

Picking Up The Best Goal

An Analytical Study in Defeasible Logic

Guido Governatori¹, Francesco Olivieri^{1,2,4}, Antonino Rotolo³, Simone Scannapieco^{1,2,4}, and Matteo Cristani²

¹ NICTA, Queensland Research Laboratory, Australia

² Department of Computer Science, University of Verona, Italy

³ CIRSIFID and DSG, University of Bologna, Italy

⁴ Institute for Integrated and Intelligent Systems, Griffith University, Australia

Abstract. In this paper we analyse different notions of the concept of goal starting from the idea of sequences of “alternative acceptable outcomes”. We study the relationships between goals and concepts like agent’s beliefs, norms and desires, and we propose a computationally oriented formalisation using Defeasible Logic that will be able to provide a computationally feasible approach. The resulting system is able to capture various nuances of the notion of goal against different normative domains, for which the right decision is not only context-dependent, but it must be chosen from a pool of alternatives as wide as possible.

1 Motivation and Basic Intuitions

The BDI architecture is a prominent approach to model rational agents [4, 16]. BDI agents are means-ends reasoners equipped with: (i.) Desires, Goals, Intentions (or Tasks); (ii.) a description of the current state of the environment (Beliefs); (iii.) Actions. That is the case, the key tenet of this architecture is that the agent behaviour is the outcome of a rational balance among different mental states. Typically these mental attitudes are taken as primitive and independent from each other, even though some mutual influences are considered (e.g., intentions are seen as desires satisfied up to commitment). Previous work implementing the BDI paradigm in Defeasible Logic [6, 11, 9] follow this schema.

We work here on a different perspective and to provide a fresh and efficient rule-based framework that considers goals, desires and intentions as facets of the same phenomenon (all of them being *goal-like attitudes*): the notion of *expected outcome*, which is simply something an agent would like or is expected to achieve. An advantage of the proposed framework is that it allows agents to compute different degrees of motivational attitudes, and degrees of commitment that take into account other factors, such as beliefs and norms.

Consider the following example. Alice, during her holidays, plans to pay a visit to her friend John, who lives close to her parents. The plan can be described by the sentence *I shall come over to John’s place to visit him on Monday, but if he is not home or the visit is not possible, I am going to visit my parents. If this is not possible as well, I shall take some rest at home.* This idea can be easily implemented by building for

each alternatives a sequence of other alternatives A_1, \dots, A_n that are preferred when the first choice is no longer feasible. Normally, each alternatives is the result of a specific context C . This scenario can be represented as $C \Rightarrow A_1, \dots, A_n$ which is closely related of contrary-to-duty obligations [10], where a norm is represented by a rule of the type represented as follows:

$$r_1 : \text{drive_car} \Rightarrow_O \neg \text{damage} \odot \text{compensate} \odot \text{foreclosure}.$$

Rule r_1 states that, if one agent drives a car, she has the obligation not to cause any damage to others; if this happens, she is obliged to compensate; if she fails to compensate, there is an obligation of foreclosure. The previous setting can be rewritten as:

$$r_2 : \text{holiday} \Rightarrow_U \text{visit_friend} \odot \text{visit_parents} \odot \text{stay_home}.$$

In both examples (rules r_1 and r_2), sequences express a preference ordering among outcomes, which means that also *stay_home* and *foreclosure*, though not the best options, still correspond to acceptable situations.

Besides rules for outcomes and obligations, we also have rules for beliefs such as

$$r_3 : \text{friend_away} \Rightarrow_B \neg \text{visit_friend}$$

for which we assume there is no preference ordering, since they do not express expected outcomes but simply describe how the world is.

These building blocks allow us to introduce different types of goal-like attitudes and degrees of commitment to outcomes: desires, goals, intentions, and social intentions.

Desires as acceptable outcomes Suppose an agent is equipped with the following outcome rules expressing two preference orderings:

$$r : a_1, \dots, a_n \Rightarrow_U b_1 \odot \dots \odot b_m \quad s : a'_1, \dots, a'_n \Rightarrow_U b'_1 \odot \dots \odot b'_k$$

and that the situation described by a_1, \dots, a_n and a'_1, \dots, a'_n are mutually compatible but b_1 and b'_1 are not, namely $b_1 = \neg b'_1$. In this case $b_1, \dots, b_m, b'_1, \dots, b'_k$ are anyway all *acceptable outcomes*, including the incompatible outcomes b_1 and b'_1 . *Desires* are expected or acceptable outcomes, independently of whether they are compatible with other expected or acceptable outcomes.

Goals as preferred outcomes For rule r alone the preferred outcome is b_1 , and for rule s alone it is b'_1 . But if both rules are applicable, then a state where both b_1 and b'_1 hold is not possible. Hence, the agent would not be rational if she considers both b_1 and $\neg b_1$ as her preferred outcomes. Then the agent has to decide if she prefers a state where b_1 holds to one where b'_1 (i.e., $\neg b_1$) holds, or *vice versa*. If the agent cannot make up her mind, i.e., she has no way to decide which is the most suitable option for her, then neither the chain of r nor that of s can produce preferred outcomes. Suppose that the agent opts for the latter option; this can be done if the agent establishes that the second rule overrides the first one, i.e., $s > r$. Accordingly, the preferred outcome is b'_1 for the chain of outcomes defined by s , and b_2 is the preferred outcome of r . b_2 is the second best alternative according to rule r : in fact b_1 has been discarded as an acceptable outcome given that s prevails over r .

Two degrees of commitment: intentions and social intentions Let us clarify when an agent commits to acceptable outcomes. If the agent values some outcomes more than others, she should strive for the best, i.e., for the most preferred outcomes. Let us start by considering the case where only rule r applies. Here, the agent should commit to the outcome she values the most, i.e., b_1 . But what if the agent *believes* that b_1 cannot be achieved in the environment where she is currently situated in, or she knows that $\neg b_1$ holds? Committing to b_1 would result in a waste of agent's resources; rationally, she should target the next best outcomes, in this case b_2 . Suppose, now, that b_2 is *forbidden*, and the agent is social (an agent is social if the agent would not knowingly commit to anything that is forbidden [11]). Once again, in this situation the agent has to lower her expectation and settle for b_3 , which is the next acceptable outcome.

To complete the analysis, consider the situation where both rules r and s apply and the agent prefers s to r . As we have seen before, $\neg b_1$ (b'_1) and b_2 are the preferred outcomes based on the preference of the agent over the two rules. Assume that, this time, the agent knows she cannot achieve $\neg b_1$ (or equivalently, b_1 holds). If the agent is rational, she cannot commit to $\neg b_1$. Thus, the best option for her is to commit to b'_2 and b_1 , where she is guaranteed to be successful. In this scenario, the best course of action for the agent is where she commits herself to some outcomes that are not her preferred ones, or even that she would consider not acceptable based only on her preferences, but such that they influence her decision process given that they represent relevant external factors (either her beliefs or the norms that apply to her).

The layout of the paper is as follows. Section 2 presents the new logical framework. Section 3 illustrates how to employ this framework to characterise various types of agent compliance. Section 4 describes the algorithms to prove that the logic has linear complexity. Section 5 ends the paper with some conclusions and a discussion of related work.

2 Logic

Defeasible Logic (DL) [1] is a simple but flexible and efficient rule based non-monotonic formalism. The strength of DL lays in the constructive proof theory, which has an argumentation-like structure and allows us to draw meaningful conclusions from (potentially) conflicting and incomplete knowledge bases. The framework provided by the proof theory accounts for the possibility of extensions of the logic, in particular extensions with modal operators. Several extensions have been proposed, which resulted in particular in applications in the area of normative reasoning [8] and modelling agents [11, 12, 9]. Efficient implementations of the logic (including the modal variants) are developed which are able to handle very large knowledge bases [13, 2, 18].

2.1 Language

The main aim of this subsection is to establish an inference process to compute factual knowledge, desires, intentions, goals and obligations from existing facts, primitive desires, intentions, goals and unconditional obligations. As a first step, we introduce the language adopted. Let PROP be a set of propositional atoms, $MOD = \{B, O, D, G, I, SI\}$

the set of modal operators¹ and Lbl be a set of arbitrary labels. The set $\text{Lit} = \text{PROP} \cup \{\neg p \mid p \in \text{PROP}\}$ denotes the set of *literals*. The *complementary* of a literal q is denoted by $\sim q$; if q is a positive literal p , then $\sim q$ is $\neg p$, and if q is a negative literal $\neg p$ then $\sim q$ is p . The set of *modal literals* is $\text{ModLit} = \{\Box l, \neg \Box l \mid l \in \text{Lit}, \Box \in \{\text{O}, \text{D}, \text{G}, \text{I}, \text{SI}\}\}$. We assume that the “ \Box ” modal operator for belief B is the empty modal operator, thus a modal literal Bl is equivalent to literal l . Accordingly, we state that the complementary of Bl is $\sim \text{Bl}$ as well as $\neg \text{Bl}$ is $\sim l$.

We define a *defeasible theory* D as a structure $(F, R, >)$, where (i.) F is a set of *facts* or indisputable statements, (ii.) R contains three sets of *rules*: for beliefs, obligations, and outcomes and (iii.) $> \subseteq R \times R$ is a *superiority relation* to determine the relative strength of conflicting rules. *Belief rules* are used to relate the factual knowledge of an agent (her vision of the environment), and defines the relationships between states of the world. As such, provability for beliefs does not generate modal literals. *Obligation rules* determine when and which obligations are in force. The conclusions generated by obligation rules are modalised with obligation. Finally, *outcome rules* establish the possible outcomes of an agent depending on the particular context. Apart from obligation rules, outcome rules are used to derive conclusions for all modes representing possible types of outcomes: desires, goals, intentions, and social intentions.

Following ideas given in [10], rules can gain more expressiveness when a *preference operator* \odot is used: an expression like $a \odot b$ means that if a is possible, then a is the first choice and b is the second one; if $\neg a$ holds, then the first choice is not attainable and b is the actual choice. This operator is used to build chains of preferences, called \odot -expressions. The formation rules for \odot -expressions are: (i.) every literal is an \odot -expression, (ii.) if A is an \odot -expression and b is a literal then $A \odot b$ is an \odot -expression. In addition we stipulate that \odot obeys to the following properties: (i.) $a \odot (b \odot c) = (a \odot b) \odot c$ (associativity); (ii.) $\odot_{i=1}^n a_i = (\odot_{i=1}^{k-1} a_i) \odot (\odot_{i=k+1}^n a_i)$ where exists j such that $a_j = a_k$ and $j < k$ (duplication and contraction on the right). \odot -expressions are given by the agent designer, or obtained through *construction rules* based on the particular logic [10].

In this paper we exploit the classical definition of *defeasible rule* in DL [1]. A defeasible rule is an expression $r : A(r) \Rightarrow_{\Box} C(r)$, where

1. $r \in \text{Lbl}$ is the name of the rule;
2. $A(r) = \{a_1, \dots, a_n\}$ with $a_i \in \text{Lit} \cup \text{ModLit}$ is the set of the premises (or the *antecedent*) of the rule;
3. $\Box \in \{\text{B}, \text{O}, \text{U}\}$ represents the *mode* of the rule (from now on, we omit the subscript B in rules for beliefs, i.e., \Rightarrow is used as a shortcut for \Rightarrow_{B});
4. $C(r)$ is the *consequent* (or *head*) of the rule, which is a single literal if $\Box = \text{B}$, or an \odot -expression otherwise².

¹ The reading of the modal operators is B for *belief*, O for *obligation*, D for *desire*, I for *intention* and SI for *social intention*.

² It is worth noting that modal literals can occur only in the antecedent of rules: the reason is that the rules are used to derive modal conclusions and we do not conceptually need to iterate modalities. The motivation of a single literal as a consequent for belief rules is dictated by the intended reading of the belief rules, where these rules are used to describe the environment.

We use the following abbreviations on sets of rules: R^\square ($R^\square[q]$) denotes all rules of mode \square (with consequent q), and $R[q] = \bigcup_{\square \in \{B, O, U\}} R^\square[q]$. $R[q, i]$ denotes the set of rules whose head is $\odot_{j=1}^n c_j$ and $c_i = q$, with $1 \leq i \leq n$.

Most of the terminology defined so far appears in [11], where an extension of DL with modal operators is introduced to differentiate modal and factual rules. However, labelling the rules of DL produces nothing more but a simple treatment of the modalities, thus two interaction strategies between modal operators are analysed.

Rule conversions It is sometimes meaningful to use rules for a modality X as they were for another modality Y , i.e., to convert one type of conclusions into a different one. For example, if ‘a car industry has the purpose of assembling perfectly working cars’ and ‘it is known that in every working car there is a working engine’, then ‘a car industry has also the purpose of assembling working engines in every car produced’. Formally, we define an asymmetric binary relation $\text{Convert} \subseteq \text{MOD} \times \text{MOD}$ such that $\text{Convert}(X, Y)$ means ‘a rule of mode X can be used also to produce conclusions of mode Y ’. This intuitively corresponds to the following logical schema:

$$\frac{Ya_1, \dots, Ya_n \quad a_1, \dots, a_n \Rightarrow_X b}{Yb} \text{Convert}(X, Y).$$

In our framework obligations and goal-like attitudes cannot change what the agent believes or how she perceives the world, thus we only consider conversion with mode for belief as the first element of the relation (i.e., $\text{Convert}(B, X)$ with $X \in \{O, D, G, I, SI\}$).

Conflict-detection/resolution It is crucial to identify criteria for detecting and solving conflicts between different modalities. Formally, we define an asymmetric binary relation $\text{Conflict} \subseteq \text{MOD} \times \text{MOD}$ such that $\text{Conflict}(X, Y)$ means ‘modes X and Y are in conflict and mode X prevails over Y ’. Consider the following theory:

$$\begin{aligned} F &= \{ \text{sunny_day}, \text{school_day} \}, \\ R &= \{ r_1 : \text{Sunny_day} \Rightarrow_U \text{go_outside}, r_2 : \text{school_day} \Rightarrow_O \neg \text{go_outside} \}. \end{aligned}$$

Even if there is a sunny day, a responsible parent would not go outside and play with her kid but will bring him to school; this behaviour is captured by $\text{Conflict}(O, SI)$, which means that the rule that forbids to go outside prevents the agent from obtaining the (social) intention of going outside, and the parent will not derive the (social) intention.

In our framework, we consider conflicts between beliefs and intentions, beliefs and social intentions, and obligations and social intentions. In other words, we have:

- $\text{Conflict}(B, I)$, $\text{Conflict}(B, SI)$ meaning that the agents are realistic (cf. [5]), and
- $\text{Conflict}(O, SI)$ meaning that the agents are social (cf. [11]).

Observation 1 *Convert and Conflict relations behave differently in our framework than the usual deployed in the literature [11]. Typically, there is a bijective correspondence between a mode and the type of rule “representing” it. For example, there are rules with mode O to derive obligations, or rules with mode I to derive intentions. This is not the case in our logic where outcome rules are used to derive conclusions for all goal-like attitudes. Thus, we can have $\text{Conflict}(O, SI)$ exhibiting the sociality of the agent but not $\text{Conflict}(O, U)$ since desires and obligation do not attack each other.*

There are two applications of the *superiority relation*: the first considers rules of the same mode; the latter compares rule of different mode. Given $r \in R^X$ and $s \in R^Y$, notice that $r > s$ iff r converts X into Y , or s converts Y into X , i.e., the superiority relation is used when rules, each with a different mode, are used to produce complementary conclusions of the same mode. Consider the following theory with $\text{Convert}(B, G)$:

$$\begin{aligned} F &= \{ \text{go_to_Rome}, \text{parent_anniversary}, \text{August} \}, \\ R &= \{ r_1 : \text{go_to_Rome} \Rightarrow_B \text{go_to_Italy} \\ &\quad r_2 : \text{parent_anniversary} \Rightarrow_U \text{go_to_Rome} \\ &\quad r_3 : \text{August} \Rightarrow_U \neg \text{go_to_Italy} \}, \\ &=> \{ (r_1, r_3) \}. \end{aligned}$$

Typically, I have the goal not to go to Italy in August since the weather is too hot and it is too crowded. However, it is my parents' anniversary and they are going to celebrate it this August in Rome, which is the capital of Italy. Nonetheless, I have the goal to go to Italy for my parents' wedding anniversary, since I am a good son. Here, the superiority applies because we use r_1 through a conversion from belief to goal.

2.2 Inferential Mechanism

A *proof* P of length n is a finite sequence $P(1), \dots, P(n)$ of *tagged literals* of the type $+\partial_X q$ and $-\partial_X q$, where $X \in \text{MOD}$. The proof conditions below define the logical meaning of such tagged literals. As a conventional notation, $P(1..i)$ denotes the initial part of the sequence P of length i . Given a defeasible theory D , $+\partial_X q$ means that q is defeasibly provable in D with the mode X , and $-\partial_X q$ that it has been proved in D that q is not defeasibly provable in D with the mode X . As usual, we use $D \vdash \pm \partial_\square l$ iff there is a proof P in D such that $P(n) = \pm \partial_\square l$ for an index n .

In order to characterise the notions of provability for beliefs ($\pm \partial_B$), obligations ($\pm \partial_O$), desires ($\pm \partial_D$), goals ($\pm \partial_G$), intentions ($\pm \partial_I$) and social intentions ($\pm \partial_{SI}$), it is essential to define when a rule is *applicable* or *discarded*. To this end, the preliminary notion of when a rule is *body-applicable/discarded* must be introduced, stating that each literal in the body of the rule must be proved/rejected with the suitable mode.

Definition 1. Let P be a proof and $\square \in \{O, D, G, I, SI\}$. A rule $r \in R$ is *body-applicable* (at step $n+1$) iff for all $a_i \in A(r)$:

1. if $a_i = \square l$ then $+\partial_\square l \in P(1..n)$,
2. if $a_i = \neg \square l$ then $-\partial_\square l \in P(1..n)$,
3. if $a_i = l \in \text{Lit}$ then $+\partial l \in P(1..n)$.

A rule $r \in R$ is *body-discarded* (at step $n+1$) iff there is $a_i \in A(r)$ such that

1. $a_i = \square l$ and $-\partial_\square l \in P(1..n)$, or
2. $a_i = \neg \square l$ and $+\partial_\square l \in P(1..n)$, or
3. $a_i = l \in \text{Lit}$ and $-\partial l \in P(1..n)$.

As already stated, belief rules allow us to derive literals with different modes. The applicability mechanism must take into account this constraint.

Definition 2. Let P be a proof. A rule $r \in R$ is 1. Conv-applicable, 2. Conv-discarded (at step $n + 1$) for X iff

1. $r \in R^B$, $A(r) \neq \emptyset$ and for all $a \in A(r)$, $+\partial_X a \in P(1..n)$;
2. $r \notin R^B$ or $A(r) = \emptyset$ or $\exists a \in A(r)$, $-\partial_X a \in P(1..n)$.

Let us consider the following theory

$$F = \{a, b, Oc\}, \quad R = \{r_1 : a \Rightarrow_O b, r_2 : b, c \Rightarrow d\},$$

r_1 is applicable, while r_2 is not since c is not proved as a belief. Instead, r_2 is *Conv-applicable* in the condition for $\pm\partial_O$, since Oc is a fact and r_1 proves Ob .

The notion of applicability gives guidelines on how to consider the next element in a given chain. Since a rule for belief cannot generate reparative chains but only single literals, we can conclude that the applicability condition for belief collapses into body-applicability. The same happens to desires, where we also consider the Convert relation. For obligations, each element before the current one must be a violated obligation. A literal is a candidate to be a goal only if none of the previous elements in the chain have been proved as a goal. For intentions, the elements of the chain must pass the wishful thinking filter, while social intentions are also constrained not to violate any norm.

Definition 3. Given a proof P , $r \in R[q, i]$ is applicable (at index i and step $n + 1$) for

1. B iff $r \in R^B$ and is body-applicable.
2. O iff either: (2.1.1) $r \in R^O$ and is body-applicable, (2.1.2) $\forall c_k \in C(r), k < i$, $+\partial_O c_k \in P(1..n)$ and $-\partial c_k \in P(1..n)$, or (2.2) r is Conv-applicable.
3. D iff either: (3.1) $r \in R^D$ and is body-applicable, or (3.2) Conv-applicable.
4. X, $X \in \{G, I, SI\}$ iff either: (4.1.1) $r \in R^X$ and is body-applicable, (4.1.2) $\forall c_k \in C(r), k < i$, $+\partial_Y \sim c_k \in P(1..n)$ for some Y such that $\text{Conflict}(Y, X)$ and $-\partial_X c_k \in P(1..n)$ or (4.2) r is Conv-applicable.

For G there are no conflicts; for I we have $\text{Conflict}(B, I)$, and for SI we have $\text{Conflict}(B, SI)$ and $\text{Conflict}(O, SI)$.

Conditions to establish that a rule is discarded correspond to the constructive failure to prove that the same rule is applicable, and follow the principle of *strong negation*.³

We can now describe the proof conditions for the various modal operators; we start with those for desires:

$+\partial_D$: If $P(n + 1) = +\partial_D q$ then

- (1) $Dq \in F$ or
- (2) (2.1) $\neg Dq \notin F$ and
 - (2.2) $\exists r \in R[q, i]$: r is applicable for D and
 - (2.3) $\forall s \in R[\sim q, j]$ either
 - (2.3.1) s is discarded for D , or
 - (2.3.2) $s \not\prec r$.

³ The strong negation principle is closely related to the function that simplifies a formula by moving all negations to an inner most position in the resulting formula, and replaces the positive tags with the respective negative tags, and the other way around [9].

We say that a *desire* is each element in a chain of an outcome rule for which there is no stronger argument for the opposite desire. The proof conditions for $+\partial_X$, with $X \in \{B, O, G, I, SI\}$ are as follows:

$+\partial_X$: If $P(n+1) = +\partial_X q$ then

- (1) $Xq \in F$ or
- (2) (2.1) $\neg Yq \notin F$ for $Y = X$ or $\text{Convert}(Y, X)$ and
 - (2.2) $\exists r \in R[q, i]$: r is applicable for X and
 - (2.3) $\forall s \in R^Y[\sim q, j]$ either
 - (2.3.1) s is discarded for Y , or
 - (2.3.2) $\exists t \in R^T[q, k]$: t is applicable for T and either
 - (2.3.2.1) $t > s$ if $Y = T$, $\text{Convert}(Y, T)$, or $\text{Convert}(T, Y)$; or
 - (2.3.2.2) $\text{Conflict}(T, Y)$.

To show that a literal q is defeasibly provable with modality X we have two choices: (1) modal literal Xq is a fact; or (2) we need to argue using the defeasible part of D . In this case, we require that a complementary literal (of the same modality, or of a conflictual modality) does not appear in the set of facts (2.1), and that there must be an applicable rule for q for mode X (2.2). Moreover, each possible attack brought by a rule s for $\sim q$ has to be either discarded (3.1), or successfully counterattacked by another stronger rule t for q (2.3.2). We recall that the superiority relation combines rules of the same mode, rules with different modes that produce complementary conclusion of the same mode through conversion (both considered in clause (2.3.2.1)), and conflictual modalities (clause 2.3.2.2). Obviously, if $\Box = B$, then the proof conditions reduce to those of classical defeasible logic [1].

Again, the negative counterparts ($-\partial_D$ and $-\partial_X$) are derived by strong negation applied to conditions for $+\partial_D$ and $+\partial_X$, respectively. As an example, consider the theory:

$$F = \{\neg b_1, O\neg b_2, SIb_4\} \quad R = \{r : \Rightarrow_U b_1 \odot b_2 \odot b_3 \odot b_4\}.$$

Then r is trivially applicable for D and $+\partial_D b_i$ holds, for $1 \leq i \leq 4$. Moreover, we have $+\partial_G b_1$ and r is discarded for G after b_1 . Since $+\partial\neg b_1$, $-\partial b_1$ holds (as well as $-\partial_{SI} b_1$); the rule is applicable for I and b_2 , and we are able to prove $+\partial b_2$, thus the rule becomes discarded for I after b_2 . Given that $O\neg b_2$ is a fact, r is discarded for SI and b_2 and $-\partial_{SI} b_2$ is proved, which in turn makes the rule applicable for SI at b_3 , proving $+\partial_{SI} b_3$. As we have argued before, this would make the rule discarded for b_4 . Nevertheless, b_4 is still provable with mode SI (in this case because it is a fact, but in other theories there could be more rules with b_4 in their head).

The logic resulting enjoys properties describing the appropriate behaviour of the modal operators.

Definition 4. A defeasible theory $D = (F, R, >)$ is consistent iff $>$ is acyclic and F does not contain pairs of complementary (modal) literals, that is pairs like (i.) l and $\sim l$, (ii.) $\Box l$ and $\neg \Box l$, $\Box \in \text{MOD}$, and (iii.) $\Box l$ and $\Box \sim l$, $\Box \in \text{MOD} \setminus \{D\}$.

Proposition 1 Let D be a consistent modal defeasible theory. For any literal l , it is not possible to have both

1. $D \vdash +\partial_{\Box} l$ and $D \vdash -\partial_{\Box} l$ with $\Box \in \text{MOD}$;

2. $D \vdash +\partial_{\Box} l$ and $D \vdash +\partial_{\Box} \sim l$ with $\Box \in \text{MOD} \setminus \{D\}$.

Moreover, given $\Box \in \text{MOD} \setminus \{D\}$, then:

3. if $D \vdash +\partial_{\Box} l$, then $D \vdash -\partial_{\Box} \sim l$.

3 Norm and Outcome Compliance

There are two fundamental strategies to characterize the concept of *norm compliance* in agent systems [7, 19]:

- Compliance is achieved by adopting the so-called norm regimentation strategy, which can amount to designing the system in such a way as illegal states are ruled out and made impossible in it, or by imposing that the occurrence of any illegal states is in theory possible but it leads to the system global failure;
- Norms are soft constraints and so do not limit in advance agents' behavior. Compliance is then ensured by system mechanisms stating that violations should result in sanctions or other normative effects which are supposed to recover from violations.

The second perspective is the one captured in this paper by devising rules for obligations that introduce preference orderings among normative outcomes: in fact, subsequent obligations in any \odot -expressions are meant to be countermeasures for the violation of previous obligations in the given sequence. In other words, \odot -constructions identify situations that are not ideal but still acceptable.

In BDI-like systems, norm compliance can be checked in the agents' deliberative stage against different types of motivational attitudes. In this paper, we model the case where obligations prevail over conflicting intentions in such a way that those intentions that are still provable are called *social intentions*. In this way, obligations are a mechanism for filtering expected outcomes in agents' minds.

On the other hand, since agents are oriented to the achievement of their own expected outcomes, and these outcomes are ranked as well in \odot -sequences, we may also speak of different types of *outcome compliance*, which means that, given the the environment where agents act, they deliberate and commit to achieve a state of the world where all the preferred outcomes are obtained. For example, if we want to go to the airport and not to be hungry, the actions of catching the train from our place to the airport, and stopping in a fast food to buy an hamburger, result in a state of the world where all our goals are fulfilled. In this perspective, \odot -sequences allow us to identify outcomes 'compensating' the non-achievement of other outcomes.

Informally, an agent is (a) *norm compliant* if each applicable norm is not violated, or if so, it is compensated, (b) *outcome compliant* if she deliberates to commit to at least all the minimal objectives, which means that there exists at least one possible way to satisfy at least one of all preferred outcomes in each \odot -sequence that is fired. Formally, those notions of compliance can be captured by establishing that all \odot -chains contain one last null element \perp , such that all of them have this general form: $b_1 \odot \dots \odot b_n \odot \perp$. With this done, various types of compliance are defined as follows:

Definition 5 (Norm and Outcome Compliance). *Let T be a defeasible theory. T is norm compliant if $T \vdash -\partial_{\odot} \perp$ and it is outcome compliant if $T \vdash -\partial_X \perp$, where $X \in \{G, I, SI\}$.*

4 Algorithmic Results

We now present the algorithms apt to compute the *extension* of a *finite* defeasible theory, i.e., with finite set of facts and rules, in order to bind the complexity of the logic introduced in the previous sections. The algorithms are inspired by ideas of [15, 14].

For the sake of clarity, from now on \blacksquare denotes a generic mode in MOD, \diamond a generic mode in $\text{MOD} \setminus \{B\}$, and \square a fixed mode chosen in \blacksquare . Moreover, we will treat literals $\square l$ and l as synonyms whenever $\square = B$. To accommodate the Convert relation to the algorithms, we denote with $R^{B, \diamond}$ the set of belief rules with non-empty body that can be used for a conversion to mode \diamond . Furthermore, for each literal l , l_{\blacksquare} is the set (initially empty) such that $\pm \square \in l_{\blacksquare}$ iff $D \vdash \pm \partial_{\square} l$. Given a modal defeasible theory D , a set of rules R , and a rule $r \in R^{\square}[l]$, we expand $>$ by incorporating the Conflict relation into it; this will ease the computation. Then, we define: (i.) $r_{sup} = \{s \in R : (s, r) \in >\}$ and $r_{inf} = \{s \in R : (r, s) \in >\}$ for any $r \in R$; (ii.) HB_D as the set of literals such that the literal or its complement appears in D , where ‘appears’ means that it is a sub-formula of a modal literal occurring in D ; (iii.) the modal Herbrand Base of D as $HB = \{\square l \mid \square \in \text{MOD}, l \in HB_D\}$. Accordingly, the extension of a defeasible theory is defined as follows.

Definition 6. *Given a modal defeasible theory D , the defeasible extension of D is defined as $E(D) = (+\partial_{\square}, -\partial_{\square})$ where $\pm \partial_{\square} = \{l \in HB_D : D \vdash \pm \partial_{\square} l\}$ with $\square \in \text{MOD}$. Two defeasible theories D and D' are equivalent whenever $E(D) = E(D')$.*

The next definition extends the concept of complement presented in Section 2 for modal literals and establishes the logical connection among proved and refuted literals.

Definition 7. *The complement of a given modal literal l , denoted by \tilde{l} , is:*

1. *if $l = Dm$, then $\tilde{l} = \{\neg Dm\}$;*
2. *if $l = \square m$, then $\tilde{l} = \{\neg \square m, \square \sim m\}$, with $\square \in \{O, G, I, SI\}$;*
3. *if $l = \neg \square m$, then $\tilde{l} = \{\square m\}$.*

Truncation and removal are two syntactical operations on the consequent of rules.

Definition 8. *Let $c_1 = a_1 \odot \dots \odot a_{i-1}$ and $c_2 = a_{i+1} \odot \dots \odot a_n$ be two (possibly empty) \odot -expressions such that a_i does not occur in them, and $c = c_1 \odot a_i \odot c_2$ is an \odot -expression. Let r be a rule with form $A(r) \Rightarrow_X c$. We define the*

- *truncation of the consequent c at a_i as $A(r) \Rightarrow_X c!a_i = A(r) \Rightarrow_X c_1 \odot a_i$;*
- *removal of a_i from the consequent c as $A(r) \Rightarrow_X c \ominus a_i = A(r) \Rightarrow_X c_1 \odot c_2$.*

Given $\square \in \text{MOD}$, the sets $\pm \partial_{\square}$ denote the global sets of defeasible conclusions (i.e., the set of literals for which condition $\pm \partial_{\square}$ holds), while ∂_{\square}^{\pm} are the corresponding temporary sets. Moreover, to simplify the calculus we do not operate on outcome rules: for each rule $r \in R^U$ we create instead a new rule for all the other goal-like modes (resp. r^D , r^G , r^I , and r^{SI}). Accordingly, we will use expressions like “the intention rule” as a shortcut for “the clone of outcome rule used to derive intentions”.

The idea of all algorithms is to use the operations of truncation and elimination in order to obtain, step after step, a simpler but equivalent theory. Indeed, proving a literal

does not give just local information about the element itself, but reveals which rules will be applicable, discarded, or reduced in their head or tail.

Observation 2 Assume that, at a given step, the algorithm proves l . At the next step,

1. the applicability of any rule r with l in its antecedent $A(r)$ does not depend on l any longer. Accordingly, we can safely remove l from $A(r)$.
2. Any rule s where \bar{l} is in its antecedent $A(s)$ is discarded. Consequently, any superiority tuple involving this rule is now meaningless and can be removed from the superiority relation as well.
3. We can shorten chains by exploiting conditions of Definition 3. For example, if $l = Om$, we can truncate chains for obligations at $\sim m$ and eliminate $\sim m$.

Algorithm 1 DEFEASIBLEEXTENSION

```

1:  $+\partial_{\blacksquare}, \partial_{\blacksquare}^+ \leftarrow \emptyset; -\partial_{\blacksquare}, \partial_{\blacksquare}^- \leftarrow \emptyset$ 
2:  $R \leftarrow R \cup \{r^\square : A(r) \Rightarrow_\square C(r) | r \in R^U\} \setminus R^U$ , with  $\square \in \{D, G, I, SI\}$ 
3:  $R^{B,\diamond} \leftarrow \{r^\diamond : \diamond a_1, \dots, \diamond a_n \Rightarrow_\diamond C(r) | r \in R^B, A(r) \neq \emptyset, A(r) \subseteq \text{Lit}, a_i \in A(r)\}$ 
4:  $>\leftarrow > \cup \{(r^\diamond, s^\diamond) | r^\diamond, s^\diamond \in R^{B,\diamond}, r > s\} \cup \{(r, s) | r \in R^\blacksquare, s \in R^\diamond \cup R^{B,\diamond}, \text{Conflict}(\blacksquare, \diamond)\}$ 
5: for  $l \in F$  do
6:   if  $l = \square m$  then  $\text{PROVED}(m, \square)$ 
7:   if  $l = \neg \square m \wedge \square \neq D$  then  $\text{REFUTED}(m, \square)$ 
8: end for
9:  $+\partial_{\blacksquare} \leftarrow +\partial_{\blacksquare} \cup \partial_{\blacksquare}^+; -\partial_{\blacksquare} \leftarrow -\partial_{\blacksquare} \cup \partial_{\blacksquare}^-$ 
10:  $R_{\text{inf}} \leftarrow \emptyset$ 
11: repeat
12:    $\partial_{\blacksquare}^+ \leftarrow \emptyset; \partial_{\blacksquare}^- \leftarrow \emptyset$ 
13:   for  $\square l \in HB$  do
14:     if  $R^\square[l] \cup R^{B,\square}[l] = \emptyset$  then  $\text{REFUTED}(l, \square)$ 
15:   end for
16:   for  $r \in R^\square \cup R^{B,\square}$  do
17:     if  $A(r) = \emptyset$  then
18:        $r_{\text{inf}} \leftarrow \{r \in R : (r, s) \in >, s \in R\}; r_{\text{sup}} \leftarrow \{s \in R : (s, r) \in >\}$ 
19:        $R_{\text{inf}} \leftarrow R_{\text{inf}} \cup r_{\text{inf}}$ 
20:       Let  $l$  be the first literal of  $C(r)$  in  $HB$ 
21:       if  $r_{\text{sup}} = \emptyset$  then
22:         if  $\square = D$  then
23:            $\text{PROVED}(m, D)$ 
24:         else
25:            $\text{REFUTED}(\sim l, \square)$ 
26:            $\text{REFUTED}(\sim l, \diamond)$  for  $\diamond$  s.t.  $\text{Conflict}(\square, \diamond)$ 
27:           if  $R^\square[\sim l] \cup R^{B,\square}[\sim l] \cup R^\blacksquare[\sim l] \setminus R_{\text{inf}} \subseteq r_{\text{inf}}$ , for  $\blacksquare$  s.t.  $\text{Conflict}(\blacksquare, \square)$  then
28:              $\text{PROVED}(m, \square)$ 
29:           end if
30:         end if
31:       end if
32:     end if
33:   end for
34:    $\partial_{\blacksquare}^+ \leftarrow \partial_{\blacksquare}^+ \setminus +\partial_{\blacksquare}; \partial_{\blacksquare}^- \leftarrow \partial_{\blacksquare}^- \setminus -\partial_{\blacksquare}$ 
35:    $+\partial_{\blacksquare} \leftarrow +\partial_{\blacksquare} \cup \partial_{\blacksquare}^+; -\partial_{\blacksquare} \leftarrow -\partial_{\blacksquare} \cup \partial_{\blacksquare}^-$ 
36: until  $\partial_{\blacksquare}^+ = \emptyset$  and  $\partial_{\blacksquare}^- = \emptyset$ 
37: return  $(+\partial_{\blacksquare}, -\partial_{\blacksquare})$ 

```

Algorithm 1 DEFEASIBLEEXTENSION is the core to compute the extension of a defeasible theory. The first part (lines 1–4) sets up the data structure needed for the computation. Lines 5–8 are to handle facts as immediately provable literals. The main idea of the algorithm is to check whether there are rules whose body is empty. Since defeasible rules can have \odot -expressions as their head, the literal we are interested in is the first element of the \odot -expression (loop **for** at lines 16–33 and **if** condition at line 17). Such rules are clearly applicable and they can produce conclusions with the right modality. However, before asserting that the first element of the conclusion is provable, we have to check whether there are no rules for the complement (again with the appropriate mode), otherwise such rules for the complement must be weaker than the applicable rules. This information is stored in R_{inf_d} inspired by the technique of [14]. If no rule stronger than the current one exists, the complementary conclusion must be refuted by condition (2.3) of $-\partial_{\Box}$ (line 25). A straightforward consequence of $D \vdash -\partial_{\Box} l$ is that literal l is also refutable in D with any modality conflicting with \Box (line 26). Notice that this reasoning does not hold for desires: since the logic allows to have Dl and $D\sim l$ at the same time, when $\Box = D$ the algorithm directly invokes procedure 2 PROVED (line 23).

The next step is to check whether there exist rules for the complement of the literal with the same (or conflicting) mode. The rules for the complement should not be defeated by an applicable rule, i.e., they should not be in R_{inf_d} . If all these rules are defeated by r (line 27), then conditions for deriving $+\partial_{\Box}$ are satisfied. If a literal is assessed to be provable (with the appropriate modality) the algorithm calls procedure 2 PROVED, otherwise the procedure 3 REFUTED is invoked. The algorithm finally returns the extension of the input theory when no modifications are done on sets $\partial_{\blacksquare}^{\pm}$.

Algorithm 2 PROVED is invoked when literal l is proved with modality \Box . The computation starts by updating the relative positive extension set for modality \Box and the local information on literal l (line 2); l is then removed from HB at line 3. Proposition 1 Part 3. defines the modes with which literal $\sim l$ can be refuted (**if** condition at line 4). Lines 5 to 7 modifies the sets of rules R and $R^{B, \Box}$, and the superiority relation accordingly to ideas of Observation 2.

Depending on the modality \Box of l , we have to perform some specific operations on chains (condition **switch** at lines 8–27). Entering into the detail of each **case** would be redundant without giving more information than conditions of a rule being applicable or discarded in Section 2. Therefore, we propose one significative example by considering the scenario where l has been proved as a belief (**case** at lines 9–13). Here, chains of obligation (resp. intention) rules can be truncated after l since are discarded for all following elements (line 10). Analogously, condition (4.1.2) of Definition 3 allows to eliminate $\sim l$ from intention and social intention rules (line 11). If $+\partial_O \sim l$ has been already proved, then we eliminate $\sim l$ since it represents a violated obligation. Vice versa, if $-\partial_O \sim l$ is the case, then each element after l cannot be a social intention (**if** conditions at lines 12 and 13, respectively).

Algorithm 3 REFUTED performs all necessary operations in case literal l is refuted with mode \Box . The initialisation steps at lines 2–6 follow the same schema exploited at lines 2–7 of Algorithm 2 PROVED. Again, the operations to be performed on chains vary according to the current mode \Box (**switch** at lines 7–19). For example, if $\Box = B$ (lines 8–

Algorithm 2 PROVED

```
1: procedure PROVED( $l \in \text{Lit}, \square \in \text{MOD}$ )
2:    $\partial_{\square}^+ \leftarrow \partial_{\square}^+ \cup \{l\}; l_{\blacksquare} \leftarrow l_{\blacksquare} \cup \{+\square\}$ 
3:    $HB \leftarrow HB \setminus \{\square l\}$ 
4:   if  $\square \neq \text{D}$  then REFUTED( $\sim l, \square$ )
5:    $R \leftarrow \{r : A(r) \setminus \{\square l, \neg \square \sim l\} \hookrightarrow C(r) \mid r \in R, A(r) \cap \widetilde{\square l} = \emptyset\}$ 
6:    $R^{\text{B}, \square} \leftarrow \{r : A(r) \setminus \{\square l\} \hookrightarrow C(r) \mid r \in R^{\text{B}, \square}, A(r) \cap \widetilde{\square l} = \emptyset\}$ 
7:    $\triangleright \triangleleft \triangleright \setminus \{(r, s), (s, r) \in \triangleright \mid A(r) \cap \widetilde{\square l} \neq \emptyset\}$ 
8:   switch ( $\square$ )
9:     case B:
10:       $R^X \leftarrow \{A(r) \Rightarrow_X C(r)! \mid r \in R^X[l, n]\}$  with  $X \in \{\text{O}, \text{I}\}$ 
11:       $R^X \leftarrow \{A(r) \Rightarrow_X C(r) \ominus \sim l \mid r \in R^X[\sim l, n]\}$  with  $X \in \{\text{I}, \text{SI}\}$ 
12:      if  $+\text{O} \in \sim l_{\blacksquare}$  then  $R^{\text{O}} \leftarrow \{A(r) \Rightarrow_{\text{O}} C(r) \ominus \sim l \mid r \in R^{\text{O}}[\sim l, n]\}$ 
13:      if  $-\text{O} \in \sim l_{\blacksquare}$  then  $R^{\text{SI}} \leftarrow \{A(r) \Rightarrow_{\text{SI}} C(r)! \mid r \in R^{\text{SI}}[l, n]\}$ 
14:     case O:
15:       $R^{\text{O}} \leftarrow \{A(r) \Rightarrow_{\text{O}} C(r)! \ominus l \ominus \sim l \mid r \in R^{\text{O}}[\sim l, n]\}$ 
16:       $R^{\text{SI}} \leftarrow \{A(r) \Rightarrow_{\text{SI}} C(r) \ominus \sim l \mid r \in R^{\text{SI}}[\sim l, n]\}$ 
17:      if  $-\text{B} \in l_{\blacksquare}$  then  $R^{\text{O}} \leftarrow \{A(r) \Rightarrow_{\text{O}} C(r) \ominus l \mid r \in R^{\text{O}}[l, n]\}$ 
18:      if  $-\text{B} \in \sim l_{\blacksquare}$  then  $R^{\text{SI}} \leftarrow \{A(r) \Rightarrow_{\text{SI}} C(r)! \mid r \in R^{\text{SI}}[l, n]\}$ 
19:     case D:
20:      if  $+\text{D} \in \sim l_{\blacksquare}$  then
21:         $R^{\text{G}} \leftarrow \{A(r) \Rightarrow_{\text{G}} C(r)! \ominus l \mid r \in R^{\text{G}}[l, n]\}$ 
22:         $R^{\text{G}} \leftarrow \{A(r) \Rightarrow_{\text{G}} C(r)! \sim l \ominus \sim l \mid r \in R^{\text{G}}[\sim l, n]\}$ 
23:      end if
24:     otherwise:
25:       $R^{\square} \leftarrow \{A(r) \Rightarrow_{\square} C(r)! \mid r \in R^{\square}[l, n]\}$ 
26:       $R^{\square} \leftarrow \{A(r) \Rightarrow_{\square} C(r) \ominus \sim l \mid r \in R^{\square}[\sim l, n]\}$ 
27:   end switch
28: end procedure
```

Algorithm 3 REFUTED

```
1: procedure REFUTED( $l \in \text{Lit}, \square \in \text{MOD}$ )
2:    $\partial_{\square}^- \leftarrow \partial_{\square}^- \cup \{l\}; l_{\blacksquare} \leftarrow l_{\blacksquare} \cup \{-\square\}$ 
3:    $HB \leftarrow HB \setminus \{\square l\}$ 
4:    $R \leftarrow \{r : A(r) \setminus \{\neg \square l\} \hookrightarrow C(r) \mid r \in R, \square l \notin A(r)\}$ 
5:    $R^{\text{B}, \square} \leftarrow R^{\text{B}, \square} \setminus \{r \in R^{\text{B}, \square} : \square l \in A(r)\}$ 
6:    $\triangleright \triangleleft \triangleright \setminus \{(r, s), (s, r) \in \triangleright \mid \square l \in A(r)\}$ 
7:   switch ( $\square$ )
8:     case B:
9:       $R^{\text{I}} \leftarrow \{A(r) \Rightarrow_{\text{I}} C(r)! \sim l \mid r \in R^{\text{I}}[\sim l, n]\}$ 
10:      if  $+\text{O} \in l_{\blacksquare}$  then  $R^{\text{O}} \leftarrow \{A(r) \Rightarrow_{\text{O}} C(r) \ominus l \mid r \in R^{\text{O}}[l, n]\}$ 
11:      if  $-\text{O} \in l_{\blacksquare}$  then  $R^{\text{SI}} \leftarrow \{A(r) \Rightarrow_{\text{SI}} C(r)! \sim l \mid r \in R^{\text{SI}}[\sim l, n]\}$ 
12:     case O:
13:       $R^{\text{O}} \leftarrow \{A(r) \Rightarrow_{\text{O}} C(r)! \ominus l \mid r \in R^{\text{O}}[l, n]\}$ 
14:      if  $-\text{B} \in l_{\blacksquare}$  then  $R^{\text{SI}} \leftarrow \{A(r) \Rightarrow_{\text{SI}} C(r)! \sim l \mid r \in R^{\text{SI}}[\sim l, n]\}$ 
15:     case D:
16:       $R^X \leftarrow \{A(r) \Rightarrow_X C(r) \ominus l \mid r \in R^X[l, n]\}$  with  $X \in \{\text{D}, \text{G}\}$ 
17:     otherwise:
18:       $R^{\square} \leftarrow \{A(r) \Rightarrow_{\square} C(r) \ominus l \mid r \in R^{\square}[l, n]\}$ 
19:   end switch
20: end procedure
```

11), then applicability condition (4.1.2) for $\pm\partial_1$ cannot be satisfied for any literal after $\sim l$ in chains for intentions, and such chains can be truncated at $\sim l$. Furthermore, if the algorithm has already proven $+\partial_O l$, then l represents a violated obligation. Thus, l can be removed from all chains for obligations. If instead $-\partial_O l$ holds, then the elements after $\sim l$ in chains for social intentions do not satisfy applicability condition (4.1.2) of $\pm\partial_{S_I}$ for $\sim l$, and the algorithm removes them.

We conclude by showing the computational properties of the algorithms proposed.

Theorem 1. *Algorithm 1 DEFEASIBLEEXTENSION terminates and its computational complexity is $O(|R| * |HB|)$.*

Proof Sketch. Termination is ensured since at every iteration either no modification occurs and line 36 ends the computation, or a literal is removed from HB . This set is finite, since the sets of facts and rules are finite, thus eventually the process empties HB . This bounds also the complexity to the number of rules and literals in HB , since each modal literal is processed once and every time we scan the set of rules.

Theorem 2. *Algorithm 1 DEFEASIBLEEXTENSION is sound and complete.*

5 Conclusions and Related Work

We provided a fresh characterisation for motivational states as the concept of goals, intentions, and social intentions obtained through a deliberative process based on various types of preferences among desired outcomes. In this sense, this contribution has strong connections with [6, 11, 9] but presents significant improvements in at least two respects. First, while in those works the agent deliberation is the result of the derivation of mental states from *precisely* the corresponding rules of the logic, here the proof theory is more aligned with the BDI intuition, according to which intentions and goals are the results of the manipulation of desires. This allow us to encode this idea within a logical language and a proof theory, by exploiting the different interaction patterns between the basic mental states, as well as the derived ones. Hence, our framework is more expressive than BOLD [5], which uses different rules to derive the corresponding mental states and proposes simple criteria to solve conflicts between rule types.

Second, the framework proposes a rich language expressing two orthogonal concepts of preference among motivational attitudes. One is encoded within \odot sequences, which state reparative orders among homogeneous mental states or motivations, and which are contextual. The second type of preference is encoded via the superiority relation between rules which can work locally, as well as via the Conflict relation. The interplay between these mechanisms can help us to isolate complex ways for deriving mental states, but the resulting logical machinery is still computationally tractable.

Finally, since the preferences allow us to determine what preferred outcomes can be chosen by an agent (in a specific scenario) when previous goals in \odot -sequences are not (or no longer) feasible, our logic in fact provides an abstract semantics for several types of goal and intention reconsideration. Intention reconsideration was expected to play a crucial role in the BDI paradigm since intentions obey the law of inertia and resist retraction or revision, but they can be reconsidered when new relevant information

comes in [3]. Despite that, the problem of revising intentions in BDI frameworks has received little attention. A very sophisticated exception is [20], where revisiting intentions mainly depends on the dynamics of beliefs but the process is incorporated in a very complex framework for reasoning about mental states. Recently, [17] discussed how to revise the commitments to planned activities because of mutually conflicting intentions, which interestingly has connections with our work. How to employ our logic to give a semantics for intention reconsideration is not the main goal of the paper and is left to future work.

References

1. G. Antoniou, D. Billington, G. Governatori, and M. J. Maher. Representation results for defeasible logic. *ACM Transactions on Computational Logic*, 2(2):255–287, 2001.
2. N. Bassiliades, G. Antoniou, and I. P. Vlahavas. A defeasible logic reasoner for the semantic web. *Int. J. Semantic Web Inf. Syst.*, 2(1):1–41, 2006.
3. M. Bratman. *Intentions, Plans and Practical Reason*. Harvard University Press, 1987.
4. M. Bratman, D. Israel, and M. Pollack. Plans and resource-bounded practical reasoning. *Computational Intelligence*, 4:349–355, 1988.
5. J. Broersen, M. Dastani, J. Hulstijn, and L. van der Torre. Goal generation in the BOID architecture. *Cognitive Science Quarterly*, 2(3-4):428–447, 2002.
6. M. Dastani, G. Governatori, A. Rotolo, and L. van der Torre. Programming cognitive agents in defeasible logic. In *Proc. LPAR 2005*, 2005.
7. A. Garcia-Camino, J. Rodriguez-Aguilar, C. Sierra, and W. Vasconcelos. Constraint rule-based programming of norms for electronic institutions. *JAAMAS*, 18:186–217, 2009.
8. G. Governatori. Representing business contracts in RuleML. *International Journal of Co-operative Information Systems*, 14(2-3):181–216, 2005.
9. G. Governatori, V. Padmanabhan, A. Rotolo, and A. Sattar. A defeasible logic for modelling policy-based intentions and motivational attitudes. *Logic Journal of the IGPL*, 17(3):227–265, 2009.
10. G. Governatori and A. Rotolo. Logic of violations: A Gentzen system for reasoning with contrary-to-duty obligations. *Australasian Journal of Logic*, 4:193–215, 2006.
11. G. Governatori and A. Rotolo. BIO logical agents: Norms, beliefs, intentions in defeasible logic. *Autonomous Agents and Multi-Agent Systems*, 17(1):36–69, 2008.
12. K. Kravari, C. Papatheodorou, G. Antoniou, and N. Bassiliades. Reasoning and proofing services for semantic web agents. In *Proc. IJCAI 2011*, 2011.
13. H.-P. Lam and G. Governatori. The making of SPINdle. In *Proc. RuleML 2009*, 2009.
14. H.-P. Lam and G. Governatori. What are the Necessity Rules in Defeasible Reasoning? In *Proc. LPNMR-11*. 2011.
15. M. J. Maher. Propositional defeasible logic has linear complexity. *Theory and Practice of Logic Programming*, 1(6):691–711, 2001.
16. A. S. Rao and M. P. Georgeff. Modelling rational agents within a BDI-architecture. In *Proc. KR'91*, 1991.
17. S. Shapiro, S. Sardina, J. Thangarajah, L. Cavedon, and L. Padgham. Revising conflicting intention sets in bdi agents. In *Proc. AAMAS '12*, 2012.
18. I. Tachmazidis, G. Antoniou, G. Flouris, S. Kotoulas, and L. McCluskey. Large-scale parallel stratified defeasible reasoning. In *Proc. ECAI 2012*, 2012.
19. N. Tinnemeier, M. Dastani, and J.-J. Meyer. Roles and norms for programming agent organizations. In *Proc. AAMAS '09*, 2009.
20. W. van der Hoek, W. Jamroga, and M. Wooldridge. Towards a theory of intention revision. *Synthese*, 155(2):265–90, 2007.