ , i.e. the reduct of  $P$ .

Now, since such  $E$  satisfies  $IC$ , which means that, for  $S = claims(E)$ ,

$$\forall r_{ic} \in IC \text{ if } body(r_{ic})^- \cap S = \emptyset, \text{ then } body(r_{ic})^+ \not\subseteq S,$$

the answer set  $S$  of  $P^S$  satisfies  $body(r_{ic})^+ = \{L_1, \dots, L_m\} \not\subseteq S$  if  $body(r_{ic})^- \cap S = \{L_{m+1}, \dots, L_n\} \cap S = \emptyset$  for any integrity constraint  $r_{ic} \in IC$  as follows:

$$r_{ic} : \quad \leftarrow L_1, \dots, L_m, not L_{m+1}, \dots, not L_n.$$

Therefore it is concluded that  $S$  is an answer set of  $(P \cup IC)^S$ . Hence  $S = claims(E)$  is an answer set of  $P \cup IC$ .

( $\Rightarrow$ ) The converse is also proved similarly. □

After preparing the following lemma, we show the proof of Theorem 3.

**Lemma 1.** *For a PLP  $(P \cup IC, \Phi)$ , let  $PAF(P, IC, \Phi)$  be  $(Args_P, attacks_P, IC, \leq)$ ,  $CAF(P, IC)$  be  $(Args_P, attacks_P, IC)$ ,  $E_1, E_2$  be stable  $\mathcal{C}$ -extensions of  $CAF(P, IC)$ , and  $S_1, S_2$  be answer sets of  $P \cup IC$ . Then it holds that,*

$$E_1 \sqsubseteq_{ex} E_2 \quad \text{iff} \quad S_1 \sqsubseteq_{as} S_2 \quad \text{for } S_1 = claims(E_1) \text{ and } S_2 = claims(E_2).$$

*Proof*

Suppose  $E$  is a stable  $\mathcal{C}$ -extension of  $CAF(P, IC)$ . Then according to Theorem 2,  $claims(E)$  coincides with an answer set  $S$  of  $P \cup IC$ . Moreover, for an argument  $Ag \in Args_P$  and its claim  $e \in Lit_P$ , i.e.  $e = claim(Ag)$ , it holds that,

$$Ag \in E \quad \text{iff} \quad e \in S, \quad \text{and} \quad Ag \notin E \quad \text{iff} \quad e \notin S. \quad (3)$$

Now with respect to stable  $\mathcal{C}$ -extensions  $E_1, E_2$  of  $CAF(P, IC)$  whose claims are  $S_1 = \text{claims}(E_1)$ ,  $S_2 = \text{claims}(E_2)$  respectively, it holds that, due to (3), for a literal  $e_2 \in \text{Lit}_P$  such that  $e_2 = \text{claim}(Ag_2)$ ,

$$\begin{aligned} Ag_2 \in E_2 \setminus E_1 \quad \text{iff} \quad Ag_2 \in E_2 \text{ and } Ag_2 \notin E_1 \quad \text{iff} \quad e_2 \in S_2 \text{ and } e_2 \notin S_1 \\ \text{iff} \quad e_2 \in S_2 \setminus S_1. \end{aligned} \quad (4)$$

Similarly for a literal  $e_1 \in \text{Lit}_P$  such that  $e_1 = \text{claim}(Ag_1)$ , it holds that,

$$Ag_1 \in E_1 \setminus E_2 \quad \text{iff} \quad e_1 \in S_1 \setminus S_2. \quad (5)$$

On the other hand, according to Definition 22,

$$Ag_1 \leq Ag_2 \quad \text{iff} \quad e_1 \preceq e_2 \in \Phi^* \quad \text{for} \quad \text{claim}(Ag_1) = e_1 \quad \text{and} \quad \text{claim}(Ag_2) = e_2. \quad (6)$$

Thus due to (4), (5), (6), it holds that,

$$\begin{aligned} \exists Ag_2 \in E_2 \setminus E_1 \quad \text{and} \quad \exists Ag_1 \in E_1 \setminus E_2 \quad \text{such that} \quad Ag_1 \leq Ag_2 \\ \text{iff} \quad \exists e_2 \in S_2 \setminus S_1 \quad \text{and} \quad \exists e_1 \in S_1 \setminus S_2 \quad \text{such that} \quad e_1 \preceq e_2 \in \Phi^*. \end{aligned} \quad (7)$$

Therefore by extending (7), it is obviously derived that,

$$\begin{aligned} \exists Ag_2 \in E_2 \setminus E_1 [ \exists Ag_1 \in E_1 \setminus E_2 \quad \text{such that} \quad Ag_1 \leq Ag_2 \\ \wedge \neg \exists Ag_3 \in E_1 \setminus E_2 \quad \text{s.t.} \quad Ag_2 < Ag_3 \quad \text{w.r.t.} \quad \leq ], \\ \text{iff} \quad \exists e_2 \in S_2 \setminus S_1 [ \exists e_1 \in S_1 \setminus S_2 \quad \text{such that} \quad e_1 \preceq e_2 \in \Phi^* \\ \wedge \neg \exists e_3 \in S_1 \setminus S_2 \quad \text{s.t.} \quad e_2 \prec e_3 \in \Phi^* ] \end{aligned} \quad (8)$$

where  $S_i = \text{claims}(E_i)$  and  $e_j = \text{claim}(Ag_j)$  ( $1 \leq i \leq 2, 1 \leq j \leq 3$ ).

(8) means that  $E_1 \sqsubseteq_{ex} E_2$  iff  $S_1 \sqsubseteq_{as} S_2$  for  $S_1 = \text{claims}(E_1)$  and  $S_2 = \text{claims}(E_2)$  w.r.t. the item no.2 of Definition 23 and that of Definition 5. Since both  $\sqsubseteq_{ex}$  and  $\sqsubseteq_{as}$  are reflexive and transitive, it also holds that,  $E_1 \sqsubseteq_{ex} E_2$  iff  $S_1 \sqsubseteq_{as} S_2$  w.r.t. items no.1 and no.3 of these definitions.  $\square$

### Proof of Theorem 3

*Proof:* For a PLP( $P \cup IC, \Phi$ ), let  $AS$  be the set of all answer sets of  $P \cup IC$  and  $\mathcal{E}$  be the set of all stable  $\mathcal{C}$ -extensions of  $CAF(P, IC) = (Args_P, attacks_P, IC)$ . Then, it follows that,

$E \in \mathcal{E}$  is a stable  $\mathcal{P}$ -extensions of  $PAF(P, IC, \leq)$  built on a PLP( $P \cup IC, \Phi$ )  
 iff  $E \sqsubseteq_{ex} E'$  implies  $E' \sqsubseteq_{ex} E$  (with respect to  $\leq$ ) for any  $E' \in \mathcal{E}$   
 iff w.r.t.  $S = \text{claims}(E) \in AS$ ,

$S \sqsubseteq_{as} S'$  implies  $S' \sqsubseteq_{as} S$  (with respect to  $\Phi$ ) for any  $S' = \text{claims}(E') \in AS$   
 due to Theorem 2 and Lemma 1,

iff  $S = \text{claims}(E) \in AS$  is a preferred answer set of  $(P, \Phi)$ .  $\square$

# Author Index

- Abdel-Naby, Sameh 141  
Atkinson, Katie 12
- Beaufils, Bruno 141  
Black, Elizabeth 12  
Botti, Vicente 123
- Capobianco, Marcela 171
- Dignum, F. 68  
Dimopoulos, Yannis 105
- Emele, Chukwuemeka David 86
- Falappa, Marcelo A. 228
- García, Alejandro J. 286  
Grossi, Davide 190  
Groza, Adrian 209  
Guerin, Frank 86
- Hadidi, Nabila 105  
Heras, Stella 123  
Hitchcock, David 1
- Julián, Vicente 123
- Kok, Eric M. 31
- Letia, Ioan Alfred 209
- Marcos, M. Julieta 228  
McBurney, Peter 159  
Meyer, John-Jules Ch. 31, 68  
Moraitis, Pavlos 105  
Morge, Maxime 141
- Norman, Timothy J. 86, 268
- Okuno, Kenichi 248  
Ontañón, Santiago 49
- Parsons, Simon 86, 159, 268  
Plaza, Enric 49  
Prakken, Henry 31, 68
- Simari, Guillermo R. 171, 228  
Sklar, Elizabeth 159
- Takahashi, Kazuko 248  
Tang, Yuqing 268  
Thimm, Matthias 286
- van der Weide, T.L. 68  
Vreeswijk, Gerard A.W. 31, 68
- Wakaki, Toshiko 306