

# Interaction between Normative Systems and Cognitive Agents in Temporal Modal Defeasible Logic

Regis Riveret<sup>1</sup>, Antonino Rotolo<sup>1</sup>, Guido Governatori<sup>3</sup>

<sup>1</sup> CIRSIFD and Law Faculty  
Via Galliera 3, 40121, Bologna, Italy  
[rriveret@cirsfid.unibo.it](mailto:rriveret@cirsfid.unibo.it)  
[rotolo@cirsfid.unibo.it](mailto:rotolo@cirsfid.unibo.it)

<sup>2</sup> University of Queensland  
Brisbane, Queensland, QLD 4072, Australia  
[guido@itee.uq.edu.au](mailto:guido@itee.uq.edu.au)

**Abstract.** While some recent frameworks on cognitive agents addressed the combination of mental attitudes with deontic concepts, they commonly ignore the representation of time. An exception is [1], which also manages some temporal aspects with regard to both cognition and deontic provisions. We propose in this paper a variant of the logic presented in [1] to deal in particular with temporal intervals.

**Keywords.** Time, Norm, Temporal Modal Defeasible Logic

## 1 Introduction

A common approach in the agent literature for programming cognitive agents in a BDI (belief, desire, intention) framework is the use of rules to represent or manipulate the agents mental attitudes. In addition to the three mental attitudes of beliefs, desires and intentions, some works include deontic concepts to denote norms, commitments of social agents and social rationality [2,3,4,5,6]. However, these frameworks commonly ignore the representation of time. An exception is [1], which adopts the rule-based approach of [7,8,9] and extends it to accommodate temporal aspects. Time is integrated by pairing assertions with instants representing the time at which assertions hold and by discriminating transient and persistent conclusions. Persistent conclusions persists until some interrupting event occurs. Pairing assertions with instants is unsatisfactory for at least two reasons: (i) some properties may end at a certain time not associated to any explicit external event, (ii) we may like to represent rules where conditions have to hold for certain temporal intervals. To remedy these issues, in this paper, we increase the expressive power of the logic presented in [1] with temporal intervals. The framework presented is based on Temporal Defeasible Logic (TDL), an umbrella expression designating extensions of Defeasible Logic to capture time. Beside [1], TDL has proved useful in modelling temporal aspects of normative reasoning, such as temporalised normative provisions [10]; in addition, the notion of temporal viewpoints -the temporal positions from which things are viewed- allows for a logical account of retroactive norms and norm modifications [11].

The paper is organised as follows. In section 2, we introduce the general conceptual model behind the framework. Section 3 provides an outline of basic Defeasible Logic. Section 4 describes a variant of modal TDL that formalises the model of cognition.

## 2 Time, norms and mental attitudes

Our model aims to give an account of some temporal aspects with regard to both mental attitudes and deontic provisions. The starting point is the acknowledgement that, on the one hand, recent works shows that reasoning about agents can be embedded in frameworks based on non-monotonic logic, as the most interesting problems concern cases where the agent’s mental attitudes are in conflict or when they are incompatible with deontic provisions. On the other hand, in a temporal setting, non-monotonicity can also be used to conclude that mental attitudes or deontic provisions persist up to some future time unless there is a reason for it not to persist. One can thus argue that a type of non-monotonicity concerns situations where mental attitudes are in conflict or when they are incompatible with some deontic provisions, while another type of non-monotonicity concerns temporal aspects. Our model is based on these two types of non-monotonicity.

We adopt the model of [1] that extends the works of [7,8,9] with time. These later works are themselves inspired by Bratman’s analysis of so-called policy-based attitudes. In Bratman’s view intentions are used to choose partial plans for realisation of a goal and have a close relation to mean-ends, whereas [7,8,9] intentions are related not only to means-ends but also to their consequences. This notion is particularly relevant with deontic and normative notions, for example if we want to say that an agent is legally for A if the A is a side effect and if the agent did A with the intention to do A. [7,8,9] extends this policy-based approach to other attitudes and motivational factors as beliefs, intentions and obligations. An agent types correspond to the different ways through which conflicts are detected and solved: a realistic agent thus corresponds to a conflict-resolution type in which beliefs override all other factors, while other agent types, such as simple minded, selfish or social ones adopt different orders of overruling.

[1] is on the same line of research of [7,8,9] and focus on some temporal aspects. [1] is based on Bratman’s [12] which in his pursuit for a temporally extended rational agency exposed a principle that can be roughly stated as follows:

- At  $t_0$ , agent A consider the policy to adopt with respect a certain range of activities. On this basis, agent A forms a general intention to  $\varphi$  in circumstances of type  $\psi$ .
- From  $t_0$  to  $t_1$ , A retains this general intention.
- At  $t_1$ , A notes that he/she is or will be in circumstance  $\psi$  at  $t_2$ .
- Based on the previous steps, A forms the intention at  $t_1$  to  $\varphi$  at  $t_2$ .

Given the temporal nature of Bratman’s historical principle, and the idea that some intentions can be retained from one moment to another, [1] accounts for two types of temporal intentions: transient intentions which hold only for an instant of time, and persistent intentions which an agent is going to retain unless some interrupting event occurs that forces the agent to reconsider them. This event can be just a brute fact or it can be a modification of the policy of the agent.

The expressive power of [1] is unsatisfactory for at least two reasons. First some properties may end at a certain time not associated to any explicit external event, for example, an obligation or a norm in force may hold until some specific temporal reference. Secondly, we may like to represent rules where conditions have to hold for certain temporal intervals. To remedy these issues, in this paper, we increase the expressive power of the logic presented in [1] with temporal intervals.

Ordinarily, intervals are defined as sets of instants between two indicated instants. Doing so, some difficulties may arise when we want to express that an event (for example) occurs in an interval. This refers to the non-homogeneity or transient character of events: if an event occurs in an interval conceived as a set of instants, then it would also occur in the set of instants that defines it and this would conflict with the transient characterisation of events. Hence, we deviate somewhat to the standart definition of intervals as a set of instants, and define an interval as a pair of instants of the form  $[t_i, t_f]$  and usually denote them by  $T$  (plus eventual subscript). We identify two subsets of interval to differentiate intervals in which an associated property holds at any instant between the boundaries and intervals in which an associated property holds at least one instant between the boundaries. We shall call the firsts A-interval and the seconds B-intervals. A-intervals are represented by expressions of the form  $\overline{[t_i, t_f]}$  and are usually denoted by  $\overline{T}$  while B-intervals are represented by expressions of the form  $\widehat{[t_i, t_f]}$  and denoted by  $\widehat{T}$ . If the wide hat or the line over an interval is omitted then it is either an A-interval or a B-interval.

Mental attitudes and normative provisions are related to temporal references and the passage of time allows change of these elements. This is in accordance with the commonly accepted opinion that in a static system where nothing changes, the temporal dimension does not provide more understanding. Our references are intervals and allows us to temporalise literals and rules. In its simplest form, a temporal literal is an expression of the form  $l:T$  where  $l$  is a literal and  $T$  is either an A-interval or a B-interval. Intuitively,  $l:\overline{T}$  means that  $l$  holds for all instants between the boundaries of  $\overline{T}$  while  $l:\widehat{T}$  means that  $l$  holds for at least an instant between the boundaries of  $\widehat{T}$ . For example,  $adult(bob):\overline{[1973, max]}$  means that Bob has legally reach adulthood in 1973. Similarly, rules are temporalised by associating to it a time interval, and so a temporal rule is an expression of the form:

$$(r: a_1:T_1 \dots a_n:T_n \leftrightarrow b:\overline{T}): \overline{T}_r$$

The time labels allow us to deal formally with the different temporal dimensions of a normative system. The temporal intervals labelling the antecedent of a rule, the consequent of the rule and the overall rule are interpreted respectively as the intervals of *efficacy*, *applicability* and *time of force* of the represented provision. These different temporal dimensions are in line with the legal temporal model developed in [13], and that allows us to give an accurate account of temporal aspects of norms and therefore to be consistent with legal principles. Note that the interval  $\overline{T}_r$  labelling the entire rule is an A-interval because the force of a provision is generally an homogeneous property. Similarly, we constraint for the sake of simplicity the interval labelling the literal in the head of the rule to be an A-interval. Intervals in the body can be A-intervals or B-intervals. An example of a temporal rule is:

$$(r: \text{born}(X):[t, t] \rightarrow \text{major}(X):[\overline{t + 18, \text{max}}]):[\overline{1970, \text{max}}]$$

This rule formalises the provision in force in 1970 and later that people legally reach adulthood at 18. Consequently of the different temporal dimensions, a conclusion can be associated to two temporal intervals. The first interval is the interval with which the consequent of the rule is labelled while the second interval corresponds to the time of force interval associated to the rule. We represent such temporalisation of conclusion by concatenation of intervals by means of the symbol ':' and we call such concatenation chain of viewpoints. For example, giving the rule  $r$  and the fact that Bob was born in 1960, then one can conclude  $\text{major}(\text{bob}):[1978, \text{max}]:[1970, \text{max}]$ , that is, Bob is legally adult in 1978 (and later) from somebody reasoning in 1970 (and later).

Chain of viewpoints are of the utmost importance when one has to deal with the retroactivity of norms. Retroactivity usually occurs when the effects of a rule  $r$  apply to an interval  $[t_i, t_f]$  which begins before the interval  $[t'_i, t'_f]$  attached to the antecedent of  $r$ , that is,  $t_i < t'_i$ . Another case of retroactivity is when the consequence of a rule  $r'$  in force in  $[t_{ri}, t_{rf}]$  has as interval of applicability  $[t_i, t_f]$  and  $t_i < t_{ri}$ . For an illustration of the utility of chain of viewpoints with respects to retroactivity, consider the following rules:

$$(r1: \text{Income} > 90:[\widehat{1\text{Mar}06, 1\text{Jun}06}] \Rightarrow_{\text{OBL}} \neg \text{Tax}:[\overline{1\text{Jan}06, 1\text{Jun}06}]):[\overline{15\text{Jan}06, 1\text{Jun}06}]$$

$$(r2: \text{Income} > 100:[\widehat{1\text{Mar}06, 1\text{Jun}06}] \Rightarrow_{\text{OBL}} \text{Tax}:[\overline{1\text{Jan}06, 1\text{Jun}06}]):[\overline{1\text{Apr}06, 1\text{Jun}06}]$$

Rule  $r1$  states that if the income of a person is in excess of ninety thousand between the 1st March 2006 and the 1st June 2006 then she has not to pay the tax from 1st January 2006 to 1st June 2006 with the policy being in force from 15 January 2006 to 1st June 2006. This means that the norm is part of the tax regulation from 15 January 2006 to 1st June 2006. The second rule, in force from 1st April 2006, establishes a tax returns lodged after 1st April 2006. These two rules illustrate the concept of viewpoints. Consider that the conditions in the antecedent of both rules hold, then one would derive  $\neg \text{Tax}:[\overline{1\text{Jan}06, 1\text{Jun}06}]:[\widehat{15\text{Jan}06, 1\text{Jun}06}]$  but  $\text{Tax}:[\overline{1\text{Jan}06, 1\text{Jun}06}]:[\overline{1\text{Apr}06, 1\text{Jun}06}]$ , that is, if one reason from a point of view between the 15 January and the 1st April then the tax is due while if one reason from a point of view between the 1st April and he 1st June 2006 then no tax is due. Even though trivial cases of the phenomenon of retroactivity are captured by rules such as  $r1$  and  $r2$ , we should be able to detect retroactivity also in other scenarios, where normative effects are in fact applied retroactively to some conditions as a result of complex arguments that involve many rules. This problem is of great importance not only because the designer of a normative system may have the goal to state retroactive effects in more articulated scenarios, but also because she should be able to check whether such effects are not obtained when certain regulations regard matters for which retroactivity is not in general permitted. This is the case of criminal law, where the principle -Nullum crimen, nulla poena sine praevia lege poenali- is valid.

### 3 Defeasible Logic

Our system is formalised in an extension of Defeasible Logic. We provide in this section a brief recall of it. Defeasible Logic [14,15,16] is based on a logic programming-like language and it is a simple, efficient but flexible non-monotonic formalism capable of dealing with many different intuitions of non-monotonic reasoning. An argumentation semantics exists [17] that makes its use possible in argumentation systems. DL has a linear complexity [18] and also has several efficient implementations [19].

A Defeasible Logic theory is a structure  $D = (F, R, \prec)$  where  $F$  is a finite set of facts,  $R$  a finite set of rules, and  $\prec$  a superiority relation on  $R$ . Facts are indisputable statements, for example, “Bob is a minor,” formally written as  $minor(bob)$ . Rules can be strict, defeasible, or defeaters. Strict rules are rules in the classical sense; whenever the premises are indisputable, so is the conclusion. An example of a strict rule is “Minors are persons,” formally written as  $r1: minor(X) \rightarrow person(X)$ . Defeasible rules are rules that can be defeated by contrary evidence. An example of a defeasible rule is “Persons have legal capacity”; formally,  $r2: person(X) \Rightarrow hasLegalCapacity(X)$ . Defeaters are rules that cannot be used to draw any conclusion. Their only use is to prevent some conclusions by defeating some defeasible rules. An example of this kind of rule is “Minors might not have legal capacity,” formally expressed as  $r3: minor(X) \rightsquigarrow \neg hasLegalCapacity(X)$ . The idea here is that even if we know that someone is a minor, this is not sufficient evidence for the conclusion that he or she does not have legal capacity. The superiority relation between rules indicates the relative strength of each rule. That is, stronger rules override the conclusions of weaker rules. For example, if  $r3 \succ r2$ , then the rule  $r3$  overrides  $r2$ , and we can derive neither the conclusion that Bob has legal capacity nor the conclusion that he does have legal capacity.

Given a set  $R$  of rules, we denote the set of all strict rules in  $R$  by  $R_s$ , the set of defeasible rules in  $R$  by  $R_d$ , the set of strict and defeasible rules in  $R$  by  $R_{sd}$ , and the set of defeaters in  $R$  by  $R_{df}$ .  $R[q]$  denotes the set of rules in  $R$  with consequent  $q$ . In the following  $\sim p$  denotes the complement of  $p$ , that is,  $\sim p$  is  $\neg p$  if  $p$  is an atom, and  $\sim p$  is  $q$  if  $p$  is  $\neg q$ . For a rule  $r$  we will use  $A(r)$  to indicate the body or antecedent of the rule and  $C(r)$  for the head or consequent of the rule. A rule  $r$  consists of its antecedent  $A(r)$  (written on the left;  $A(r)$  may be omitted if it is the empty set), which is a finite set of literals; an arrow; and its consequent  $C(r)$ , which is a literal. In writing rules we omit set notation for antecedents. Conclusions are tagged according to whether they have been derived using defeasible rules or strict rules only. So, a conclusion of a theory  $D$  is a tagged literal having one of the following four forms:

- $+\Delta q$  meaning that  $q$  is definitely provable in  $D$ .
- $-\Delta q$  meaning that  $q$  is not definitely provable in  $D$ .
- $+\partial q$  meaning that  $q$  is defeasibly provable in  $D$ .
- $-\partial q$  meaning that  $q$  is not defeasibly provable in  $D$ .

These different notions of provability come of use here because they enable the system to label a suggestion as stronger or weaker depending on the kind of proof associated with it. Provability is based on the concept of a derivation (or proof) in  $D$ . A derivation

is a finite sequence  $P = (P(1), \dots, P(n))$  of tagged literals. Each tagged literal satisfies some proof conditions. A proof condition corresponds to the inference rules that refer to one of the four kinds of conclusions we have mentioned above.  $P(1..n)$  denotes the initial part of the sequence  $P$  of length  $n$ . We state below the conditions for defeasibly derivable conclusions:

If  $P(i+1) = +\partial q$  then

- (1)  $+\Delta q \in P(1..i)$  or
- (2) (2.1)  $\exists r \in R_{sd}[q] \forall a \in A(r) : +\partial a \in P(1..i)$  and
  - (2.2)  $-\Delta \sim q \in P(1..i)$  and
  - (2.3)  $\forall s \in R[\sim q]$  either
    - (2.3.1)  $\exists a \in A(s) : -\partial a \in P(1..i)$  or
    - (2.3.2)  $\exists t \in R_{sd}[q]$  such that
      - $\forall a \in A(t) : +\partial a \in P(1..i)$  and  $t \succ s$ .

If  $P(i+1) = -\partial q$  then

- (1)  $-\Delta q \in P(1..i)$  and
- (2) (2.1)  $\forall r \in R_{sd}[q] \exists a \in A(r) : -\partial a \in P(1..i)$  or
  - (2.2)  $+\Delta \sim q \in P(1..i)$  or
  - (2.3)  $\exists s \in R[\sim q]$  such that
    - (2.3.1)  $\forall a \in A(s) : +\partial a \in P(1..i)$  and
    - (2.3.2)  $\forall t \in R_{sd}[q]$  either
      - $\exists a \in A(t) : -\partial a \in P(1..i)$  or  $t \not\succeq s$ .

Informally, a defeasible derivation for a provable literal consists of three phases: First, we propose an argument in favour of the literal we want to prove. In the simplest case, this consists of an applicable rule for the conclusion (a rule is applicable if its antecedent has already been proved). Second, we examine all counter-arguments (rules for the opposite conclusion). Third, we rebut all the counter-arguments (the counter-argument is weaker than the pro-argument) or we undercut them (some of the premises of the counterargument are not provable).

## 4 Temporal Modal Defeasible Logic

Defeasible Logic allows us to deal with defeasibility but as such does not provide any mean to deal with modalities and temporal aspects. Temporal Modal Defeasible Logic is an umbrella expression to designate possible extensions of Defeasible Logic to capture modalities and time. We present in this section an extension of [1] with intervals as exposed in the model (see Section 2).

### 4.1 Modal Domain

The combination of mental attitudes and obligations are framed in extending Defeasible Logic following the works of [7,8,9] and capture some basic facets of the modal notions of knowledge, intentions, action and obligation.

To extend Defeasible Logic with modal operators, new types of rules relative to modal operator are introduced: arrows of the rules are labelled by the different modalities we want to deal with. This solution leads to distinguishing different modes through which the literals can be derived using rules. How such types of derivation are related to the introduction of the corresponding modalised literals can be expressed as follows: if  $X \in \{\text{KNOW}, \text{INT}, \text{ACT}, \text{OBL}\}$ , then

$$\frac{\Gamma \quad \Gamma \Rightarrow_X \psi}{\Gamma \sim X\psi} \text{ MI}$$

We make an exception when rules for knowledge are concerned. The reason for this is that we assume that beliefs are conceived of as the knowledge the agent has of the environment, and so they are used by the agent to make inferences about how the world is: in this perspective, belief conclusions correspond to factual knowledge and do not need to be modalised. But besides this exception, which can be removed if required, schema MI captures the basic logical behaviour of our modal rules. Notice, also, that actions are successful and intentional and so, when  $\text{ACT}\psi$  is derived, this also implies that  $\psi$  and  $\text{INT}\psi$  are the case.

Other relations between modalities are captured by means of *rule conversions* and *conflicts*.

The notion of *rule conversion* permits to use rules for a modality  $X$  as they were for another modality  $Y$ . Suppose that a rule of a specific type is given and also suppose that all the literals in the antecedent of a rule are provable in one and the same modality, then it is arguable that the conclusion of the rule inherits the modality of the antecedent. For example, consider the following formalisation of the Yale Shooting Problem.

$$\text{load}:\overline{[t,t]}, \text{shoot}:\overline{[t,t]} \Rightarrow_{\text{KNOW}} \text{kill}:\overline{[t,t]}$$

This rule encodes the knowledge of an agent that knows that loading the gun with live ammunitions, and then shooting will kill her friend. This example clearly shows that the qualification of the conclusions depends on the modalities relative to the individual acts “load” and “shoot”. In particular, if we obtain that the agent intends to load and to shoot the gun ( $\text{INT}(\text{load}), \text{INT}(\text{shoot})$ ), then, since she knows that the consequence of these actions is the death of her friend, she intends to kill him. However, if shooting was not intended, then we have *prima facie* to say that killing, too, was not intentional. To define the admitted conversions we introduce a binary relation *Convert* over the modalities of the language. When we write  $\text{Convert}(\text{KNOW}, \text{OBL})$  this means that a knowledge rule  $r$  can be used to derive an obligation (of course, provided that all its antecedents are derived as obligations):  $r$  can thus be converted into a rule for intention.

Beside conversions, *Conflicts* play an important role in the current context and it is crucial to establish criteria for detecting and solving conflicts between the different components which characterise the cognitive profiles of agent’s deliberation, and, above all between mental states and normative provisions. Conflicts are detected and solved by a similar strategy than basic Defeasible Logic, i.e, by following a pattern such that (i) in a first phase an argument supporting the conclusion is advanced (ii) in the second phase any possible attack are considers, and (iii) finally the counter-attack for each attack. Accordingly we introduce a ternary relation *Attack* over the set of modalities

that defines which types of rules are in conflict and which are the stronger ones. For example, if we write  $\text{Attack}(\text{OBL}, \text{INT}, \text{ACT})$  this means that, in the reasoning pattern illustrated above, obligations in general override intentions, which in turn override actions.

The relation  $\text{Attack}$  is explicitly linked to that of agent type. Classically, agent types are characterised by stating conflict resolution types in terms of orders of overruling between rules [3,7,9,8]. In this perspective, agent types are meaningful within a non-monotonic setting and are nothing but general strategies to detect and solve conflicts between the different components of the cognitive profiles of agent’s deliberation. In [3] 24 possible types are identified while, in [8], based on a different framework, 20 combinations are proposed. Typically, rational agents are assumed to be at least *realistic*: a realistic agent, in fact, is such that rules for knowledge override all other components. If the realistic condition is abandoned, we may have various forms of wishful thinking. Given the minimal assumption that a rational agent should be realistic, we may further constrain agent’s deliberation in order not to violate obligations: a *social agent* type requires that obligations are stronger than the other motivational components with the exception of beliefs. Other agent types can be specified, for which see [7,8,9].

## 4.2 Temporal domain

Approaches in temporal reasoning are traditionally based on either instants, intervals or both by representing one through the other. We represent intervals by means of instants. Formally, we consider a totally ordered discrete set  $\mathcal{T}$  of points of time termed “instants” and over it the order relation  $> \subseteq \mathcal{T} \times \mathcal{T}$ . We usually denote the variables ranging over the members of  $\mathcal{T}$  by  $t$  and its eventual subscripts, and the minimal unit by  $u$ .

Ordinarily, intervals are defined as sets of instants between two indicated instants. Here we deviate to this definition because of the non-homogeneity or transient character of events: if an event occurs in an interval conceived as a set of instants, then it would also occur in the set of instants that defines it and this would conflict with the transient characterisation of events. Hence, we define an interval as a pair of instants. Formally, an interval is a member of the set  $\text{Inter} = \{[t_1, t_2] \in \mathcal{T} \times \mathcal{T} \mid t_1 \leq t_2\}$ . As can be noted, this definition allows “punctual intervals”, i.e. intervals of the form  $[t, t]$ . Among the set  $\text{Inter}$ , we identify two subsets of interval to differentiate intervals in which an associated property holds at any instant between the boundaries and intervals in which an associated property holds at least one instant between the boundaries. We shall call the first A-interval and the second B-intervals. The set of A-intervals is denoted  $\text{AInter}$  while the set of B-intervals is denoted  $\text{BInter}$ . We shall usually denote intervals by  $T$ , A-intervals by  $\bar{T}$  and B-intervals by  $\hat{T}$  (plus eventual subscripts). We consider the functions  $\text{start}()$  and  $\text{end}()$  that returns respectively the lower bound and upper bound of an interval.

As explained in section 2, a conclusion can be associated to two temporal intervals consequently of the different temporal dimensions. The first interval is the interval of applicability with which the consequent of the rule is labelled while the second interval corresponds to the time of force interval associated to the rule. Each interval can be



assimilated to temporal Russian-dolled viewpoints from which conclusions are considered. We represent such temporalisation of conclusion by concatenation of intervals by means of the symbol ':' and we call such concatenation chain of viewpoints. Chain of viewpoints are denoted by  $V$  (plus eventual subscripts).

Temporal calculi are driven by operators over intervals. In the literature, one can find many relations that hold between intervals. For example, [20] proposes an algebra of intervals with thirteen mutually exclusive relations between two intervals. For our purpose, we consider the set of relations to "subinterval" denoted  $\sqsubseteq$ , "over" denoted  $\over$ , "meet" denoted  $\text{meet}$ , "start in" denoted  $\text{si}$ , "start before end" denoted  $\text{sbe}$ , and "start before start" denoted  $\text{sbs}$ .

**Definition 1.** *Let two intervals  $T \in \text{Inter}$  and  $T' \in \text{Inter}$ ,*

$T \sqsubseteq T'$  *iff*  $\text{start}(T') \leq \text{start}(T)$  *and*  $\text{end}(T) \leq \text{end}(T')$ .

$\over(T, T')$  *iff*  $\text{start}(T') \leq \text{start}(T) \leq \text{end}(T')$  *or*  
 $\text{start}(T') \leq \text{end}(T) \leq \text{end}(T')$  *or*  $\text{start}(T) \leq \text{start}(T') \leq \text{end}(T)$ .

$\text{meet}(T, T')$  *iff*  $\text{end}(T) + u = \text{start}(T')$ .

$\text{si}(T, T')$  *iff*  $\text{start}(T') \leq \text{start}(T) \leq \text{end}(T')$ .

$\text{sbe}(T, T')$  *iff*  $\text{start}(T) \leq \text{end}(T')$ .

$\text{sbs}(T, T')$  *iff*  $\text{start}(T) \leq \text{start}(T')$ .

Note that  $T \sqsubseteq T'$ ,  $\text{si}(T, T')$  or  $\text{sbe}(T, T')$  implies  $\over(T, T')$ , that  $T \sqsubseteq T'$  implies  $\text{si}(T, T')$  and that  $\over(T, T')$  implies  $\over(T', T)$ .

In order to lighten the paper, we may use the abbreviation consisting in placing chain of viewpoints as arguments of the previous operators, such that for example,

- $T \sqsubseteq T' : T''$  stands for  $T \sqsubseteq T'$  and  $T \sqsubseteq T''$ .
- $T : T' \sqsubseteq T'' : T'''$  stands for  $T \sqsubseteq T''$  and  $T' \sqsubseteq T'''$ .
- $T : T' \sqsubseteq T''$  stands for  $T \sqsubseteq T''$  and  $T' \sqsubseteq T''$ .

and similarly for other operators. Finally, we also use some abbreviations with regard to the function  $\text{start}()$  and  $\text{end}()$ , such that for example,

- $\text{start}(T) = \text{end}(T'' : T''')$  stands for  $\text{start}(T) = \text{end}(T'')$  and  $\text{start}(T) = \text{end}(T''')$ .
- $\text{start}(T : T') = \text{end}(T'' : T''')$  stands for  $\text{start}(T) = \text{end}(T'')$  and  $\text{start}(T') = \text{end}(T''')$ .
- $\text{start}(T : T') = \text{end}(T'')$  stands for  $\text{start}(T) = \text{end}(T'')$  and  $\text{start}(T') = \text{end}(T'')$ .

and similarly for others combinations of relation between  $\text{start}()$  and  $\text{end}()$ .

### 4.3 The Language

A temporal defeasible agent theory consists of a discrete totally ordered set of instants of time, a set of *facts* or indisputable statements, four sets of rules for knowledge, intentions, intentional actions, and obligations, and a *superiority relation*  $>$  among rules saying when a single rule may override the conclusion of another rule. For  $X \in \{\text{KNOW}, \text{INT}, \text{ACT}, \text{OBL}\}$ , a temporal *strict rule* is an expression of the form  $(\phi_1, \dots, \phi_n \rightarrow_X \psi) : \overline{T}_r$  such that whenever the premises  $\phi_1 : \widehat{T}_r, \dots, \phi_n : \widehat{T}_r$  are indisputable so is the conclusion  $\psi : \overline{T}_r$ . A *defeasible rule* is an expression of the form  $(\phi_1, \dots, \phi_n \Rightarrow_X \psi) : \overline{T}_r$  whose conclusion can be defeated by contrary evidence. An expression  $(\phi_1, \dots, \phi_n \rightsquigarrow_X \psi) : \overline{T}_r$  is a *defeater* used to defeat some defeasible rules by producing evidence to the contrary. It is worth noting that modalised literals can occur only in the antecedent of rules: the reason of this is that the rules are used to derive modalised conclusions while we do not conceptually need to iterate modalities. This limitation makes the system more manageable.

**Definition 2 (Language).** Let  $\mathcal{T}$  a discrete totally ordered set of instants of time, Prop be a set of propositional atoms, Mod = {KNOW, INT, ACT, OBL} be the set of modal operators, and Lab be a set of labels. The sets below are the smallest sets closed under the following rules:

#### Literals

$$\text{Lit} = \text{Prop} \cup \{\neg p \mid p \in \text{Prop}\}$$

#### Modal Literals

$$\text{ModLit} = \{Xl, \neg Xl \mid l \in \text{Lit}, X \in \{\text{INT}, \text{ACT}, \text{OBL}\}\};$$

#### Intervals

$$\text{Inter} = \{T = [t1, t2] \mid t1, t2 \in \mathcal{T}, t1 \leq t2\};$$

#### A-Intervals

$$\text{AInter} = \{\overline{T} = \overline{[t1, t2]} \mid t1, t2 \in \mathcal{T}, t1 \leq t2\};$$

#### B-Intervals

$$\text{BInter} = \{\widehat{T} = \widehat{[t1, t2]} \mid t1, t2 \in \mathcal{T}, t1 \leq t2\};$$

#### Chain of Viewpoints

$$\text{ChainView} = \{V = T1, V' = T1 : T2 \mid T1, T2 \in \text{AInter} \cup \text{BInter}\};$$

#### Temporal Literals

$$\text{TempLit} = \{l : T \mid l \in \text{Lit}, T \in \text{AInter} \cup \text{BInter}\};$$

#### Multi-Temporal Literals

$$\text{MTempLit} = \{l : V \mid l \in \text{Lit}, V \in \text{ChainView}\};$$

#### Temporal Modal literals

$$\text{TempModLit} = \{Xl : T \mid Xl \in \text{ModLit}, T \in \text{AInter} \cup \text{BInter}\};$$

### Multi-Temporal Modal literals

$$\text{MTempModLit} = \{Xl : V \mid Xl \in \text{ModLit}, V \in \text{ChainView}\};$$

### Temporal Rules

$$\begin{aligned} \text{Rule}_s &= \{(r : \phi_1, \dots, \phi_n \rightarrow_X \psi) : T \mid \\ &\quad r \in \text{Lab}, A(r) \subseteq \text{TempLit} \cup \text{TempModLit}, X \in \text{Mod}, \psi \in \text{TempLit}, T \in \text{AInter}\} \\ \text{Rule}_d &= \{(r : \phi_1, \dots, \phi_n \Rightarrow_X \psi) : T \mid \\ &\quad r \in \text{Lab}, A(r) \subseteq \text{TempLit} \cup \text{TempModLit}, X \in \text{Mod}, \psi \in \text{TempLit}, T \in \text{AInter}\} \\ \text{Rule}_{dfi} &= \{(r : \phi_1, \dots, \phi_n \rightsquigarrow_X \psi) : T \mid \\ &\quad r \in \text{Lab}, A(r) \subseteq \text{TempLit} \cup \text{TempModLit}, X \in \text{Mod}, \psi \in \text{TempLit}, T \in \text{AInter}\} \\ \text{Rule} &= \text{Rule}_s \cup \text{Rule}_d \cup \text{Rule}_{dfi} \end{aligned}$$

We use some abbreviations:  $A(r)$  denotes the set  $\{\phi_1, \dots, \phi_n\}$  of *antecedents* of the rule  $r$ , and  $C(r)$  to denote the *consequent*  $\psi$  of the rule  $r$ . We use also superscript for mental attitude, subscript for type of rule, and  $\text{Rule}[\phi]$  for rules whose consequent is  $\phi$ . If one does not refer to the content of the rule, a temporal rule can be written as  $r:\bar{T}$  where  $r$  is the label of the rule and  $\bar{T}$  is a temporal interval. If  $q$  is a literal,  $\sim q$  denotes the complementary literal (if  $q$  is a positive literal  $p$  then  $\sim q$  is  $\neg p$ ; and if  $q$  is  $\neg p$ , then  $\sim q$  is  $p$ );

**Definition 3 (Defeasible Agent Theory).** A defeasible agent theory is a structure

$$D = (\mathcal{T}, F, R^{\text{KNOW}}, R^{\text{INT}}, R^{\text{ACT}}, R^{\text{OBL}}, >, \mathcal{C}, \mathcal{A})$$

where

- $\mathcal{T}$  a discrete totally ordered set of instants of time;
- $F \subseteq \text{TempLit} \cup \text{TempModLit}$  is a finite set of facts;
- $R^{\text{KNOW}} \subseteq \text{Rule}^{\text{KNOW}}, R^{\text{INT}} \subseteq \text{Rule}^{\text{INT}}, R^{\text{ACT}} \subseteq \text{Rule}^{\text{ACT}}, R^{\text{OBL}} \subseteq \text{Rule}^{\text{OBL}}$  are four finite sets of rules such that each rule has a unique label;
- $> \subseteq R^{\text{KNOW} \cup \text{INT} \cup \text{ACT} \cup \text{OBL}} \times R^{\text{KNOW} \cup \text{INT} \cup \text{ACT} \cup \text{OBL}}$  is an acyclic superiority relation.
- $\mathcal{C} \subseteq \{\text{Convert}(X, Y) \mid X, Y \in \text{Mod}\}$  is a set of conversions.
- $\mathcal{A} \subseteq \{\text{Attack}(X, Y, Z) \mid X, Y, Z \in \text{Mod}\}$  is a set of attack relation.

## 4.4 Proof Theory

The formalism we have introduced allows us to temporalise rules, thus we have to admit the possibility that rules are not only given but can be proved to hold for certain span of time. Accordingly we have to give conditions that allow us to derive rules instead of literals. A conclusion of a theory  $D$  is a tagged temporal literal or rule having one of the following forms:

- + $\Delta\gamma:V$  meaning that  $\gamma:V$  is definitely provable in  $D$ .
- $\Delta\gamma:V$  meaning that  $\gamma:V$  is not definitely provable in  $D$ .

$+\partial\gamma:V$  meaning that  $\gamma:V$  is defeasible provable in  $D$ .

$-\partial\gamma:V$  meaning that  $\gamma:V$  is not defeasible provable in  $D$ .

Provability is based on the concept of a derivation (or proof) in  $D$ . A derivation is a finite sequence  $P = (P(1), \dots, P(n))$  of tagged modal literals or rules temporalised by chain of viewpoints. Each tagged temporal modal literal or rule satisfies some proof conditions, which correspond to inference rules for the four kinds of conclusions we have mentioned above. In order to lighten the presentation of the proof conditions, we present separately the condition for applicability of rules:

If  $\text{Convert}(Y, X)$  and  $r:\overline{T}_r$  is  $\Delta$ -*applicable* in the proof condition for  $\pm\Delta_X$  then

- (1)  $+\Delta r:\overline{T}_r \in P(1..i)$ , and either
  - (2)  $r:\overline{T}_r \in R^X$ ,
    - (2.1)  $\forall\alpha:\overline{T}_\alpha \in A(r:\overline{T}_r)$ ,
      - (2.1.1)  $+\Delta_{\text{KNOW}}\alpha : \overline{T}_\alpha \in P(1..i)$ , or  $+\Delta_{\text{KNOW}}\alpha : \overline{T}_\alpha:\widehat{T}_r \in P(1..i)$ , or
      - (2.1.2)  $+\Delta_{\text{ACT}}\alpha : \overline{T}_\alpha \in P(1..i)$ , or  $+\Delta_{\text{ACT}}\alpha : \overline{T}_\alpha:\widehat{T}_r \in P(1..i)$ ,
    - (2.2)  $\forall\alpha:\widehat{T}_\alpha \in A(r:\overline{T}_r)$ ,
      - (2.2.1)  $+\Delta_{\text{KNOW}}\alpha : \widehat{T}_\alpha \in P(1..i)$ , or  $+\Delta_{\text{KNOW}}\alpha : \widehat{T}_\alpha:\widehat{T}_r \in P(1..i)$ , or
      - (2.2.2)  $+\Delta_{\text{ACT}}\alpha : \widehat{T}_\alpha \in P(1..i)$ , or  $+\Delta_{\text{ACT}}\alpha : \widehat{T}_\alpha:\widehat{T}_r \in P(1..i)$ , and
  - (2.3)  $\forall Z\alpha:\overline{T}_\alpha \in A(r:\overline{T}_r)$ ,  $+\Delta_Z\alpha : \overline{T}_\alpha \in P(1..i)$ , or  $+\Delta_Z\alpha : \overline{T}_\alpha:\widehat{T}_r \in P(1..i)$ , and
  - (2.4)  $\forall Z\alpha:\widehat{T}_\alpha \in A(r:\overline{T}_r)$ ,  $+\Delta_Z\alpha : \widehat{T}_\alpha \in P(1..i)$ , or  $+\Delta_Z\alpha : \widehat{T}_\alpha:\widehat{T}_r \in P(1..i)$ , or
- (3)  $r:\overline{T}_r \in R^Y$ ,
  - (3.1)  $\forall\alpha:\overline{T}_\alpha \in A(r:\overline{T}_r)$ ,  $+\Delta_X\alpha : \overline{T}_\alpha \in P(1..i)$ , or  $+\Delta_X\alpha : \overline{T}_\alpha:\widehat{T}_r \in P(1..i)$ , and
  - (3.2)  $\forall\alpha:\widehat{T}_\alpha \in A(r:\overline{T}_r)$ ,  $+\Delta_X\alpha : \widehat{T}_\alpha \in P(1..i)$ , or  $+\Delta_X\alpha : \widehat{T}_\alpha:\widehat{T}_r \in P(1..i)$ .

The conditions for a rule  $r$  to be  $\partial$ -*applicable* are the same as those for  $\Delta$ -*applicable*, but where we replace  $\Delta$  with  $\partial$ .

If  $\text{Convert}(Y, X)$  and  $r:\overline{T}_r$  is  $\Delta$ -*discarded* in the proof condition for  $\pm\Delta_X$  then

- (1)  $-\Delta r:\overline{T}_r \in P(1..i)$ , or either
  - (2)  $r:\overline{T}_r \in R^X$ ,
    - (2.1)  $\exists\alpha:\overline{T}_\alpha \in A(r:\overline{T}_r)$ ,
      - (2.1.1)  $-\Delta_{\text{KNOW}}\alpha : \overline{T}_\alpha \in P(1..i)$ , and  $-\Delta_{\text{KNOW}}\alpha : \overline{T}_\alpha:\widehat{T}_r \in P(1..i)$ , and
      - (2.1.2)  $-\Delta_{\text{ACT}}\alpha : \overline{T}_\alpha \in P(1..i)$ , and  $-\Delta_{\text{ACT}}\alpha : \overline{T}_\alpha:\widehat{T}_r \in P(1..i)$ , or
    - (2.2)  $\exists\alpha:\widehat{T}_\alpha \in A(r:\overline{T}_r)$ ,
      - (2.2.1)  $-\Delta_{\text{KNOW}}\alpha : \widehat{T}_\alpha \in P(1..i)$ , and  $-\Delta_{\text{KNOW}}\alpha : \widehat{T}_\alpha:\widehat{T}_r \in P(1..i)$ , and
      - (2.2.2)  $-\Delta_{\text{ACT}}\alpha : \widehat{T}_\alpha \in P(1..i)$ , and  $-\Delta_{\text{ACT}}\alpha : \widehat{T}_\alpha:\widehat{T}_r \in P(1..i)$ , or
  - (2.3)  $\exists Z\alpha:\overline{T}_\alpha \in A(r:\overline{T}_r)$ ,  $-\Delta_Z\alpha : \overline{T}_\alpha \in P(1..i)$ , and  $-\Delta_Z\alpha : \overline{T}_\alpha:\widehat{T}_r \in P(1..i)$ , or
  - (2.4)  $\exists Z\alpha:\widehat{T}_\alpha \in A(r:\overline{T}_r)$ ,  $-\Delta_Z\alpha : \widehat{T}_\alpha \in P(1..i)$ , and  $-\Delta_Z\alpha : \widehat{T}_\alpha:\widehat{T}_r \in P(1..i)$ , or
- (3)  $r:\overline{T}_r \in R^Y$ ,
  - (3.1)  $\exists\alpha:\overline{T}_\alpha \in A(r:\overline{T}_r)$ ,  $-\Delta_X\alpha : \overline{T}_\alpha \in P(1..i)$ , and  $-\Delta_X\alpha : \overline{T}_\alpha:\widehat{T}_r \in P(1..i)$ , or
  - (3.2)  $\exists\alpha:\widehat{T}_\alpha \in A(r:\overline{T}_r)$ ,  $-\Delta_X\alpha : \widehat{T}_\alpha \in P(1..i)$ , and  $-\Delta_X\alpha : \widehat{T}_\alpha:\widehat{T}_r \in P(1..i)$ .

The conditions for a rule  $r:\overline{T}_r$  to be  $\partial$ -*discarded* are the same as those for  $\Delta$ -*discarded*, but where we replace  $\Delta$  with  $\partial$ .

We are now ready to define the proof theory that is, the inference conditions to derive tagged conclusions from a given theory  $D$ . We begin with the proof conditions to determine whether a rule is a definite conclusion of a theory  $D$ . A temporal rule  $r:\bar{T}$  is definitely provable ( $+\Delta$ ) if (1) there exists a rule  $r:\bar{T}_r$  in the set of rule such that  $\bar{T} \sqsubseteq \bar{T}_r$ , or (2)  $r$  is defined in two intervals  $\bar{T}_{r1}$  and  $\bar{T}_{r2}$  that make up  $\bar{T}$ . Formally:

If  $P(i+1) = +\Delta r:\bar{T}$  then

- (1)  $\exists \bar{T}_r, \bar{T} \sqsubseteq \bar{T}_r, r:\bar{T}_r \in R, \bar{T}$ , or
- (2)  $\exists \bar{T}_{r1}, \exists \bar{T}_{r2}, \text{meets}(\bar{T}_{r1}, \bar{T}_{r2}), \text{start}(\bar{T}_{r1}) = \text{start}(\bar{T}), \text{end}(\bar{T}_{r2}) = \text{end}(\bar{T}), r:\bar{T}_{r1} \in R$  and  $r:\bar{T}_{r2} \in R$ .

A rule  $r$  is not definitely provable at interval  $\bar{T}$  if (1) there is not such rule in the set rules defined in a larger interval (2)  $r$  is not defined in any intervals  $\bar{T}_{r1}$  and  $\bar{T}_{r2}$  that make up  $\bar{T}$ . Formally:

If  $P(i+1) = -\Delta r:\bar{T}$  then

- (1)  $\forall \bar{T}_r, \bar{T} \sqsubseteq \bar{T}_r, r:\bar{T}_r \notin R, \bar{T}$ , and
- (2)  $\forall \bar{T}_{r1}, \forall \bar{T}_{r2}, \text{meets}(\bar{T}_{r1}, \bar{T}_{r2}), \text{start}(\bar{T}_{r1}) = \text{start}(\bar{T}), \text{end}(\bar{T}_{r2}) = \text{end}(\bar{T}), r:\bar{T}_{r1} \notin R$  or  $r:\bar{T}_{r2} \notin R$ .

A temporal rule  $r:\hat{T}$  is definitely provable ( $+\Delta$ ) if there exists a rule  $r:\bar{T}_r$  in the set of rule such that  $\text{over}(\hat{T}, \bar{T}_r)$ . Formally:

If  $P(i+1) = +\Delta r:\hat{T}$  then  $\exists \bar{T}_r, \text{over}(\hat{T}, \bar{T}_r), r:\bar{T}_r \in R$ .

If  $P(i+1) = -\Delta r:\hat{T}$  then  $\forall \bar{T}_r, \text{over}(\hat{T}, \bar{T}_r), r:\bar{T}_r \notin R$ .

We can now move to definite conclusion of temporal literals. We begin with literals temporalised by a chain of viewpoints  $\bar{V}$ , i.e. temporalised by  $\bar{T}$  or  $\bar{T}:\bar{T}^i$ .

If  $P(i+1) = +\Delta_X \gamma:\bar{V}$  then

- (1)  $\exists \bar{T}_\gamma, \bar{V} \sqsubseteq \bar{T}_\gamma, X\gamma:\bar{T}_\gamma \in F$ , or
- (2) if  $X = \text{KNOW}$  then  $\exists \bar{T}_\gamma, \bar{V} \sqsubseteq \bar{T}_\gamma, \gamma:\bar{T}_\gamma \in F$ , or
- (3) if  $X = \text{INT}$  then  $\exists \bar{V}_\gamma, \bar{V} \sqsubseteq \bar{V}_\gamma, +\Delta_{\text{ACT}}\gamma:\bar{V}_\gamma$ , or
- (4)  $\exists r:\bar{T}_r \in R_s[\gamma:\bar{T}_\gamma], \bar{V} \sqsubseteq \bar{T}_\gamma:\bar{T}_r, r:\bar{T}_r$  is  $\Delta$ -applicable, or
- (5)  $\exists \bar{V}_{\gamma1}, \exists \bar{V}_{\gamma2}, \text{meets}(\bar{V}_{\gamma1}, \bar{V}_{\gamma2}), \text{start}(\bar{V}_{\gamma1}) = \text{start}(\bar{V}), \text{end}(\bar{V}_{\gamma2}) = \text{end}(\bar{V})$   
 $+\Delta_X \gamma:\bar{V}_{\gamma1} \in P(1..i)$  and  $+\Delta_X \gamma:\bar{V}_{\gamma2} \in P(1..i)$ .

To prove that a modal literal temporalised by a chain of viewpoints is not definitely provable we have to show that any attempt to give a definite proof fails.

If  $P(i+1) = -\Delta_X \gamma:\bar{V}$  then

- (1)  $\forall \bar{T}_\gamma, \bar{V} \sqsubseteq \bar{T}_\gamma, X\gamma:\bar{T}_\gamma \notin F$ , and
- (2) if  $X = \text{KNOW}$  then  $\forall \bar{T}_\gamma, \bar{V} \sqsubseteq \bar{T}_\gamma, \gamma:\bar{T}_\gamma \notin F$ , and
- (3) if  $X = \text{INT}$  then  $\forall \bar{V}_\gamma, \bar{V} \sqsubseteq \bar{V}_\gamma, -\Delta_{\text{ACT}}\gamma:\bar{V}_\gamma$ , and
- (4)  $\forall r:\bar{T}_r \in R_s[\gamma:\bar{T}_\gamma], \bar{V} \sqsubseteq \bar{T}_\gamma:\bar{T}_r, r:\bar{T}_r$  is  $\Delta$ -discarded, and
- (5)  $\forall \bar{V}_{\gamma1}, \forall \bar{V}_{\gamma2}, \text{meets}(\bar{V}_{\gamma1}, \bar{V}_{\gamma2}), \text{start}(\bar{V}_{\gamma1}) = \text{start}(\bar{V}), \text{end}(\bar{V}_{\gamma2}) = \text{end}(\bar{V})$   
 $-\Delta_X \gamma:\bar{V}_{\gamma1} \in P(1..i)$  or  $-\Delta_X \gamma:\bar{V}_{\gamma2} \in P(1..i)$ .

The conditions for a temporal literal  $\gamma:\widehat{V}$  (i.e.  $\gamma:\widehat{T}$  or  $\gamma:\widehat{T}:\widehat{T}'$ ) to be not definitely provable with modality X ( $\pm\Delta_X$ ) are formally expressed below.

If  $P(i+1) = +\Delta_X\gamma:\widehat{V}$  then  $\exists\widehat{V}_\gamma, \text{over}(\widehat{V}, \widehat{V}_\gamma), +\Delta_X\gamma:\widehat{V}_\gamma$ .

If  $P(i+1) = -\Delta_X\gamma:\widehat{V}$  then  $\forall\widehat{V}_\gamma, \text{over}(\widehat{V}, \widehat{V}_\gamma), -\Delta_X\gamma:\widehat{V}_\gamma$ .

The definition of  $\Delta$ -*applicable*, and  $\Delta$ -*discarded* of rules contains the definite (un)provability of modal literals temporalised by chain of viewpoint of the type  $\overline{T}:\widehat{T}_r$ . We cater for such cases in the two next proof conditions.

If  $P(i+1) = +\Delta_X\gamma:\overline{T}:\widehat{T}_r$  then  
 $\exists\overline{T}_{\gamma 1}, \exists\overline{T}_{r 1}, \overline{T} \sqsubseteq \overline{T}_{\gamma 1}, \text{over}(\widehat{T}_r, \overline{T}_{r 1}), +\Delta_X\gamma:\overline{T}_{\gamma 1}:\overline{T}_{r 1} \in P(1..i)$ .

If  $P(i+1) = -\Delta_X\gamma:\overline{T}:\widehat{T}_r$  then  
 $\forall\overline{T}_{\gamma 1}, \forall\overline{T}_{r 1}, \overline{T} \sqsubseteq \overline{T}_{\gamma 1}, \text{over}(\widehat{T}_r, \overline{T}_{r 1}), -\Delta_X\gamma:\overline{T}_{\gamma 1}:\overline{T}_{r 1} \in P(1..i)$ .

We now turn our attention to defeasible derivations, that is, derivations giving a temporal assertion  $\gamma:V$  as a defeasible conclusion of a theory  $D$ . We begin with the proof conditions to determine whether a rule is a defeasible conclusion.

If  $P(i+1) = +\partial r:\overline{T}$  then  $+\Delta r:\overline{T} \in P(1..i)$

If  $P(i+1) = +\partial r:\widehat{T}$  then  $+\Delta r:\widehat{T} \in P(1..i)$ .

Defeasible provability ( $+\partial$ ) for temporal literals consists of three phases. In the first phase, we put forward a supported reason for the temporal assertion that we want to prove. Then in the second phase, we consider all possible attacks against the desired conclusion. Finally in the last phase, we have to counter-attack the attacks considered in the second phase.

If  $P(i+1) = +\partial_X\gamma:\overline{V}$  and  $\text{Convert}(Y, X)$  and  $\text{Attack}(W, Z, X)$  then

(1)  $+\Delta_X\gamma:\overline{V} \in P(1..i)$ , or

(2)  $-\Delta_X\sim\gamma:\widehat{V} \in P(1..i)$ , and

(2.1) if  $X = \text{INT}$  then  $\exists\overline{V}_\gamma, \overline{V} \sqsubseteq \overline{V}_\gamma, +\partial_{\text{ACT}}\gamma:\overline{V}_\gamma$ , or

(2.2)  $\exists r:\overline{T}_r \in R^{X \cup Y}[\gamma:\overline{T}_\gamma], \overline{V} \sqsubseteq \overline{T}_\gamma:\overline{T}_r, r:\overline{T}_r$  is  $\partial$ -applicable,

(2.3)  $\forall s:\overline{T}_s \in R^{W \cup Z \cup X \cup Y}[\sim\gamma:\overline{T}_\sim\gamma], \text{si}(\overline{T}_\sim\gamma:\overline{T}_s, \overline{T}_\gamma:\overline{T}_r), \text{sbe}(\overline{T}_\sim\gamma:\overline{T}_s, \overline{V})$ ,

(2.3.1)  $s:\overline{T}_s$  is  $\partial$ -discarded, or

(2.3.2)  $\exists w:\overline{T}_w \in R^K[\gamma:\overline{T}_w\gamma], \overline{V} \sqsubseteq \overline{T}_w\gamma:\overline{T}_w$ ,

(2.3.2.1)  $w:\overline{T}_w$  is  $+\partial$ -applicable, and either

(2.3.2.2)  $s:\overline{T}_s \in R^{X \cup Y}$ ,

(2.3.2.2.1)  $w:\overline{T}_w \in R^{W \cup Z}$ , or

(2.3.2.2.2)  $w:\overline{T}_w \in R^{X \cup Y}, w:\overline{T}_w \succ s:\overline{T}_s$ , or

(2.3.2.3)  $s:\overline{T}_s \in R^Z$ ,

(2.3.2.3.1)  $w:\overline{T}_w \in R^W$ , or

(2.3.2.3.2)  $w:\overline{T}_w \in R^Z, w:\overline{T}_w \succ s:\overline{T}_s$ , or

(2.3.2.4)  $s:\overline{T}_s \in R^W, w:\overline{T}_w \in R^W, w:\overline{T}_w \succ s:\overline{T}_s$ , or

(3)  $\exists\overline{V}_{\gamma 1}, \exists\overline{V}_{\gamma 2}, \text{meets}(\overline{V}_{\gamma 1}, \overline{V}_{\gamma 2}), \text{start}(\overline{V}_{\gamma 1}) = \text{start}(\overline{V}), \text{end}(\overline{V}_{\gamma 2}) = \text{end}(\widehat{V})$   
 $+\partial_X\gamma:\overline{V}_{\gamma 1} \in P(1..i)$  and  $+\partial_X\gamma:\overline{V}_{\gamma 2} \in P(1..i)$ .

If  $P(i+1) = -\partial_X \gamma: \bar{V}$  and  $\text{Convert}(Y, X)$  and  $\text{Attack}(W, Z, X)$  then

- (1)  $-\Delta_X \gamma: \bar{V} \in P(1..i)$ , and
- (2)  $+\Delta_X \sim \gamma: \hat{V} \in P(1..i)$ , or
  - (2.1) if  $X = \text{INT}$  then  $\forall \bar{V}_\gamma, \bar{V} \sqsubseteq \bar{V}_\gamma, -\partial_{\text{ACT}} \gamma: \bar{V}_\gamma$ , and
  - (2.2)  $\forall r: \bar{T}_r \in R_{sd}^{XUY}[\gamma: \bar{T}_\gamma], \bar{V} \sqsubseteq \bar{T}_\gamma: \bar{T}_r, r: \bar{T}_r$  is  $\partial$ -applicable,
  - (2.3)  $\exists s: \bar{T}_s \in R^{WUZUXUY}[\sim \gamma: \bar{T}_{\sim \gamma}], \text{si}(\bar{T}_{\sim \gamma}: \bar{T}_s, \bar{T}_\gamma: \bar{T}_r), \text{sbe}(\bar{T}_{\sim \gamma}: \bar{T}_s, \bar{V})$ ,
    - (2.3.1)  $s: \bar{T}_s$  is  $\partial$ -applicable, and
    - (2.3.2)  $\forall w: \bar{T}_w \in R[\gamma: \bar{T}_w \gamma], \bar{T} \sqsubseteq \bar{T}_w \gamma: \bar{T}_w$ , either
      - (2.3.2.1)  $w: \bar{T}_w$  is  $\partial$ -discarded, or
      - (2.3.2.2)  $s: \bar{T}_s \in R^{XUY}$ ,
        - (2.3.2.2.1)  $w: \bar{T}_w \notin R^{WUZ}$ , and
        - (2.3.2.2.2)  $w: \bar{T}_w \in R^{XUY}, w: \bar{T}_w \not\prec s: \bar{T}_s$ , and
      - (2.3.2.3)  $s: \bar{T}_s \in R^Z$ 
        - (2.3.2.3.1)  $w: \bar{T}_w \notin R^W$ , and
        - (2.3.2.3.2)  $w: \bar{T}_w \in R^Z, w: \bar{T}_w \not\prec s: \bar{T}_s$ , and
      - (2.3.2.4)  $s: \bar{T}_s \in R^W, w: \bar{T}_w \in R^W, w: \bar{T}_w \not\prec s: \bar{T}_s$ , and
- (3)  $\forall \bar{V}_{\gamma 1}, \forall \bar{V}_{\gamma 2}, \text{meets}(\bar{V}_{\gamma 1}, \bar{V}_{\gamma 2}), \text{start}(\bar{V}_{\gamma 1}) = \text{start}(\bar{V}), \text{end}(\bar{V}_{\gamma 2}) = \text{end}(\bar{V})$   
 $-\partial_X \gamma: \bar{V}_{\gamma 1} \in P(1..i)$  and  $-\partial_X \gamma: \bar{V}_{\gamma 2} \in P(1..i)$

Let us illustrate the proof condition of the defeasible provability of  $X\gamma: \bar{V}$ . We have two cases: 1) We show that  $X\gamma: \bar{V}$  is already definitely provable; or 2) we need to argue using the defeasible part of  $D$ . In this second case, to prove  $X\gamma: \bar{V}$  defeasibly we must show that  $X\sim \gamma: \hat{V}$  is not definitely provable (2). We require then there must be a strict or defeasible rule  $r: \bar{T}_r \in R^{XUY}$  which can be applied and with head  $\gamma: \bar{T}_\gamma$  such that  $\bar{V} \sqsubseteq \bar{T}_\gamma: \bar{T}_r$  (2.1). But now we need to consider possible attacks, i.e., reasoning chains in support of  $\sim \gamma: \bar{V}$ , that is, any rule  $s: \bar{T}_s \in R^{WUZUXUY}$  which has head  $\sim \gamma: \bar{T}_{\sim \gamma}$  such that  $\text{si}(\bar{T}_{\sim \gamma}: \bar{T}_s, \bar{T}_\gamma: \bar{T}_r)$ , and  $\text{sbe}(\bar{T}_{\sim \gamma}: \bar{T}_s, \bar{V})$ . Note that here we consider defeaters, too, whereas they could not be used to support the conclusion  $X\gamma: \bar{V}$ ; this is in line with the motivation of defeaters given earlier. These attacking rules  $s: \bar{T}_s$  have to be discarded (2.3.1), or must be counterattacked by a stronger rule  $w: \bar{T}_w$  which has a head  $\gamma: \bar{T}_w \gamma$  such that  $\bar{V}$  is contained in  $\bar{T}_w \gamma: \bar{T}_w$  (2.3.2). Finally, we have to cater for the case where  $X\gamma$  is defeasible provable on  $\bar{V}_{\gamma 1}$  and  $\bar{V}_{\gamma 2}$  that make up  $\bar{V}$  (3).

The defeasible proof for a temporal literal to hold in some instants of a chain of viewpoints  $\hat{V}$  is less demanding since it is sufficient to prove it for at least an instant in  $\bar{V}$ .

If  $P(i+1) = +\partial_X \gamma: \hat{V}$  and  $\text{Convert}(Y, X)$  and  $\text{Attack}(W, Z, X)$  then  
 $\exists \bar{V}_\gamma, \text{over}(\hat{V}, \bar{V}_\gamma), +\partial_X \gamma: \bar{V}_\gamma \in P(1..i)$ .

If  $P(i+1) = -\partial_X \gamma: \hat{V}$  and  $\text{Convert}(Y, X)$  and  $\text{Attack}(W, Z, X)$  then  
 $\forall \bar{V}_\gamma, \text{over}(\hat{V}, \bar{V}_\gamma), -\partial_X \gamma: \bar{V}_\gamma \in P(1..i)$ .

Similarly to definite conclusions, the definition of  $\partial$ -applicable, and  $\partial$ -discarded of rules contains the defeasible (un)provability of modal literals temporalised by a chain of viewpoints of the type  $\bar{T}: \hat{T}$ . We cater for such cases by these two final proof conditions.

If  $P(i+1) = +\partial_X \gamma: \overline{T}:\widehat{T}_r$  then  
 $\exists \overline{T}_{\gamma_1}, \exists \overline{T}_{r_1}, \overline{T} \sqsubseteq \overline{T}_{\gamma_1}$ , over  $(\widehat{T}_r, \overline{T}_{r_1})$ ,  $+\partial_X \gamma: \overline{T}_{\gamma_1}:\overline{T}_{r_1} \in \mathbf{P}(1..i)$ .

If  $P(i+1) = -\partial_X \gamma: \overline{T}:\widehat{T}_r$  then  
 $\forall \overline{T}_{\gamma_1}, \forall \overline{T}_{r_1}, \overline{T} \sqsubseteq \overline{T}_{\gamma_1}$ , over  $(\widehat{T}_r, \overline{T}_{r_1})$ ,  $-\partial_X \gamma: \overline{T}_{\gamma_1}:\overline{T}_{r_1} \in \mathbf{P}(1..i)$ .

Proof conditions for modal literals temporalised by chain of viewpoints of the type  $\widehat{T}:\overline{T}_r$  are not presented here but follows similar schema.

## 5 Conclusions

In this paper we extended the logic presented in [1] with temporal intervals in order to express its expressive power. Doing so, we have extended the programming cognitive agents approach with modal literals and rules temporalised with intervals. This makes the resulting logic more expressive and more suitable for the task at hand. In addition we have considered the notion of viewpoint. The deliberation of an agent based on a policy depends not only on the environment but also on the rules in force in the policy at the time of deliberation and at the time when the plan resulting from the deliberation will be executed. These two aspects are neglected in the literature on agent planning. An aspect we did not consider here is how revise theories in the same way as complex modification of normative codes [11].

## References

1. Governatori, G., Padmanabhan, V., Rotolo, A.: Rule-based agents in temporalised defeasible logic. In Yang, Q., Webb, G., eds.: Ninth Pacific Rim International Conference on Artificial Intelligence. Number 4099 in LNAI, Berlin, Springer (2006) 31–40
2. Dastani, M., van der Torre, L.W.N.: Programming boid-plan agents: Deliberating about conflicts among defeasible mental attitudes and plans. In: AAMAS. (2004) 706–713
3. Broersen, J., Dastani, M., Hulstijn, J., van der Torre, L.: Goal generation in the BOID architecture. *Cognitive Science Quarterly* **2** (2002) 428–447
4. Dastani, M., de Boer, F., Dignum, F., Meyer, J.J.: Programming agent deliberation: an approach illustrated using the 3apl language. In: AAMAS '03: Proceedings of the second international joint conference on Autonomous agents and multiagent systems, New York, NY, USA, ACM Press (2003) 97–104
5. Dastani, M., van Riemsdijk, B., Dignum, F., Meyer, J.: A programming language for cognitive agents: Goal-directed 3apl. In: PROMAS. (2003) 111–130
6. van Riemsdijk, M.B., Dastani, M., Meyer, J.J.C.: Semantics of declarative goals in agent programming. In: AAMAS '05: Proceedings of the fourth international joint conference on Autonomous agents and multiagent systems, New York, NY, USA, ACM Press (2005) 133–140
7. Governatori, G., Rotolo, A.: Defeasible logic: Agency, intention and obligation. In Lomuscio, A., Nute, D., eds.: *Deontic Logic in Computer Science*. Number 3065 in LNAI, Berlin, Springer-Verlag (2004) 114–128
8. Dastani, M., Governatori, G., Rotolo, A., van der Torre, L.: Programming cognitive agents in defeasible logic. In Sutcliffe, G., Voronkov, A., eds.: *Proc. LPAR 2005*. Volume 3835 of LNAI., Springer (2005) 621–636



9. Dastani, M., Governatori, G., Rotolo, A., van der Torre, L.: Preferences of agents in defeasible logic. In Zhang, S., Jarvis, R., eds.: Proc. Australian AI05. Volume 3809 of LNAI, Springer (2005) 695–704
10. Governatori, G., Rotolo, A., Sartor, G.: Temporalised normative positions in defeasible logic. In Gardner, A., ed.: 10th International Conference on Artificial Intelligence and Law (ICAIL05), ACM Press (2005) 25–34
11. Governatori, G., Palmirani, M., Riveret, R., Rotolo, A., Sartor, G.: Norm modifications in defeasible logic. In Moens, M.F., Spyns, P., eds.: Legal Knowledge and Information Systems. Number 134 in Frontiers in Artificial Intelligence and Applications. IOS Press, Amsterdam (2005) 13–22
12. Bratman, M.E.: Intentions, Plans and Practical Reason. Harvard University Press, Cambridge, MA (1987)
13. Palmirani, M., Brighi, R.: Time model for managing the dynamic of normative system. In: Proceedings of EGOV 2006, Berlin, Springer (2006) 207–218
14. Nute, D.: Defeasible reasoning. In: Proceedings of 20th Hawaii International Conference on System Science, IEEE press (1987)
15. Nute, D.: Defeasible logic. In Gabbay, D., Hogger, C., Robinson, J., eds.: Handbook of Logic in Artificial Intelligence and Logic Programming. Volume 3. Oxford University Press (1993) 353–395
16. Antoniou, G., Billington, D., Governatori, G., Maher, M.J.: Representation results for defeasible logic. *ACM Transactions on Computational Logic* **2** (2001) 255–287
17. Governatori, G., Maher, M.J., Billington, D., Antoniou, G.: Argumentation semantics for defeasible logics. *Journal of Logic and Computation* **14** (2004) 675–702
18. Maher, M.J.: Propositional defeasible logic has linear complexity. *Theory and Practice of Logic Programming* **1** (2001) 691–711
19. Bassiliades, N., Antoniou, G., Vlahavas, I.: DR-DEVICE: A defeasible logic system for the Semantic Web. In Ohlbach, H.J., Schaffert, S., eds.: 2nd Workshop on Principles and Practice of Semantic Web Reasoning. Number 3208 in LNCS, Springer (2004) 134–148
20. Allen, J.F.: Towards a general theory of action and time. *Artif. Intell.* **23** (1984) 123–154