

# Formalising Motivational Attitudes of Agents Using the KARO Framework

B. van Linder J.-J. Ch. Meyer W. van der Hoek  
Utrecht University  
Department of Computer Science  
P.O. Box 80.089  
3508 TB Utrecht  
The Netherlands  
Email: jj@cs.ruu.nl

## Abstract

In this paper we present a formalisation of motivational attitudes, the attitudes that are the driving forces behind the actions of agents. We consider the statics of these attitudes both at the assertion level, i.e. ranging over propositions, and at the practition<sup>1</sup> level, i.e. ranging over actions, as well as the dynamics of these attitudes, i.e. how they change over time. Starting from an agent's wishes, which form the primitive, most fundamental motivational attitude, we define its goals as induced by those wishes that do not yet hold, i.e. are unfulfilled, but are within the agent's practical possibility to bring about, i.e. are implementable for the agent. Among these unfulfilled, implementable wishes the agent selects those that qualify as its goals. Based on its knowledge on its goals and practical possibilities, an agent may make certain commitments. In particular, an agent may commit itself to actions that it knows to be correct and feasible to bring about some of its known goals. As soon as it no longer knows its commitments to be useful, i.e. leading to fulfilment of some goal, and practically possible, an agent is able to undo these commitments. Both the act of committing as well as that of undoing commitments is modelled as a special model-transforming action in our framework, which extends the usual state-transition paradigm of Propositional Dynamic Logic. In between making and undoing commitments, an agent is committed to all the actions that are known to be identical for all practical purposes to the ones in its agenda. By modifying the agent's agenda during the execution of actions in a straightforward way, it is ensured that commitments display an acceptable behaviour with regard to composite actions.

## 1. Introduction

The formalisation of rational agents is a topic of continuing interest in Artificial Intelligence. Research on this subject has held the limelight ever since the pioneering work of Moore [26] in which knowledge and actions are considered. Over the years important contributions have been made on both *informational* attitudes like knowledge and belief [11], and *motivational* attitudes like commitments and obligations [5]. Recent developments include the Belief-Desire-Intention architecture [30], logics for the specification and verification of multi-agent systems [40], and cognitive robotics [18].

---

<sup>1</sup>The term 'practition' is due to Castañeda [2].

In a series of papers [14, 15, 21, 16, 23, 22, 24] we defined a *theorist* logic for rational agents, i.e. a formal system that may be used to *specify, analyse* and *reason about* the behaviour of rational agents. In the basic framework [14, 15], the *knowledge, belief* and *abilities* of agents, as well as the *opportunities* for and the *results* of their actions are formalised. In this so-called KARO framework it can for instance be modelled that an agent knows that some action is *correct* to bring about some state of affairs since it knows that performing the action will lead to the state of affairs, and that it knows that an action is *feasible* in the sense that the agent knows of its ability to perform the action.

Having dealt with both informational attitudes and various aspects of action in previous work, this paper is aimed at providing a formalisation of the *motivational* attitudes of rational agents. In the last decade various formalisations of different kinds of motivational attitudes have been proposed [5, 29, 37]. The approach presented in this paper makes three main contributions to the theory of formalising motivational attitudes. Firstly, we consider a fairly wide scope of motivational attitudes, situated at two different levels. At the *assertion* level, this is the level where operators deal with assertions, we consider *wishes* and *goals*. At the *practition* level, where operators range over actions, we define *commitments*. With respect to these commitments we introduce both an operator modelling the commitments that an agent has made, and an action which models the act of committing. The notions that we formalise avoid (most of) the well-known problems that plague formalisations of motivational attitudes. Secondly, our formalisation of the various notions is strictly *bottom up*. That is, after defining the fundamental notion of wishes, goals are defined in terms of wishes, and commitments are introduced using the notion of goals. In this way, we provide a formalisation of motivational attitudes that does not have to resort to ‘tricks’ like (circularly) defining the intention to do an action in terms of the goal to have done it. Lastly, in our formalisation we aim at being faithful — to a certain degree — to the insights on motivational attitudes as they have been gained in the philosophical research on practical reasoning.

The rest of the paper is organised as follows. To sketch the context and the area of application of this research, we start in Section 2 with the (re)introduction of some of our ideas on knowledge, belief, abilities, opportunities, and results; furthermore the definition of our formal language is given. In Section 3 we discuss the motivational attitudes that we will treat in this paper, and extend the basic language with extra operators to express these attitudes. Next we turn to the semantics of our language: in Section 4 we present the models of our basic framework, in Section 5 we consider the semantics of the operator for expressing wishes and in Section 6 and Section 7 the notions of a goal and a commitment, respectively, are treated formally. We will be able to formalise goals in such a way that problems of ‘logical omniscience’ that plague most formalisations of motivational attitudes are avoided. Finally, in Section 8 we summarise and conclude. Selected proofs are provided in an appendix.

## 2. Knowledge, abilities, opportunities, and results

The main informational attitude that we consider is that of *knowledge*. In representing knowledge we follow the approach common in epistemic logic [11, 13]: the formula  $\mathbf{K}_i\varphi$  denotes the fact that agent  $i$  knows  $\varphi$ , and is interpreted in a Kripke-style possible worlds semantics.

At the action level we consider *results, abilities* and *opportunities*. Slightly simplifying ideas of Von Wright [41], we consider any aspect of the state of affairs brought about by the execution of an action by an agent in some state to be among the results of the event consisting of the execution of that particular action by that particular agent, in that particular state. An

important aspect of any investigation of action is the relation that exists between ability and opportunity. In order to successfully complete an action, both the opportunity and the ability to perform the action are necessary. Although these notions are interconnected, they are surely not identical: the abilities of agents comprise mental and physical powers, moral capacities, and physical possibility, whereas the opportunity to perform actions is best described by the notion of circumstantial possibility (cf. [17]). The abilities of agents are formalised via the  $\mathbf{A}_i$  operator; the formula  $\mathbf{A}_i\alpha$  denotes the fact that agent  $i$  has the ability to do  $\alpha$ . When using the descriptions of opportunities and results as given above, the framework of (propositional) dynamic logic provides an excellent means to formalise these notions. Using events  $\text{do}_i(\alpha)$  to refer to the performance of the action  $\alpha$  by the agent  $i$ , we consider the formulae  $\langle \text{do}_i(\alpha) \rangle \varphi$  and  $[\text{do}_i(\alpha)] \varphi$ . As we shall only encounter deterministic actions in this paper,  $\langle \text{do}_i(\alpha) \rangle \varphi$  is the stronger of these formulae; it represents the fact that agent  $i$  has the opportunity to do  $\alpha$  and that doing  $\alpha$  leads to  $\varphi$ . The formula  $[\text{do}_i(\alpha)] \varphi$  is noncommittal about the opportunity of the agent to do  $\alpha$  but states that if the opportunity to do  $\alpha$  is indeed present, doing  $\alpha$  results in  $\varphi$ .

2.1. DEFINITION. Let denumerable sets  $A = \{1, \dots, n\}$  of agents,  $\Pi$  of propositional symbols and  $\text{At}$  of atomic actions be given. The language  $L$  is the smallest superset of  $\Pi$  such that:

- if  $\varphi, \psi \in L, i \in A, \alpha \in \text{Ac}$  then  $\neg\varphi, \varphi \vee \psi, \mathbf{K}_i\varphi, \langle \text{do}_i(\alpha) \rangle \varphi, \mathbf{A}_i\alpha \in L$

where  $\text{Ac}$  is the smallest superset of  $\text{At}$  such that if  $\varphi \in L, \alpha, \alpha_1, \alpha_2 \in \text{Ac}$  then

- **confirm**  $\varphi \in \text{Ac}$  *confirmations*
- $\alpha_1; \alpha_2 \in \text{Ac}$  *sequential composition*
- **if**  $\varphi$  **then**  $\alpha_1$  **else**  $\alpha_2$  **fi**  $\in \text{Ac}$  *conditional composition*
- **while**  $\varphi$  **do**  $\alpha$  **od**  $\in \text{Ac}$  *repetitive composition*

The intuitive interpretation of the **confirm**  $\varphi$  action is the verification whether  $\varphi$  holds: if it does, execution can be continued with the following action; if it doesn't, execution fails (remains pending without yielding a next state). Actions that are either atomic or a confirmation, are called *semi-atomic*. The class of semi-atomic actions is denoted by  $\text{At}^+$ . The sequential composition  $\alpha_1; \alpha_2$  is interpreted as  $\alpha_1$  followed by  $\alpha_2$ . The conditional composition **if**  $\varphi$  **then**  $\alpha_1$  **else**  $\alpha_2$  **fi** means the execution of  $\alpha_1$  if  $\varphi$  holds and that of  $\alpha_2$  otherwise. The repetitive composition **while**  $\varphi$  **do**  $\alpha$  **od** is interpreted as executing  $\alpha$  as long as  $\varphi$  holds. The formal semantics of these actions will be treated in the next section. The purely propositional fragment of  $L$  is denoted by  $L_0$ . Constructs  $\wedge, \rightarrow, \leftrightarrow, \top$  and  $\perp$  are defined in the usual way, and  $[\text{do}_i(\alpha)] \varphi$  is used as an abbreviation of  $\neg \langle \text{do}_i(\alpha) \rangle \neg \varphi$ . Finally, we define the class  $\text{Acseq}$  as the smallest superset of  $\text{At}^+$  closed under sequential composition. Thus  $\text{Acseq}$  consists of sequences of basic actions, where we consider both atomic actions and confirm actions as basic.

## 2.1. The Can-predicate and the Cannot-predicate

To formalise the knowledge of agents on their practical (im)possibilities, we introduced the so-called Can-predicate and Cannot-predicate. These are binary predicates, pertaining to a pair consisting of an action and a proposition, and denoting that an agent knows that performing the action constitutes a practical (im)possibility to bring about the proposition. We consider practical possibility to consist of two parts, viz. correctness and feasibility: action  $\alpha$  is *correct* with respect to  $\varphi$  iff  $\langle \text{do}_i(\alpha) \rangle \varphi$  holds and  $\alpha$  is *feasible* iff  $\mathbf{A}_i\alpha$  holds.

2.2. DEFINITION. The Can-predicate and the Cannot-predicate are, for all agents  $i$ , actions  $\alpha$  and formulae  $\varphi$ , defined as follows.

- $\mathbf{PracPoss}_i(\alpha, \varphi) =^{\text{def}} (\text{do}_i(\alpha))\varphi \wedge \mathbf{A}_i\alpha$
- $\mathbf{Can}_i(\alpha, \varphi) =^{\text{def}} \mathbf{K}_i\mathbf{PracPoss}_i(\alpha, \varphi)$
- $\mathbf{Cannot}_i(\alpha, \varphi) =^{\text{def}} \mathbf{K}_i\neg\mathbf{PracPoss}_i(\alpha, \varphi)$

Thus the Can-predicate and the Cannot-predicate express the agent’s knowledge about its practical possibilities and impossibilities, respectively. Therefore these predicates are important for the agent’s planning of actions.

### 3. Motivational attitudes: wishes, goals and commitments

Motivational attitudes constitute what probably is the most fundamental, primitive and essential characteristic of agency. These attitudes provide the motive for any act on behalf of the agents, i.e. the acting of agents is driven by their motivational attitudes. Typical examples of motivational attitudes are amongst others wishes, desires, concerns, ambitions, goals, intentions and commitments. The meaning of most of these terms is intuitively much less clear than that of the informational attitudes of knowledge and belief, or of the aspects of action (result, opportunity, ability) that we considered in the previous section. It is therefore also not clear which of the aforementioned motivational attitudes are relevant, and worth formalising, when modelling rational agents. In their BDI-architecture, Rao & Georgeff [30] consider desires and intentions to be primitive, and define a notion of commitment in terms of these, Cohen & Levesque [5] consider goals to be primitive and define intentions using goals, and Shoham [35] restricts himself to formalising commitments. In our opinion each of these formalisations lacks some of the aspects that are vital to modelling motivational attitudes. Firstly, psychological evidence seems to suggest that notions like goals, intentions and commitments are not primitive, but rather induced by some more fundamental notion. There is an ongoing debate in the psychological literature on the actual nature of this notion. Aristotle, and in his footsteps the adherents of cognitive motivation theories, proposed that the pursuit of knowledge is men’s most primitive motivational attitude. Freud distinguished in his psychoanalytical theory two fundamental human motivations, viz. libido, the complex of desires associated with sexuality and the erotic, and aggression, which comprises besides a proclivity towards destruction also the propensity to self-preservation. Although the theories of Aristotle and Freud are readily applicable to human agents, the question arises whether this is also the case for non-human, artificial agents. We argue that this is indeed so. Highly advanced artificial agents that are meant to behave in a human-like intelligent way, must almost necessarily be ascribed human motives. On the other hand, software agents that assist some user with some specific task, are in general not to be ascribed human motives. The primitive motivational attitudes of these agents could, however, quite well be identified with those of the user they are assisting. In this paper we will simply assume that all agents, human and artificial alike, possess some fundamental motivational attitude — whatever it may actually be.

Secondly, in addition to being faithful to insights gained in psychology, we feel that it is also important to pay attention to those gained in analytical philosophy. More specifically, we are of the opinion that the modelling of *practical reasoning* should be part of any formalisation of motivational attitudes that pretends to be an adequate one. The term ‘practical reasoning’ dates back to Aristotle, and refers to the process through which (human) agents conclude

that they should perform certain actions in order to bring about some of the things that they like to be the case. It seems very likely that for autonomous agents in AI applications, which have to act (autonomously) to achieve some of their goals, practical reasoning accounts for the most essential and most frequently used kind of information processing. Hence an adequate formalisation of motivational attitudes should pay at least some attention to this kind of reasoning.

The third essential facet of any formalisation of motivational attitudes consists of the modelling of the act of *selecting*: agents have to make choices among the things they like to be the case, thereby deciding which of these they will try to achieve next. For it might be impossible to satisfy all of an agent's wishes simultaneously, since these wishes become easily either analytically inconsistent or incompatible given the agent's resources.

In our opinion all of the aspects mentioned above should be present in an adequate formalisation of motivational attitudes, and they are indeed so in the one that is presented in this paper. The notions that are essential in this formalisation are *wishes*, *goals* and *commitments*. Of these, wishes constitute the primitive motivational attitude that models the things that an agent likes to be the case. As such, wishes naturally range over propositions, corresponding to the idea that agents wish for certain aspects of the world. We formalise wishes through a normal modal operator, i.e. an operator validating just the so-called K-axiom and the necessitation rule. Agents set their goals by selecting among their wishes. However, agents are not allowed to select arbitrary wishes as their goals, but instead may only select wishes that are unfulfilled yet implementable. Whenever an agent knows that it has some goal, it may commit itself to any action that it knows to be correct and feasible with respect to the goal. This act of committing to an action is itself formalised as a special kind of action. Commitments to actions are in general to persist until all of the goals for which the commitment was made are fulfilled. Having said so, agents should not be forced to remain committed to actions that have either become useless in that they do not lead to fulfilment of any goal any more, or impossible in that the agent no longer knows that it has the opportunity and ability to perform the action. Phrased differently, an agent should be allowed to uncommit itself whenever an action is no longer known to be correct and feasible with respect to one of the agent's goals.

To formalise wishes, goals and commitments and their associated concepts, we introduce a modal operator modelling wishes, operators modelling implementability, (made) choices and (made) commitments, and action constructors modelling the acts of selecting, committing and uncommitting.

3.1. DEFINITION. To define the language  $L^C$ , the alphabet is extended with the wish operator  $\mathbf{W}_-$ , the implementability operator  $\diamond_-$ , the selected operator  $\mathbf{C}_-$ , the commitment operator  $\mathbf{Committed}_-$  and the action constructors  $\mathbf{select}_-$ ,  $\mathbf{commit\_to}_-$  and  $\mathbf{uncommit}_-$ .

The acts of committing and uncommitting are of an essentially different nature than the regular actions, execution of which changes the state of the world. Through the former actions agents (un)commit themselves to actions of the latter kind. Intuitively it does not make much sense to allow agents to commit themselves to making commitments: it is not at all clear how a statement like '*i* is committed to commit itself to do  $\alpha$ ' is to be interpreted. Also statements like 'it is implementable for agent *i* to become committed' seem to be of a rather questionable nature. To avoid these kinds of counterintuitive situations, we define the language  $L^C$  on top of the language  $L$  that we defined previously. That is, the operators modelling wishes,

implementability and selections are defined in such a way that they range over formulae from  $L$  rather than those from  $L^C$ . The operator modelling the commitments that an agent has made is defined to range over the actions from  $Ac$ , the class of actions associated with  $L$ , and not over  $Ac^C$ . Analogously, the special actions in  $Ac^C$ , as there are the action modelling the act of selecting and those modelling the making and undoing of commitments, take as arguments elements from  $L$  and  $Ac$  rather than  $L^C$  and  $Ac^C$ .

3.2. DEFINITION. The language  $L$  and the class  $Ac$  of actions are as in Definition 2.1, i.e.  $L$  is the smallest superset of  $\Pi$  closed under the core clauses of Definition 2.1 and  $Ac$  is the smallest superset of  $At$  satisfying the core clauses of Definition 2.1 again.

The language  $L^C$  is the smallest superset of  $\Pi$  such that the core clauses of Definition 2.1 are validated and furthermore

- if  $\varphi \in L$  and  $i \in A$  then  $\mathbf{W}_i\varphi \in L^C$
- if  $\varphi \in L$  and  $i \in A$  then  $\diamond_i\varphi \in L^C$
- if  $\varphi \in L$  and  $i \in A$  then  $\mathbf{C}_i\varphi \in L^C$
- if  $\alpha \in Ac$  and  $i \in A$  then  $\mathbf{Committed}_i\alpha \in L^C$

The class  $Ac^C$  is the smallest superset of  $At$  closed under the core clauses of Definition 2.1 and such that

- if  $\varphi \in L$  then  $\mathbf{select}\varphi \in Ac^C$
- if  $\alpha \in Ac$  then  $\mathbf{commit\_to}\alpha \in Ac^C$
- if  $\alpha \in Ac$  then  $\mathbf{uncommit}\alpha \in Ac^C$

3.3. DEFINITION. For  $i \in A$ ,  $\alpha \in Ac^C$  and  $\varphi \in L^C$ , the abbreviations  $\mathbf{PracPoss}_i(\alpha, \varphi)$  and  $\mathbf{Can}_i(\alpha, \varphi)$  are defined as in the language  $L$ .

## 4. Semantics

As is usual for a modal language we use Kripke-style semantics. Essentially, our models will consist of a set of possible worlds, a truth assignment function for the propositional symbols, and accessibility relations for the modal operators that we consider. For the core language defined in Definition 2.1, this means that we need accessibility relations for the knowledge operator and the operator  $\langle \text{do}_i(\alpha) \rangle$ . For technical convenience, we will use for the latter a function rather than a relation. Finally, we need a function to indicate the agent's capabilities at some possible world.

4.1. DEFINITION. The class  $\mathbf{M}$  of models contains all  $M = \langle S, \pi, R, \mathbf{r}_0, \mathbf{c}_0 \rangle$  where

- $S$  is a set of possible worlds, or states.
- $\pi : \Pi \times S \rightarrow \{\mathbf{0}, \mathbf{1}\}$  assigns a truth value to propositional symbols in states.
- $R : A \rightarrow \wp(S \times S)$  is a function that yields the epistemic accessibility relations for a given agent. It is demanded that  $R(i)$  is an equivalence relation for all  $i$ . We define  $[s]_{R(i)}$  to be  $\{s' \in S \mid (s, s') \in R(i)\}$ , the  $R(i)$ -equivalence class of  $s$ .
- $\mathbf{r}_0 : A \times At \rightarrow (S \cup \{\emptyset\}) \rightarrow (S \cup \{\emptyset\})$  is such that  $\mathbf{r}_0(i, a)(s)$  yields the (possibly empty) state transition in  $s$  caused by the event  $\text{do}_i(a)$ . For technical reasons we stipulate that always  $s \notin \mathbf{r}_0(i, a)(s)$ . This means that a successful performance of an atomic action

always results in a state transition to *another* state in the model. Of course, it may happen that this state satisfies the same formulas as the original one.<sup>2</sup>

- $c_0 : A \times At \rightarrow (S \cup \{\emptyset\}) \rightarrow \{\mathbf{0}, \mathbf{1}\}$  is the capability function such that  $c_0(i, a)(s)$  indicates whether the agent  $i$  is capable of performing the action  $a$  in  $s$ .

Furthermore, we require the following property:  $\forall s, s', t, t' \in S, a, a' \in At^+ : [(s, s') \in R(i) \ \& \ s \in r_0(i, a)(t) \ \& \ s' \in r_0(i, a')(t')] \Rightarrow [s = s' \ \& \ t = t' \ \& \ a = a']$ .

Informally, the last requirement in the above definition stipulates that in each epistemic equivalence class there is at most one state that is reachable by means of (performing) an action. Thus epistemic alternatives are strictly separated from states resulting from the execution of actions. Another way of viewing this is that actions make transitions between epistemic equivalence classes, representing the epistemic states of the agent. Thus by performing actions agents jump from the one epistemic state to another. (And the state the action results in, the ‘core’ of the epistemic equivalence class, so to speak, represents the *actual* state the agent ends up in after performing that action. Together with the requirement that successful actions really change state, this condition on the models will enable a proper recording of the agent’s commitments (by means of an ‘agenda’ function) in the sequel.

The models for the language  $L^C$  are further equipped with elements used to interpret the agents’ wishes, selections and commitments. Wishes are interpreted through an accessibility relation on worlds that denotes worlds that are (more) desirable from the agent’s point of view. Selections are straightforwardly interpreted through a set of formulae that denotes the choices that an agent has made. From a formal point of view, this set acts as a kind of *awareness* on the wishes of an agent, thereby ensuring an intuitively acceptable behaviour of goals. Originally, Fagin & Halpern [9] introduced the idea of awareness sets as a means to solve the so-called problems of logical omniscience. As we will see in Section 6, the effect of the selection sets on the behaviour of goals is similar to that of the awareness sets on the properties of knowledge. The agents’ commitments are interpreted by means of the agenda function, which yields for each agent in every state the commitments that it has made and is up to. Detailed accounts of the respective interpretations are given in the following sections.

**4.2. DEFINITION.** A model  $M$  for the language  $L^C$  is a tuple containing the core elements of Definition 4.1, the functions  $W : A \rightarrow \wp(S \times S)$ , which determines the desirability relation of an agent in a state, and  $C : A \times S \rightarrow \wp(L)$  denoting the choices made by an agent in a state, and a function  $\text{Agenda} : A \times S \rightarrow \wp(\text{Acseq})$ , which records the commitments of agents.

We interpret the acts of selecting, committing and uncommitting as model-transformations. The act of selecting changes a model by affecting the set of choices, and the act of (un)committing transforms the agent’s agenda. To account for these modifications, we introduce the set of possible result models of a given model for  $L^C$ .

**4.3. DEFINITION.** Let  $M \in \mathbf{M}^C$  be some model for  $L^C$ . The class  $\mathbf{M}_{\sim}^C \subseteq \mathbf{M}^C$  contains all models that (possibly) differ from  $M$  only in the  $C$  or the Agenda functions.

---

<sup>2</sup>This might be understood in the following sense: even an action that has no effect on the truth of the formulas such as a ‘skip’ action, still results in a state where the *history* of performed actions is changed. This idea is very similar to the notion of a (Herbrand) situation in the situation calculus ([32, 20, 34, 18]), and can be made formal by considering propositions (or variables) recording the history. In order to not complicate our models further we will not do so in this paper.

In our interpretation of actions as given in Definition 4.4 below, we generalise the standard paradigm of actions as state-transitions [12] by interpreting actions as transitions between pairs (Model, State) rather than transitions between states *per se*. Using this more general interpretation we can both account for regular actions that cause a transition between states upon execution, and special actions that transform models. Among the special actions that we considered elsewhere [24] were those modelling informational attitudes such as observations and communication; in Section 7 we will formalise the acts of committing and uncommitting in a similar way.

Furthermore, to keep in line with the idea expressed in Definition 4.1 that the performance of actions, when successful, always leads to different states, we sometimes need ‘copies’ of states, that satisfy the same formulas. In the following we assume the availability of such copies in a model  $M$ . For a state  $s$ , we denote by  $\bar{s}$  a designated state  $\neq s$  for which it holds that, for all formulas  $\varphi \in L$ ,  $M, s \models \varphi \Leftrightarrow M, \bar{s} \models \varphi$ .

4.4. DEFINITION. The binary relation  $\models$  between a formula from  $L$  and a pair  $M, s$  consisting of a model  $M \in \mathbf{M}$  and a state  $s$  in  $M$ , is for  $\varphi$  a propositional symbol, a negation, or a disjunction inductively defined as usual. For the other cases  $M, s \models \varphi$  is defined by:

$$\begin{aligned} M, s \models \mathbf{K}_i \varphi & \Leftrightarrow \forall s' \in S((s, s') \in R(i) \Rightarrow M, s' \models \varphi) \\ M, s \models \langle \text{do}_i(\alpha) \rangle \varphi & \Leftrightarrow \exists M', s'(M', s' \in \mathbf{r}(i, \alpha)(M, s) \ \& \ M', s' \models \varphi) \\ M, s \models \mathbf{A}_i \alpha & \Leftrightarrow \mathbf{c}(i, \alpha)(M, s) = \mathbf{1} \end{aligned}$$

where  $\mathbf{r}$  and  $\mathbf{c}$  are defined as follows:

$$\begin{aligned} \mathbf{r}(i, a)(M, s) & = M, \mathbf{r}_0(i, a)(s) \\ \mathbf{r}(i, \text{confirm } \varphi)(M, s) & = M, \bar{s} \text{ if } M, s \models \varphi \text{ and } \emptyset \text{ otherwise} \\ \mathbf{r}(i, \alpha_1; \alpha_2)(M, s) & = \mathbf{r}(i, \alpha_2)(\mathbf{r}(i, \alpha_1)(M, s)) \\ \mathbf{r}(i, \text{if } \varphi \text{ then } \alpha_1 & \text{ else } \alpha_2 \text{ fi})(M, s) & = \mathbf{r}(i, \alpha_1)(M, \bar{s}) \text{ if } M, s \models \varphi \text{ and} \\ & \mathbf{r}(i, \alpha_2)(M, \bar{s}) \text{ otherwise} \\ \mathbf{r}(i, \text{while } \varphi \text{ do } \alpha \text{ od})(M, s) & = M', s' \text{ such that } \exists k \in \mathbb{N} \exists M_0, s_0 \dots \exists M_k, s_k \\ & (M_0, s_0 = M, s \ \& \ M_k, s_k = M', s' \ \& \ \forall j < k \\ & (M_{j+1}, s_{j+1} = \mathbf{r}(i, \text{confirm } \varphi; \alpha)(M_j, s_j)) \ \& \ M', s' \models \neg \varphi) \\ & = \emptyset \end{aligned}$$

where  $\mathbf{r}(i, \alpha)(\emptyset) = \emptyset$

and

$$\begin{aligned} \mathbf{c}(i, a)(M, s) & = \mathbf{c}_0(i, a)(s) \\ \mathbf{c}(i, \text{confirm } \varphi)(M, s) & = \mathbf{1} \text{ if } M, s \models \varphi \text{ and } \mathbf{0} \text{ otherwise} \\ \mathbf{c}(i, \alpha_1; \alpha_2)(M, s) & = \mathbf{c}(i, \alpha_1)(M, s) \ \& \ \mathbf{c}(i, \alpha_2)(\mathbf{r}(i, \alpha_1)(M, s)) \\ \mathbf{c}(i, \text{if } \varphi \text{ then } \alpha_1 & \text{ else } \alpha_2 \text{ fi})(M, s) & = \mathbf{c}(i, \text{confirm } \varphi; \alpha_1)(M, s) \text{ or} \\ & \mathbf{c}(i, \text{confirm } \neg \varphi; \alpha_2)(M, s) \\ \mathbf{c}(i, \text{while } \varphi \text{ do } \alpha \text{ od})(M, s) & = \mathbf{1} \text{ if } \mathbf{c}(i, (\text{confirm } \varphi; \alpha)^k; \text{confirm } \neg \varphi)(M, s) = \mathbf{1} \\ & \text{for some } k \in \mathbb{N} \text{ and } \mathbf{0} \text{ otherwise} \end{aligned}$$

where  $\mathbf{c}(i, \alpha)(\emptyset) = \mathbf{0}$

Validity in a model and in a class of models is defined as usual.

With regard to the abilities of agents, the motivation for the choices made in Definition 4.4 is the following. The definition of  $\mathbf{c}(i, \text{confirm } \varphi)$  expresses that an agent is able to get confirmation for a formula  $\varphi$  iff  $\varphi$  holds. An agent is capable of performing a sequential composition  $\alpha_1; \alpha_2$  iff it is capable of performing  $\alpha_1$  (now), and it is capable of executing  $\alpha_2$  after it has performed  $\alpha_1$ . An agent is capable of performing a conditional composition, if



either it is able to get confirmation for the condition and thereafter perform the then-part, or it is able to confirm the negation of the condition and perform the else-part afterwards. An agent is capable of performing a repetitive composition `while  $\varphi$  do  $\alpha$  od` iff it is able to perform the action `(confirm  $\varphi$ ;  $\alpha_1$ )k; confirm  $\neg\varphi$`  for some natural number  $k$ , i.e. it is able to perform the  $k$ th unwinding of the while-loop.

## 5. Formalising wishes

Wishes are the most primitive, fundamental motivational attitudes, i.e. *in ultimo* agents are motivated to fulfil their wishes. As mentioned in Section 3, we formalise wishes through a plain normal modal operator, i.e. wishes are straightforwardly interpreted as a necessity operator over the accessibility relation  $W$ .

5.1. DEFINITION. The binary relation  $\models^C$  between a formula in  $L^C$  and a pair  $M, s$  consisting of a model  $M$  for  $L^C$  and a state  $s$  in  $M$  is for wishes defined as follows:

$$M, s \models^C \mathbf{W}_i\varphi \Leftrightarrow \forall s' \in S((s, s') \in W(i) \Rightarrow M, s' \models^C \varphi)$$

It is well-known that normal modal operators have certain properties that are occasionally considered undesirable for the commonsense notions that they are intended to formalise. For example, although the formal notions of knowledge and belief are closed under logical consequence, this property will in general not hold for human knowledge and belief (although it will for instance hold for the information that is recorded in a database, and it can also be defended for the knowledge and belief of an artificial agent). When formalising motivational attitudes the undesired properties induced by closure under logical consequence become even more pregnant. For agents do in general not desire all the logical consequences of their wishes, nor do they consider the logically inevitable to be among their goals. For example, an agent that wants its teeth to be restored will in general not want or wish for the pain that inevitably accompanies such a restoration. And although the sun rises in the east there will hardly be an agent that desires this to be the case. The problem embodied by the former example is known as the *side-effect* problem; the problem that all logical tautologies are wishes (goals) of an agent is known as the *transference* problem. Both in syntactical shape as in meaning, these problems are closely related to the problems of logical omniscience that have plagued formalisations of informational attitudes for many years. In terms of our framework, seven of the most (in)famous problems of logical omniscience can be formulated as follows.

5.2. DEFINITION. Let  $\varphi, \psi \in L$  be formulae, and let  $\mathbf{X}$  be some operator.

- $\models \mathbf{X}\varphi \wedge \mathbf{X}(\varphi \rightarrow \psi) \rightarrow \mathbf{X}\psi$  LO1
- $\models \varphi \Rightarrow \models \mathbf{X}\varphi$  LO2
- $\models \varphi \rightarrow \psi \Rightarrow \models \mathbf{X}\varphi \rightarrow \mathbf{X}\psi$  LO3
- $\models \varphi \leftrightarrow \psi \Rightarrow \models \mathbf{X}\varphi \leftrightarrow \mathbf{X}\psi$  LO4
- $\models (\mathbf{X}\varphi \wedge \mathbf{X}\psi) \rightarrow \mathbf{X}(\varphi \wedge \psi)$  LO5
- $\models \mathbf{X}\varphi \rightarrow \mathbf{X}(\varphi \vee \psi)$  LO6
- $\models \neg(\mathbf{X}\varphi \wedge \mathbf{X}\neg\varphi)$  LO7

Properties LO1 and LO3 as given in Definition 5.2 capture the side-effect problem, and property LO2 captures the transference problem. Of the other properties given not all are equally harmful when formalising wishes. In our opinion, property LO4 is not that harmful,

and could even be considered desirable, dependent on the demands for rationality that one is willing to make. Property LO5, which we like to think of as representing ‘the problem of *unrestricted combining*’, is in general undesirable when formalising motivational attitudes. This is for instance shown by the example of a (human) agent that likes both a sandwich with peanut butter and one with Camembert cheese, but not together. Property LO6, for which we coin the term ‘the problem of *unrestricted weakening*’, is a special instantiation of the side-effect problem. That this property is undesirable is shown by the example of an agent desiring itself to be painted green, without desiring being green or being crushed under a steam roller<sup>3</sup>. Property LO7 is unacceptable for certain kinds of motivational attitudes but a necessity for others. It is for instance perfectly possible for agents to have contradicting wishes<sup>4</sup>, but it seems hardly rational to allow agents to try and fulfil these conflicting wishes simultaneously. Thus, whereas the absence of LO7 is essential when formalising wishes, the presence is when formalising goals.

It turns out that our formalisation of wishes validates all but one of the properties of logical omniscience.

**5.3. PROPOSITION.** *All of the properties of logical omniscience formalised in Definition 5.2, with the exception of LO7, are valid for the  $\mathbf{W}_i$  operator.*

Although we argued against the properties of logical omniscience when formalising motivational attitudes, we do not consider it a serious problem that our formalisation of wishes validates (almost all of) these properties. For these wishes are both *implicit* in the terminology of Levesque [19] and *passive* in the sense of Castelfranchi *et al.* [3]. Being implicit, it will not be the case that agents *explicitly* desire all of their wishes<sup>5</sup>. Being passive, wishes in themselves do not *actively* influence the course of action that an agent is going to take. Through the act of *selecting*, agents turn some of their implicit, passive wishes into explicit, active goals. Hence even though an agent implicitly and passively desires all logical consequences of one of its wishes, it will not do so explicitly and actively. Therefore Proposition 5.3 is not taken to represent a severe problem for a formalisation of (implicit and passive) wishes, whereas it would for a formalisation of (explicit and active) goals. In the following section it will be shown how the properties of logical omniscience are avoided for goals.

## 6. Setting goals

As remarked previously, an agent’s goals are not primitive but induced by its wishes. Basically, an agent selects among its (implicit and passive) wishes those that it (explicitly and actively) aims to fulfil. Given the rationality of agents, these selected wishes should be both unfulfilled and implementable: it does not make sense for an agent to try and fulfil a wish that either already has been fulfilled or for which fulfilment is not a practical possibility. We do not take the latter constraint too stringently, i.e. we only demand wishes to be individually

---

<sup>3</sup>The problem of unrestricted weakening is intuitively related to the Ross’s paradox [33], well-known in deontic logic [1, 25]. The standard counterexample towards the desirability of LO6 in a deontic context, where the operator  $\mathbf{X}$  is interpreted as ‘being obliged to’, is that of an agent that is obliged to mail a letter while (intuitively) not being obliged to either mail the letter or burn it.

<sup>4</sup>Even stronger, human agents will almost always suffer from conflicts between their wishes.

<sup>5</sup>For the implicit belief that, in combination with awareness, constitutes explicit belief in the approach of Fagin & Halpern [9], it is also considered unproblematic that the properties of logical omniscience are validated.

implementable without requiring a simultaneous implementability of all chosen wishes. However, if desired, constraints like simultaneous implementability are easily formulated. The act of selecting is treated as a fully-fledged action by defining the opportunity, ability and result of selecting. Informally, an agent has the *opportunity* to select any of its wishes, corresponding to the idea that choices are only restricted by the elements among which is to be chosen. However, an agent is *capable* of selecting only those formulae that are unfulfilled and implementable, which can be thought of as it having a built-in aversion against selecting fulfilled or practically impossible formulae. The *result* of a selection will consist of the selected formula being marked chosen.

The notion of unfulfilledness is straightforwardly formalised as ‘not holding’, i.e. a formula  $\varphi$  is unfulfilled in a state  $s$  of some model  $M$  if and only if  $M, s \not\models^C \varphi$ . Defining implementability is a little more elaborate. Roughly speaking, we define a formula  $\varphi$  to be implementable for an agent  $i$ , denoted by  $\diamond_i \varphi$ , if  $i$  has the practical possibility to fulfil  $\varphi$  by performing an appropriate sequence of atomic actions<sup>6</sup>.

6.1. DEFINITION. The binary relation  $\models^C$  between a formula in  $L^C$  and a pair  $M, s$  consisting of a model  $M$  for  $L^C$  and a state  $s$  in  $M$  is for implementability formulae defined by:

$$M, s \models^C \diamond_i \varphi \Leftrightarrow \exists k \in \mathbb{N} \exists a_1, \dots, a_k \in \text{At}(M, s) \models^C \mathbf{PracPoss}_i(a_1; \dots; a_k, \varphi)$$

Having defined unfulfilledness and implementability, we can now formally introduce the **select** action.

6.2. DEFINITION. For  $M \in \mathbf{M}^C$  with state  $s$ ,  $i \in A$  and  $\varphi \in L$  we define:

$$\mathbf{r}^C(i, \mathbf{select} \varphi)(M, s) = \begin{cases} \emptyset & \text{if } M, s \models^C \neg \mathbf{W}_i \varphi \\ \mathbf{choose}(i, \varphi)(M, s), s & \text{if } M, s \models^C \mathbf{W}_i \varphi \end{cases}$$

where for  $M = \langle S, \pi, R, \mathbf{r}_0, \mathbf{c}_0, W, C, \text{Agenda} \rangle$  we define

$$\begin{aligned} \mathbf{choose}(i, \varphi)(M, s) &= \langle S, \pi, R, \mathbf{r}_0, \mathbf{c}_0, W, C', \text{Agenda} \rangle \text{ with} \\ C'(i', s') &= C(i', s') \text{ if } i \neq i' \text{ or } s \neq s' \\ C'(i, s) &= C(i, s) \cup \{\varphi\} \end{aligned}$$

$$\mathbf{c}^C(i, \mathbf{select} \varphi)(M, s) = \mathbf{1} \Leftrightarrow M, s \models^C \neg \varphi \wedge \diamond_i \varphi$$

The binary relation  $\models^C$  between a formula in  $L^C$  and a pair  $M, s$  consisting of a model  $M$  for  $L^C$  and a state  $s$  in  $M$  is for choices defined by:

$$M, s \models^C \mathbf{C}_i \varphi \Leftrightarrow \varphi \in C(i, s)$$

The definition of  $\mathbf{r}^C$  for the selection actions indeed provides for a correct model transformation.

---

<sup>6</sup>As was pointed out by Maarten de Rijke, defining the implementability operator in this way makes it a kind of ‘dual master modality’ (cf. [10, 38]). A formula consisting of a formula  $\varphi$  prefixed by the master modality is true in some state  $s$  of a model iff  $\varphi$  holds at all states that are reachable by any finite sequence of transitions from  $s$ . Such a formula is false iff there is some state  $s'$ , reachable by some finite sequence of transitions from  $s$ , at which  $\varphi$  does not hold. This indeed makes our implementability modality to be a dual master modality.

6.3. PROPOSITION. For all  $M \in \mathbf{M}^C$  with state  $s$ , for all  $i \in A$  and  $\varphi \in L$ , if  $M', s = \mathbf{r}^C(i, \mathbf{select} \varphi)(M, s)$  then  $M' \in \mathbf{M}_{\sim}^C$ .

Besides being correct in that well-defined models are transformed into well-defined models, our formalisation of the act of selecting is also correct with respect to minimal change. That is, the change caused by selecting some formula is minimal given that the formula is to be marked chosen, which implies that our formalisation of selections does not suffer from the frame problem. The following proposition provides a (partial) formalisation of this property.

6.4. PROPOSITION. For all  $M \in \mathbf{M}^C$  with state  $s$ , for all  $i \in A$  and  $\varphi \in L$ , if  $M', s = \mathbf{r}^C(i, \mathbf{select} \varphi)(M, s)$  then for all states  $s'$  in  $M$ ,  $M, s' \models^C \psi$  iff  $M', s' \models^C \psi$ , for all  $\psi \in L$ .

Proposition 6.4 states that all formulae from  $L$  are interpreted identically in a model  $M$  and in the one resulting from selecting some formula in an arbitrary state of  $M$ . As a direct consequence of this proposition we have the following corollary, which states that the interpretation of wishes and implementability formulae persists under selecting some formula.

6.5. COROLLARY. For all  $M \in \mathbf{M}^C$  with state  $s$ , for all  $i \in A$  and  $\varphi \in L$ , if  $M', s = \mathbf{r}^C(i, \mathbf{select} \varphi)(M, s)$  then for all states  $s'$  in  $M$  and all  $\psi \in L$ :

- $M, s' \models^C \mathbf{W}_i \psi \Leftrightarrow M', s' \models^C \mathbf{W}_i \psi$
- $M, s' \models^C \diamond_i \psi \Leftrightarrow M', s' \models^C \diamond_i \psi$

Having defined wishes and selections, one might be tempted to straightforwardly define goals to be selected wishes, i.e.  $\mathbf{Goal}_i \varphi =^{\text{def}} \mathbf{W}_i \varphi \wedge \mathbf{C}_i \varphi$ . This definition is however not adequate to formalise the idea of goals being selected, *unfulfilled*, *implementable* wishes. The reason for this is that in well-defined models from  $\mathbf{M}^C$  no relation is imposed between ‘being selected’ and ‘being unfulfilled and implementable’, i.e. one is not prevented by Definition 4.2 to come up with a well-defined model  $M$  in which for certain  $i$  and  $s$  the set  $C(i, s)$  contains formulae  $\varphi$  that are either fulfilled or not implementable. We see basically two ways of solving this problem, a semantical and a syntactical one. Semantically one could restrict the set of well-defined models for  $L^C$  to those in which the set  $C(i, s)$  contains for all agents  $i$  and states  $s$  only unfulfilled and implementable formulae, thereby ensuring beforehand that goals are unfulfilled and implementable when using the definition suggested above. Syntactically one could define goals to be only those selected wishes that are indeed unfulfilled and implementable. Hence instead of (semantically) restricting the set of well-defined models for  $L^C$  one (syntactically) expands the definition of goals. Although both the semantic and the syntactic approach are equally well applicable, we will restrict ourselves here to pursuing the syntactic one. Therefore, goals are defined to be those wishes that are unfulfilled, implementable and selected.

6.6. DEFINITION. The  $\mathbf{Goal}_i$  operator is for  $i \in A$  and  $\varphi \in L$  defined by:

$$\mathbf{Goal}_i \varphi =^{\text{def}} \mathbf{W}_i \varphi \wedge \neg \varphi \wedge \diamond_i \varphi \wedge \mathbf{C}_i \varphi$$

As mentioned above, the goals of agents, being the explicit and active notions that they are, are not to validate the properties of logical omniscience as formalised in Definition 5.2. Fortunately, though not surprisingly, this indeed turns out to be the case when defining goals as in Definition 6.6.

6.7. PROPOSITION. None of the properties of logical omniscience formalised in Definition 5.2, with the exception of *LO7*, is valid for the  $\mathbf{Goal}_i$  operator.

The only property of logical omniscience satisfied by the goal operator, viz. LO7, formalises the idea that an agent's goals are consistent. This is a highly desirable property for rational creatures. For although it is quite possible for a rational agent to have contradictory wishes, it is rather irrational to try and fulfil these simultaneously.

Besides invalidating the undesired ones among the properties of logical omniscience, particularly those embodying the side-effect and transference problem, our definition of goals and selections has some other pleasant and desirable features. The following proposition formalises some of these features together with some properties characterising the act of selecting.

6.8. PROPOSITION. *For all  $i \in \mathbf{A}$  and  $\varphi \in \mathbf{L}$  we have:*

1.  $\models^C \mathbf{W}_i\varphi \leftrightarrow \langle \text{do}_i(\text{select } \varphi) \rangle \top$
2.  $\models^C \langle \text{do}_i(\text{select } \varphi) \rangle \top \leftrightarrow \langle \text{do}_i(\text{select } \varphi) \rangle \mathbf{C}_i\varphi$
3.  $\models^C \neg \mathbf{A}_i \text{select } \varphi \rightarrow [\text{do}_i(\text{select } \varphi)] \neg \mathbf{Goal}_i\varphi$
4.  $\models^C \mathbf{PracPoss}_i(\text{select } \varphi, \top) \leftrightarrow \langle \text{do}_i(\text{select } \varphi) \rangle \mathbf{Goal}_i\varphi$
5.  $\models^C \varphi \Rightarrow \models^C \neg \mathbf{Goal}_i\varphi$
6.  $(\varphi \rightarrow \psi) \rightarrow (\mathbf{Goal}_i\varphi \rightarrow \mathbf{Goal}_i\psi)$  is not for all  $\varphi, \psi \in \mathbf{L}$  valid
7.  $\mathbf{K}_i(\varphi \rightarrow \psi) \rightarrow (\mathbf{Goal}_i\varphi \rightarrow \mathbf{Goal}_i\psi)$  is not for all  $\varphi, \psi \in \mathbf{L}$  valid

The first item of Proposition 6.8 states that agents have the opportunity to select all, and nothing but, their wishes. The second item formalises the idea that every choice for which an agent has the opportunity results in the selected wish being marked chosen. In the third item it is stated that whenever an agent is unable to select some formula, then selecting this formula will not result in it becoming one of its goals. The related item 4 states that all, and nothing but, practically possible selections result in the chosen formula being a goal. The fifth item provides a strengthening of the invalidation of the second property of logical omniscience, which embodies the transference problem. It states that no logically inevitable formula qualifies as a goal. Hence whenever a formula is valid this does not only not necessarily imply that it is a goal but it even necessarily implies that it is not. The last two items of Proposition 6.8 are related to the avoidance of the transference problem, and state that goals are neither closed under implications nor under known implications.

## 7. Formalising commitments

The last part of our formalisation of motivational attitudes concerns the agents' commitments. Commitments to actions represent promises to perform these actions, i.e. an agent that is committed to an action has promised itself to perform the action. As mentioned above, commitments may be made to plans for goals, i.e. whenever an agent is committed it should be to an action that is correct and feasible to bring about at least one of its goals.

Not only do we formalise this static aspect of made commitments, but we also consider the dynamic aspect of making and undoing commitments. The act of committing is related to, and can be seen as, an elementary implementation of practical reasoning, the process through which agents decide that they should perform certain actions (their ought-to-do's) on the basis of their wishes, desires or goals (their ought-to-be's). Ever since Aristotle, the study of practical reasoning has formed a major constituent of the research in analytical philosophy [31]. According to Von Wright [42], the essence of practical reasoning is best captured by the following syllogism:

*i* intends to make it true that  $\varphi$   
*i* thinks that, unless it does  $\alpha$ , it will not achieve this  
Therefore *i* intends to do  $\alpha$ .

The simplified version of practical reasoning that we aim to formalise through the act of committing can be described by the following syllogism,

*i* knows that  $\varphi$  is one of its *goals*  
*i* knows that  $\alpha$  is *correct* and *feasible* with respect to  $\varphi$   
Therefore *i* has the *opportunity to commit* itself to  $\alpha$

which corresponds to the idea that commitments may be made to actions that are known to be correct and feasible to achieve some of the agent’s known goals.

Commitments are formalised through the **Committed<sub>i</sub>** operator: **Committed<sub>i</sub>** $\alpha$  denotes that agent *i* is committed to the action  $\alpha$ . The act of committing is modelled by the (special) action **commit<sub>to</sub>**: **commit<sub>to</sub>** $\alpha$  represents the act of committing to the (regular) action  $\alpha$ . As mentioned in Section 3, commitments, though in general persistent, should not be maintained when having become useless or impossible, i.e. agents should have the possibility to undo useless or impossible commitments. This act of uncommitting is formalised by the **uncommit** action: **uncommit** $\alpha$  denotes the act of undoing the commitment to the action  $\alpha$ . In the sequel we successively formalise the act of committing, the commitments that have been made, and the act of uncommitting.

### 7.1. Getting committed

The act of committing, though of a special nature compared to other actions, is treated as a fully-fledged action, i.e. we define what it means to have the ability or opportunity to commit, and what the result of committing is. To start with the latter notion, given the relation between the infinitive ‘to commit’ and the past participle ‘committed’, it seems rather obvious that the act of committing should result in the agent being committed. Determining when an agent has the opportunity to perform a **commit<sub>to</sub>** $\alpha$  action is equally obvious, for it is inspired by the syllogism describing our version of practical reasoning given above. Hence agent *i* has the opportunity to perform the action **commit<sub>to</sub>** $\alpha$  if and only if it knows that  $\alpha$  is correct and feasible to bring about one of its goals. This leaves to determine the constituents of the ability of an agent to commit itself. Our definition of this ability is inspired by the observation that situations where agents are committed to two (or more) essentially different actions are highly problematic. Since ‘being committed to  $\alpha$ ’ intuitively corresponds to ‘having promised (to oneself) to perform  $\alpha$  next’, it is unclear how to interpret the case where an agent is committed to two different actions. Should both actions be performed simultaneously? But what does it mean that actions are performed simultaneously? Are they performed concurrent, interleaved or in parallel? Or should the actions be performed sequentially? If so, in which order? And what then if the commitment to perform one action does not persist under execution of the other action? It is clear that these questions have no unique answers. Here we shall assume simply that agents do not have multiple commitments. One way to ensure this is to let an agent have the ability to commit itself only if it is not up to any previously made commitments, i.e. an agent is capable to commit only if it is not already committed. (However, it should also be clear that in principle our framework allows

other, more complicated situations, if one would be interested to model that. However, since some of the results in the sequel depend on the present choice, these should be reconsidered then.)

As mentioned previously, an agent's commitments are interpreted by means of the so-called agenda function. The idea is that this function yields, for a given agent and a given state, the actions that the agent is committed to. Whenever an agent successfully commits itself to an action the agent's agenda is updated accordingly. The actual formal definition capturing this fairly unsophisticated idea is itself rather complicated. The reason for this lies in various desiderata that commitments and the act of committing should meet.

The first of these desiderata is that commitments should be known, i.e. agents should be aware of the commitments that they have made. To bring about this knowledge of commitments, epistemic equivalence classes rather than states are considered in an agenda update. Thus whenever agent  $i$  commits itself to action  $\alpha$  in some state  $s$  of a model, the agenda of all states  $s'$  that are epistemically equivalent with  $s$  is updated appropriately.

The second and very important desideratum imposed on commitments is that they behave compositionally correct, i.e. the commitment to a composite action is linked in a rational way to commitments to its constituents. It is for example desirable that an agent that is committed to an action `if  $\varphi$  then  $\alpha_1$  else  $\alpha_2$  fi` is also committed to  $\alpha_1$  whenever it knows that  $\varphi$  holds, and that an agent committed to the action  $\alpha_1; \alpha_2$  is (currently) committed to  $\alpha_1$  and committed to  $\alpha_2$  in the state of affairs that results from executing  $\alpha_1$ .

To bring about rational behaviour of commitments with respect to sequentially composed actions the actual update does not just concern the epistemic equivalence class of the current state, but also that of all the states that lay alongside the execution trajectory of the action. For example, if an agent  $i$  commits itself to  $\alpha_1; \alpha_2$  in the state  $s$  of some model, then the epistemic equivalence class of  $s$  is updated with the commitment to  $\alpha_1$ , and the epistemic equivalence class of the state  $s''$  that results from executing  $\alpha_1$  in some  $s'$  that is an element of the epistemic equivalence of  $s$  is updated with the commitment to  $\alpha_2$ .

Although we will put in the agent's agenda (syntactical) actions that represent actions it is committed to, we will, of course, also need semantical entities to link these actions to what actually happens when an agent performs (part of) the action in its agenda. To this end we will use some notions that are well-known from the area of the semantics of programming languages, viz. computation sequences, computation runs and transition relations. Let us start with the introduction of computation sequences and runs.

Since the actions from  $\text{Ac}$  are deterministic, for each event built out of these actions there is at most one finite computation sequence which consists of the (semi-)atomic actions that occur in the halting executing of the event. Or phrased differently, the set of finite computation runs of a given event  $\text{do}_i(\alpha)$  is either empty or a singleton set. This property of deterministic actions facilitates the definition of finite computation runs to a considerable extent: we simply define it to be the unique finite computation sequence for which execution terminates.

Recall the definition of  $\text{Acseq}$  from Definition 2.1.

7.1. DEFINITION. The function  $\text{CS}$ , yielding the *finite computation sequences* of a given action, is inductively defined as follows.

$$\begin{aligned} \text{CS} & : \text{Ac} \rightarrow \wp(\text{Acseq}) \\ \text{CS}(\alpha) & = \{\alpha\} \text{ if } \alpha \text{ is semi-atomic} \end{aligned}$$

$$\begin{aligned}
\text{CS}(\alpha_1; \alpha_2) &= \{\alpha'_1; \alpha'_2 \mid \alpha'_1 \in \text{CS}(\alpha_1), \alpha'_2 \in \text{CS}(\alpha_2)\} \\
\text{CS}(\text{if } \varphi \text{ then } \alpha_1 \\
&\quad \text{else } \alpha_2 \text{ fi}) &= \text{CS}(\text{confirm } \varphi; \alpha_1) \cup \text{CS}(\text{confirm } \neg\varphi; \alpha_2) \\
\text{CS}(\text{while } \varphi \text{ do } \alpha \text{ od}) &= \bigcup_{k=1}^{\infty} \text{Seq}_k(\text{while } \varphi \text{ do } \alpha \text{ od}) \cup \{\text{confirm } \neg\varphi\} \\
\text{where for } k \geq 1 \\
\text{Seq}_k(\text{while } \varphi \text{ do } \alpha \text{ od}) &= \{(\text{confirm } \varphi; \alpha'_1); \dots; (\text{confirm } \varphi; \alpha'_k); \\
&\quad \text{confirm } \neg\varphi \mid \alpha'_j \in \text{CS}(\alpha) \text{ for } j = 1, \dots, k\}
\end{aligned}$$

7.2. DEFINITION. Since  $\text{Ac}$  is closed under the core clauses of Definition 2.1 only, the function  $\text{CS} : \text{Ac} \rightarrow \wp(\text{Acseq})$  is defined as usual. For  $M \in \mathbf{M}^C$  the function  $\text{CR}_M^C : A \times \text{Ac} \times S \rightarrow \wp(\text{Acseq})$  is defined by:

$$\text{CR}_M^C(i, \alpha, s) = \{\alpha' \in \text{CS}(\alpha) \mid \mathbf{r}^C(i, \alpha')(M, s) \neq \emptyset\}$$

Note that the actions that we consider are deterministic, i.e. the set  $\text{CR}_M^C(i, \alpha, s)$  consists of at most one element, for any  $i, \alpha, s$ . This must be kept in mind below, when we consider properties of computation runs. Moreover, since actions are deterministic we will often write (sloppily)  $\text{CR}_M^C(i, \alpha, s) = \alpha'$  instead of  $\text{CR}_M^C(i, \alpha, s) = \{\alpha'\}$ .

Let, for sets  $A_1, A_2 \in \wp(\text{Acseq})$ ,  $A_1; A_2$  stand for the set  $\{\alpha_1; \alpha_2 \mid \alpha_1 \in A_1, \alpha_2 \in A_2\}$ . Then we can state the following.

7.3. PROPOSITION. *If  $\mathbf{r}^C(i, \alpha_1; \alpha_2)(M, s) \neq \emptyset$  then  $\text{CR}_M^C(i, \alpha_1; \alpha_2, s) = \text{CR}_M^C(i, \alpha_1, s); \text{CR}_M^C(i, \alpha_2, s')$  for (the unique)  $s' \in \mathbf{r}^C(i, \alpha_1)(M, s)$ .*

Next we introduce the notion of a transition relation in the spirit of Structural Operational Semantics of Plotkin ([27]), which is a very neat and elegant way to describe computation by means of (single) transition steps. This method is widely used in computer science and we can use it here fruitfully to describe what happens with the agent's agenda when it performs actions step by step (from that agenda, so to speak).

In our set-up we consider transitions of the form  $\langle \alpha, s \rangle \xrightarrow{i, a}^M \langle \alpha', s' \rangle$ , where  $M \in \mathbf{M}^C$ ,  $\alpha, \alpha' \in \text{Ac}$ ,  $i \in A$ ,  $a \in \text{At}^+$  and  $s, s' \in S$ . (To ease notation in the sequel we will drop the superscript  $M$  if the set of states  $S$  in the model  $M$  is understood.) We use the symbol  $\Lambda$  for the empty action, with as property that  $\Lambda; \alpha = \alpha$ ;  $\Lambda = \alpha$ . Furthermore, we use the projection function  $\pi_2$ , which is assumed to yield the second element of a pair.

Transitions are given by the following deductive system, often called a transition system:

7.4. DEFINITION. Let the model  $M \in \mathbf{M}^C$  be given. The transition system  $T_M$  is given by the following axioms:

- $\langle \alpha, s \rangle \xrightarrow{i, \alpha} \langle \Lambda, s' \rangle$  with  $s' = \pi_2(\mathbf{r}^C(i, \alpha)(M, s))$  if  $\alpha$  is semi-atomic.
- $\langle \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi}, s \rangle \xrightarrow{i, \text{confirm } \varphi} \langle \alpha_1, s \rangle$  if  $s \models \varphi$
- $\langle \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi}, s \rangle \xrightarrow{i, \text{confirm } \neg\varphi} \langle \alpha_2, s \rangle$  if  $s \not\models \varphi$
- $\langle \text{while } \varphi \text{ do } \alpha \text{ od}, s \rangle \xrightarrow{i, \text{confirm } \varphi} \langle \alpha; \text{while } \varphi \text{ do } \alpha \text{ od}, s \rangle$  if  $s \models \varphi$
- $\langle \text{while } \varphi \text{ do } \alpha \text{ od}, s \rangle \xrightarrow{i, \text{confirm } \neg\varphi} \langle \Lambda, s \rangle$  if  $s \not\models \varphi$

and the following rule:



$$\bullet \frac{\langle \alpha_1, s \rangle \rightarrow_{i,a} \langle \alpha'_1, s' \rangle}{\langle \alpha_1; \alpha_2, s \rangle \rightarrow_{i,a} \langle \alpha'_1; \alpha_2, s' \rangle}$$

Obviously, there is a relation between transitions and the computation runs we introduced earlier. This relation is given by the following proposition. (Here we use as a convention that  $\text{CR}_M^C(i, \Lambda, s)$  is the empty sequence of actions.)

7.5. PROPOSITION.  $\text{CR}_M^C(i, \alpha, s) = a; \text{CR}_M^C(i, \alpha', s')$  iff  $T_M \vdash \langle \alpha, s \rangle \rightarrow_{i,a} \langle \beta, s' \rangle$  for some  $\beta$  with  $\text{CR}_M^C(i, \beta, s') = \text{CR}_M^C(i, \alpha', s')$ .

This proposition gives us as a corollary how computation runs can be viewed as being generated by transitions.

7.6. COROLLARY.  $\text{CR}_M^C(i, \alpha, s) = \{a_1; a_2; \dots; a_n\}$  iff  $T_M \vdash \langle \alpha, s \rangle \rightarrow_{i,a_1} \langle \alpha_1, s_1 \rangle \rightarrow_{i,a_2} \langle \alpha_2, s_2 \rangle \rightarrow_{i,a_3} \dots \rightarrow_{i,a_n} \langle \alpha_n, s_n \rangle$  for some  $\alpha_1, \dots, \alpha_n \in \text{Ac}$ ,  $s_1, \dots, s_n \in \text{S}$ , such that  $\alpha_n = \Lambda$ .

For reasons of convenience we introduce, analogously to the Can-predicate, a so-called Intend-predicate, which is meant to formalise the intentions of agents. The definition of this predicate is based on the idea that agents (loosely) intend to do all the actions that are correct and feasible with respect to some of their goals. As such, intention provides the precondition for successful commitment<sup>7</sup>.

7.7. DEFINITION. For  $\alpha \in \text{Ac}^C$ ,  $i \in \text{A}$  and  $\varphi \in \text{L}$  we define:

$$\mathbf{Intend}_i(\alpha, \varphi) =^{\text{def}} \mathbf{Can}_i(\alpha, \varphi) \wedge \mathbf{K}_i \mathbf{Goal}_i \varphi$$

Having established the formal prerequisites, we can now present the definitions formalising the intuitive description of the act of committing as presented above.

7.8. DEFINITION. For all  $M \in \mathbf{M}^C$  with state  $s$ , for all  $i \in \text{A}$  and  $\alpha \in \text{Ac}$  we define:

$$\begin{aligned} \mathbf{r}^C(i, \text{commit\_to } \alpha)(M, s) &= \emptyset \text{ if } M, s \models^C \neg \mathbf{Intend}_i(\alpha, \varphi) \text{ for all } \varphi \in C(i, s) \\ \mathbf{r}^C(i, \text{commit\_to } \alpha)(M, s) &= M', s \text{ with } M' = \langle \text{S}, \pi, \text{R}, \mathbf{r}_0, \mathbf{c}_0, \text{W}, \text{C}, \text{Agenda}' \rangle \\ &\text{where Agenda}' \text{ is minimal such that it is closed under the following conditions:} \\ &\text{for all } s' \in [s]_{\text{R}(i)}, \text{Agenda}'(i, s') = \text{Agenda}(i, s') \cup \{\alpha\} \\ &\text{and for all } s', s'', s''' \in \text{S}, \alpha' \in \text{Agenda}'(i, s') \text{ such that, for some semi-atomic } a, \\ &T_M \vdash \langle \alpha', s' \rangle \rightarrow_{i,a} \langle \alpha'', s'' \rangle \text{ and } s''' \in [s'']_{\text{R}(i)}: \\ &\text{Agenda}'(i, s''') = \text{Agenda}(i, s''') \cup \{\alpha''\} \\ &\text{otherwise} \end{aligned}$$

$$\mathbf{c}^C(i, \text{commit\_to } \alpha)(M, s) = \mathbf{1} \text{ iff } \text{Agenda}(i, s) = \emptyset$$

This definition makes sure that in the state where the commitment to an action  $\alpha$  is made, the agent's agenda is updated with  $\alpha$ , as well as in any state that is epistemically equivalent with  $s$ , i.e. in the epistemic equivalence class of  $s$ . The latter will have as an effect that the commitment to  $\alpha$  is known to the agent, since every epistemic alternative to  $s$  will contain the same commitment in the agenda. Moreover, the rest of the definition ensures that the agenda is also updated in all states reachable from  $s$  (or elements from its equivalence class)

<sup>7</sup>Our paraphrase of Cohen & Levesque's motto 'intention is choice plus commitment' [5] could therefore be stated as 'commitments are chosen intentions'.

by performing actions  $\beta$  that are according to the fulfilment of executing the action  $\alpha$  by putting in the agenda the ‘remainder’ of  $\alpha$  after its ‘partial execution’  $\beta$ .

Note, by the way, that the agenda function is well-defined due to the restrictions we’ve put on the models in Definition 4.1: it cannot happen that the same state gets assigned two different values of the same agent’s agenda. Without the restriction mentioned there in the fourth bullet the successful performance of atomic actions might result in the same state, so that this state would get two values of the agent’s agenda (viz. that associated with the situation before the execution of the action, and one associated with that after the execution), which contradicts the agenda being a function. In Definition 4.4, an analogous problem is circumvented for confirm actions, so that this problem does not occur for any semi-atomic action. Without the other restriction mentioned in Definition 4.1 it might happen that a (or any) state in the epistemic equivalence class would get two values of the agenda, since it might then happen that there are two states in the equivalence class both resulting from performing actions but obtaining different agenda’s, and these values should both hold for the whole equivalence class, by the above definition, contradicting the ‘functionality’ of the agenda again.

A last but important point to notice is that only the agenda of  $i$  is modified, and that only in those states that are somehow, i.e. by a combination of state-transitions and epistemic accessibility relations, connected to the state in which the commitment is being made. All other elements of the model remain unchanged.

The latter aspect mentioned above, i.e. the minimality of the change caused by performing a commitment, is partly formalised in Proposition 7.10 given below. Proposition 7.9 states the correctness of the definition of  $\mathbf{r}^C$  as presented above in the sense that it yields a (unique) well-defined model when applied to a well-defined model.

**7.9. PROPOSITION.** *For all  $M \in \mathbf{M}^C$  with state  $s$ , for all  $i \in A$  and  $\alpha \in Ac$ , if  $M', s = \mathbf{r}^C(i, \text{commit\_to } \alpha)(M, s)$  then  $M' \in \mathbf{M}^C$ .*

**7.10. PROPOSITION.** *For all  $M \in \mathbf{M}^C$  with state  $s$ , for all  $i \in A$  and  $\alpha \in Ac$ , if  $M', s = \mathbf{r}^C(i, \text{commit\_to } \alpha)(M, s)$  then for all states  $s'$  in  $M$ ,  $M, s' \models^C \varphi$  iff  $M', s' \models^C \varphi$ , for all  $\varphi \in L$ .*

Additional properties related to the `commit` actions are given in 7.4.

## 7.2. Being committed

After the rather elaborate and fairly complicated definition formalising the act of committing, defining what it means to be committed is a relatively straightforward and easy job. Basically, agents are committed to all actions whose semantic essence is captured by the normalised basic action in the appropriate agenda. The only additional aspect that has to be taken into account when defining the semantics of the **Committed** <sub>$i$</sub>  operator is that agents should start at the very beginning (a very good place to start), i.e. whenever an agent’s agenda contains a normalised basic action which is not semi-atomic, then the agent is also committed to actions whose semantic essence is a prefix of this normalised basic action. This constraint is quite an obvious one: how can agents be faithfully committed to a sequentially composed action if not committed to its first constituent? Formally we ensure this behaviour by using the prefix relation on basic actions. The notation  $\text{Prefix}(\alpha, \beta)$  for  $\alpha, \beta \in Ac_{\text{seq}}$  expresses that action sequence  $\alpha$  is a prefix of action sequence  $\beta$ . The definition of  $\models^C$  for the **Committed** <sub>$i$</sub>

operator could then be informally interpreted as ‘an agent is committed to those actions of which the semantic essence is a prefix of one of the actions in its agenda’.

7.11. DEFINITION. The binary relation  $\models^C$  between a formula in  $L^C$  and a pair  $M, s$  consisting of a model  $M$  for  $L^C$  and a state  $s$  in  $M$  is for commitments defined by:

$$\begin{aligned} M, s \models^C \mathbf{Committed}_i \alpha &\Leftrightarrow \\ \forall s' \in [s]_{R(i)} \exists \alpha_1 \in CR_M^C(i, \alpha, s') \exists \alpha_2 \in \text{Agenda}(i, s') \exists \alpha'_2 \in CR_M^C(i, \alpha_2, s') (\text{Prefix}(\alpha_1, \alpha'_2)) \end{aligned}$$

An investigation of the properties of the commitment operator is postponed to 7.4.

### 7.3. Getting uncommitted

By performing an uncommit action, agents may undo previously made commitments that turned out to be either useless or impossible. That is, as soon as an agent no longer knows some commitment to be correct and feasible for at least one of its goals it may undo this commitment. Just as we did for the commit action, we have to decide upon the constituents of the result, opportunity and ability for the actions formalising the act of uncommitting. The result of such an action is obvious: agents should no longer be committed to  $\alpha$  after a successful performance of an `uncommit  $\alpha$`  action<sup>8</sup>. Defining what it means to have the opportunity and ability to uncommit represents a somewhat more arbitrary choice. We have decided to let an agent have the opportunity to undo any of its commitments, i.e. there is nothing in its circumstances that may prevent an agent to undo a commitment. Our loyal, diligent agents are however only (morally) capable of undoing commitments that have become redundant. The actual definition of the functions  $\mathbf{r}^C$  and  $\mathbf{c}^C$  consists of nothing but a formalisation of these intuitive ideas.

7.12. DEFINITION. For all  $M \in \mathbf{M}^C$  with state  $s$ , for all  $i \in A$  and  $\alpha \in Ac$  we define:  
(For technical convenience, in the following we define, simultaneously with the new agenda, the set  $\text{Reachable}_{M,s}$  of states that are reachable from  $s$  by alternatively consider epistemic alternatives and performing actions from the agenda.)

$$\begin{aligned} \mathbf{r}^C(i, \text{uncommit } \alpha)(M, s) &= \emptyset \text{ if } M, s \models^C \neg \mathbf{Committed}_i \alpha \\ \mathbf{r}^C(i, \text{uncommit } \alpha)(M, s) &= M', s \text{ with } M' = \langle S, \pi, R, \mathbf{r}_0, \mathbf{c}_0, W, C, \text{Agenda}' \rangle \\ &\text{where } \text{Agenda}' \text{ and } \text{Reachable}_{M,s} \text{ are minimal such that they are closed under the conditions:} \\ &\text{for all } s' \in [s]_{R(i)}: \\ &\quad s' \in \text{Reachable}_{M,s} \text{ and} \\ &\quad \text{Agenda}'(i, s') = \text{Agenda}(i, s') \setminus \{ \beta \mid \text{Prefix}(CR_M^C(i, \alpha, s'), CR_M^C(i, \beta, s')) \} \\ &\text{and for all } s' \in \text{Reachable}_{M,s}, s'', s''' \in S \text{ with } \alpha' \in \text{Agenda}'(i, s') \text{ and such that,} \\ &\text{for some semi-atomic } a, T_M \vdash \langle \alpha', s' \rangle \xrightarrow{i,a} \langle \alpha'', s'' \rangle \text{ and } s''' \in [s'']_{R(i)}: \\ &\quad s''' \in \text{Reachable}_{M,s} \text{ and} \\ &\quad \text{Agenda}'(i, s''') = \text{Agenda}(i, s''') \setminus \{ \beta \mid \text{Prefix}(CR_M^C(i, \alpha'', s'''), CR_M^C(i, \beta, s''')) \} \\ &\text{otherwise} \end{aligned}$$

---

<sup>8</sup>As was pointed out to us by John Fox, this description of the result of undoing a commitment comprises a major simplification. For in real life, undoing commitments involves more than just abandoning future commitments: it is also necessary to (try to) undo all the effects that followed from initially pursuing the commitment. For example, if an agent that is committed to  $\alpha_1; \alpha_2$  finds out after having done  $\alpha_1$  that its commitment to  $\alpha_2$  should be undone, then it is very plausible that it should not only remove  $\alpha_2$  from its agenda but also try to undo as many of the effects of  $\alpha_1$  as possible.

$$c^C(i, \text{uncommit } \alpha)(M, s) = 1 \text{ iff } M, s \models^C \neg \mathbf{Intend}_i(\alpha, \varphi) \text{ for all } \varphi \in C(i, s)$$

Note that the slight complication in this definition regarding prefixes: due to our (intuitively correct) definition of being committed to an action if it constitutes the partial execution of some action that is recorded in the agent's agenda, we must also take care that when we uncommit to some action  $\alpha$  that we remove from the agent's agenda all actions  $\beta$  that have  $\alpha$  as their partial execution. Formally this means that we remove all  $\beta$  that have the computation run of  $\alpha$  as a prefix of their computation run. This will ensure the desirable property below (Proposition 7.15, fifth item) stating that when the agent uncommits to an action to which it was committed, it will indeed be not committed to this action afterwards.

Our definition of  $r^C$  for the `uncommit` actions is also twofold correct: not only does performing an `uncommit` action provide for a correct model-transformation, but also does it do so while causing minimal change.

7.13. PROPOSITION. *For all  $M \in \mathbf{M}^C$  with state  $s$ , for all  $i \in A$  and  $\alpha \in Ac$ , if  $M', s = r^C(i, \text{uncommit } \alpha)(M, s)$  then  $M' \in \mathbf{M}_{\sim}^C$ .*

7.14. PROPOSITION. *For all  $M \in \mathbf{M}^C$  with state  $s$ , for all  $i \in A$  and  $\alpha \in Ac$ , if  $M', s = r^C(i, \text{uncommit } \alpha)(M, s)$  then for all states  $s'$  in  $M$ ,  $M, s' \models^C \varphi$  iff  $M', s' \models^C \varphi$ , for all  $\varphi \in L$ .*

Additional validities characterising the `uncommit` action are given below.

#### 7.4. The statics and dynamics of commitments

Here we characterise the statics and dynamics of commitments by presenting some validities for  $\models^C$ . For a start we consider a number of validities characterising the dynamics of commitments.

7.15. PROPOSITION. *For all  $i \in A$ ,  $\alpha, \beta \in Ac$  and  $\varphi \in L$  we have:*

1.  $\models^C \mathbf{Intend}_i(\alpha, \varphi) \rightarrow \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \top$
2.  $\models^C \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \top \leftrightarrow \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \mathbf{Committed}_i \alpha$
3.  $\models^C \mathbf{Committed}_i \alpha \rightarrow \neg \mathbf{A}_i \text{commit\_to } \beta$
4.  $\models^C [\text{do}_i(\text{commit\_to } \alpha)] \neg \mathbf{A}_i \text{commit\_to } \beta$
5.  $\models^C \mathbf{Committed}_i \alpha \leftrightarrow \langle \text{do}_i(\text{uncommit } \alpha) \rangle \neg \mathbf{Committed}_i \alpha$
6.  $\models^C \mathbf{Intend}_i(\alpha, \varphi) \rightarrow \neg \mathbf{A}_i \text{uncommit } \alpha$
7.  $\models^C (\mathbf{C}_i \varphi \leftrightarrow \mathbf{K}_i \mathbf{C}_i \varphi) \rightarrow (\mathbf{A}_i \text{uncommit } \alpha \leftrightarrow \mathbf{K}_i \mathbf{A}_i \text{uncommit } \alpha)$
8.  $\models^C \mathbf{Committed}_i \alpha \wedge \neg \mathbf{Can}_i(\alpha, \top) \rightarrow \mathbf{Can}_i(\text{uncommit } \alpha, \neg \mathbf{Committed}_i \alpha)$

The first two items of Proposition 7.15 jointly formalise our version of the syllogism of practical reasoning as described above. In the third item it is stated that being committed prevents an agent from having the ability to (re)commit. The fourth item states that the act of committing is ability-destructive with respect to future `commit` actions, i.e. by performing a commitment an agent loses its ability to make any other commitments. Item 5 states that being committed is a necessary and sufficient condition for having the opportunity to uncommit; as mentioned above, agents have the opportunity to undo all of their commitments. In item 6 it is stated that agents are (morally) unable to undo commitments to actions that are still known to be correct and feasible to achieve some goal. In item 7 it is formalised that

agents know of their abilities to uncommit to some action. The last item states that whenever an agent is committed to an action that is no longer known to be practically possible, it knows that it can undo this impossible commitment.

The following proposition formalises some of the desiderata for the statics of commitments that turn out to be valid in the class  $\mathbf{M}^C$  of models for  $L^C$ .

7.16. PROPOSITION. *For all  $i \in A$ ,  $\alpha, \alpha_1, \alpha_2 \in Ac$  and all  $\varphi \in L$  we have:*

1.  $\models^C \mathbf{Committed}_i \alpha \rightarrow \mathbf{K}_i \mathbf{Committed}_i \alpha$
2.  $\models^C \mathbf{Committed}_i(\alpha_1; \alpha_2) \rightarrow \mathbf{Committed}_i \alpha_1 \wedge \mathbf{K}_i[\mathbf{do}_i(\alpha_1)] \mathbf{Committed}_i \alpha_2$
3.  $\models^C \mathbf{Committed}_i \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \wedge \mathbf{K}_i \varphi \rightarrow \mathbf{Committed}_i(\mathbf{confirm} \varphi; \alpha_1)$
4.  $\models^C \mathbf{Committed}_i \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \wedge \mathbf{K}_i \neg \varphi \rightarrow \mathbf{Committed}_i(\mathbf{confirm} \neg \varphi; \alpha_2)$
5.  $\models^C \mathbf{Committed}_i \text{while } \varphi \text{ do } \alpha \text{ od} \wedge \mathbf{K}_i \varphi \rightarrow \mathbf{Committed}_i((\mathbf{confirm} \varphi; \alpha); \text{while } \varphi \text{ do } \alpha \text{ od})$

The first item of Proposition 7.16 states that commitments are known. The second item states that a commitment to a sequential composition  $\alpha_1; \alpha_2$  of actions implies a commitment to the initial part  $\alpha_1$ , and that the agent knows that after execution of this initial part  $\alpha_1$  it will be committed to the remainder  $\alpha_2$ . The third and fourth item formalise the rationality of agents with regard to their commitments to conditionally composed actions. The last item concerns the unfolding of a while-loop: if an agent is committed to a while-loop while knowing the condition of the loop to be true, then the agent is also committed to the then-part of the while-loop.

## 8. Summary and conclusions

In this paper we presented a formalisation of motivational attitudes, the attitudes that explain why agents act the way they do. This formalisation concerns operators both on the assertion level, where operators range over propositions, and on the practition level, where operators range over actions. An important feature of our formalisation is the attention paid to the acts associated with selecting between wishes and with (un)committing to actions. Starting from the primitive notion of wishes, we defined goals to be selected, unfulfilled, implementable wishes. Commitments may be made to actions that are known to be correct and feasible with respect to some goal and may be undone whenever the action to which an agent has committed itself has either become impossible or useless. Both the act of making, and the act of undoing commitments are formalised as model-transforming actions in our framework. The actions that an agent is committed to are recorded in its agenda in such a way that commitments are closed under prefix-taking and under practical identity. On the whole our formalisation is a rather expressive one, which tries to be faithful to a certain extent to both commonsense intuition and philosophical insights.

### 8.1. Future work

The major extension to the framework presented here concerns a formalisation of the actual execution of actions. Although the conditional nature of a framework based on dynamic logic makes it perhaps less suitable for an adequate formalisation of ‘doing’, one could think of a practition operator indicating which action is actually performed next. Using this predicate would enhance expressiveness in that it would be possible to formulate relations between actions that agents are committed to, and actions that they actually perform. Another way

to extend the framework would be by establishing further relations with deontic notions like obligations and violations. A combination of the ‘doing’-predicate with a deontic notion modelling violations or penalties would then allow one to model that agents should execute the actions that they are committed to if they want to avoid penalties. Research along these lines was initiated by Dignum & Van Linder [6, 7]. Furthermore, we will investigate how the model as presented may serve as a rigorous semantic framework for agent-oriented languages such as AGENT0 [36] and PLACA [39], so that, based on this model, specification methods for programs written in these languages can be obtained.

## 8.2. Related work

The formalisation of motivational attitudes has received much attention within the agent research community. Probably the most influential account of motivational attitudes is due to Cohen & Levesque [5]. Starting from the primitive notions of implicit goals and beliefs, Cohen & Levesque define so-called persistent goals, which are goals which agents give up only when they think they are either satisfied or will never be true, and intentions, both ranging over propositions and over actions. The idea underlying persistent goals is similar to that underlying our notion of goals. In the framework of Cohen & Levesque agents intend to bring about a proposition if they intend to do some action that brings about the proposition. An agent intends to do an action if it has the persistent goal to have done the action. This reduction of intentions to do actions for goals is a rather artificial and philosophically very questionable one: although intentions to actions should be related to goals, this relation should express that doing the action helps in bringing about some goal and not that doing the action in itself is a goal. Furthermore the coexistence of goals and intentions ranging over propositions seems to complicate matters unnecessarily.

Another important formalisation of motivational attitudes is proposed by Rao & Georgeff [30] in their BDI-architecture. Treating desires and intentions as primitive, Rao & Georgeff focus on the process of intention revision rather than the ‘commitment acquisition’ which is essential to our formalisation. Both desires and intentions in their framework suffer from the problems associated with logical omniscience. To avoid these problems, Cavedon *et al.* [4] propose the use of non-normal logics of intention and belief in the BDI-architecture, and more in particular Rantala’s ‘impossible worlds’ framework [28]. This ‘impossible worlds’ approach was originally proposed as a way to solve the problems of logical omniscience for informational attitudes. Hence, whereas we more or less employ the awareness approach, Cavedon *et al.* propose yet another technique developed to solve the problems of logical omniscience. It therefore may come as no surprise that the properties that Cavedon *et al.* acquire for intentions are highly similar to the properties of goals given in Section 6.

We also like to mention here the work by Singh [37], which bears some resemblance to our calculus to ‘maintain the agent’s agenda’. The main difference is that his theory rests on temporal branching-time logic. (He also has some constructs *à la* dynamic logic in his language, but these, too, have an explicit temporal interpretation.) Furthermore, the emphasis in his work is put on exploring the intricacies of intentional notions as related to issues of nondeterminism, which we have not considered in this paper at all.

The last formalisation of motivational attitudes that we would like to mention is the one proposed by Dignum *et al.* [8]. In this formalisation, which is inspired by and based on research on deontic logic as carried out by Dignum *et al.*, notions like decisions, intentions and commitments are modelled. Of these, decisions and the act of committing are interpreted as

so-called meta-actions, a notion similar to that of model-transformers. Despite its complexity, which is due to the incorporation of an algebraic semantics of actions and a trace semantics to model histories, some of the essential ideas underlying the formalisation of Dignum *et al.* are not unlike those underlying the formalisation presented here.

**Acknowledgements.** We like to thank Wieke de Vries for reading an earlier version of this paper and providing us with some very useful comments. We are also grateful to the attendants of ATAL'95 and Modelage'96 for their helpful response to this work. Furthermore, the partial support of ESPRIT Working Group MODELAGE (BRWG 8319) is gratefully acknowledged.

## A. Selected proofs

5.3. PROPOSITION. *All of the properties of logical omniscience formalised in Definition 5.2, with the exception of LO7, are valid for the  $\mathbf{W}_i$  operator.*

PROOF: Properties LO1 and LO2 state that  $\mathbf{W}_i$  is a normal modal operator and are shown as for any necessity operator. Property LO3 follows directly by combining LO1 and LO2, and LO4 is a direct consequence of LO3. Properties LO5 and LO6 are typical for necessity operators in a normal modal logic: for whenever both  $\varphi$  and  $\psi$  hold at a set of designated worlds,  $\varphi \wedge \psi$  also holds at all the worlds from that set (LO5), and if  $\varphi$  holds at all worlds from some set then  $\varphi \vee \psi$  does also (LO6). That LO7 is not valid for the  $\mathbf{W}_i$  operator is seen by considering a model M with state  $s$  such that no state  $s'$  exists with  $(s, s') \in W(i)$ . Then it holds that  $M, s \models^C \mathbf{W}_i\varphi \wedge \mathbf{W}_i\neg\varphi$ , for all  $\varphi \in L$ .

⊠

6.7. PROPOSITION. *None of the properties of logical omniscience formalised in Definition 5.2, with the exception of LO7, is valid for the  $\mathbf{Goal}_i$  operator.*

PROOF: Properties LO1, LO3, LO4, LO5 and LO6 are most easily seen not to hold for the goal operator by noting the absence of any closure properties on the set  $C(i, s)$ , for  $i \in A$  and  $s$  some state. Due to this absence it is perfectly possible that  $\varphi$  and  $\varphi \rightarrow \psi$  are both in  $C(i, s)$  while  $\psi$  is not (LO1), that  $\varphi \in C(i, s)$  and  $\psi \notin C(i, s)$  while  $\models^C \varphi \rightarrow \psi$  (LO3) or  $\models^C \varphi \leftrightarrow \psi$  (LO4), that  $\{\varphi, \psi\} \subseteq C(i, s)$  and  $\varphi \wedge \psi \notin C(i, s)$  (LO5), or that  $\varphi \in C(i, s)$  while  $\varphi \vee \psi \notin C(i, s)$  (LO6), for appropriate  $i \in A$  and  $s$  a state in some model. Property LO2 is seen not to hold by observing that  $\models^C \varphi$  implies that  $\varphi$  is fulfilled always and everywhere, which means that  $\varphi$  is not a goal. In fact, one can show that whenever  $\varphi$  is inevitable, i.e.  $\models^C \varphi$  holds, it is necessarily not a goal, i.e.  $\models^C \neg\mathbf{Goal}_i\varphi$  holds (cf. item 5 of Proposition 6.8). That LO7 holds for goals is a direct consequence of their unfulfilledness. For in any possible state  $s$  of any possible model M, either  $\varphi$  holds and thereby  $M, s \not\models^C \mathbf{Goal}_i\varphi$ , or  $\neg\varphi$  holds and thereby  $M, s \not\models^C \mathbf{Goal}_i\neg\varphi$ . Hence LO7 is a valid property for goals.

⊠

6.8. PROPOSITION. *For all  $i \in A$  and  $\varphi \in L$  we have:*

1.  $\models^C \mathbf{W}_i\varphi \leftrightarrow \langle \text{do}_i(\text{select } \varphi) \rangle \top$
2.  $\models^C \langle \text{do}_i(\text{select } \varphi) \rangle \top \leftrightarrow \langle \text{do}_i(\text{select } \varphi) \rangle \mathbf{C}_i\varphi$

3.  $\models^C \neg \mathbf{A}_i \text{select } \varphi \rightarrow [\text{do}_i(\text{select } \varphi)] \neg \mathbf{Goal}_i \varphi$
4.  $\models^C \mathbf{PracPoss}_i(\text{select } \varphi, \top) \leftrightarrow \langle \text{do}_i(\text{select } \varphi) \rangle \mathbf{Goal}_i \varphi$
5.  $\models^C \varphi \Rightarrow \models^C \neg \mathbf{Goal}_i \varphi$
6.  $(\varphi \rightarrow \psi) \rightarrow (\mathbf{Goal}_i \varphi \rightarrow \mathbf{Goal}_i \psi)$  is not for all  $\varphi, \psi \in \mathbf{L}$  valid
7.  $\mathbf{K}_i(\varphi \rightarrow \psi) \rightarrow (\mathbf{Goal}_i \varphi \rightarrow \mathbf{Goal}_i \psi)$  is not for all  $\varphi, \psi \in \mathbf{L}$  valid

PROOF: We successively show all items. Let  $M \in \mathbf{M}^C$  with state  $s$  and  $\varphi \in \mathbf{L}$  be arbitrary.

1. An easy inspection of Definition 6.2 shows that  $\mathbf{r}^C(i, \text{select } \varphi)(M, s) = \emptyset$  iff  $M, s \not\models^C \mathbf{W}_i \varphi$ . Thus  $M, s \models^C \mathbf{W}_i \varphi \leftrightarrow \langle \text{do}_i(\text{select } \varphi) \rangle \top$ , which was to be shown.
2. If  $M', s = \mathbf{r}^C(i, \text{select } \varphi)(M, s)$ , then  $M'$  is such that  $\mathbf{C}'(i, s)$  contains  $\varphi$ . Then by definition  $M', s \models^C \mathbf{C}_i \varphi$ , and thus  $M, s \models^C \langle \text{do}_i(\text{select } \varphi) \rangle \mathbf{C}_i \varphi$  if  $M, s \models^C \langle \text{do}_i(\text{select } \varphi) \rangle \top$ , which suffices to conclude item 2.
3. Suppose  $M, s \models^C \neg \mathbf{A}_i \text{select } \varphi$ , i.e.  $M, s \models^C \varphi \vee \neg \diamond_i \varphi$ . Now by definition,  $\varphi \in \mathbf{L}$ , and hence, by Proposition 6.4,  $M', s \models^C \varphi$  if  $M, s \models^C \varphi$  whenever  $M', s = \mathbf{r}^C(i, \text{select } \varphi)(M, s)$ . By Corollary 6.5 it follows that for  $M'$  as aforementioned holds that  $M', s \models^C \neg \diamond_i \varphi$  if  $M, s \models^C \neg \diamond_i \varphi$ . Thus if  $M, s \models^C \varphi \vee \neg \diamond_i \varphi$  then it holds for  $M', s = \mathbf{r}^C(i, \text{select } \varphi)(M, s)$  that  $M', s \models^C \varphi \vee \neg \diamond_i \varphi$ . By definition it then directly follows that  $M', s \models^C \neg \mathbf{Goal}_i \varphi$ , and thus  $M, s \models^C \neg \mathbf{A}_i \text{select } \varphi \rightarrow [\text{do}_i(\text{select } \varphi)] \neg \mathbf{Goal}_i \varphi$ , which was to be shown.
4. This item follows by combining item 2 of this proposition with Proposition 6.4 and Corollary 6.5.
5. If  $\models^C \varphi$  holds, then  $M, s \models^C \varphi$  for all  $M \in \mathbf{M}^C$  with state  $s$ . Hence  $M, s \models^C \neg \mathbf{Goal}_i \varphi$  for all  $M \in \mathbf{M}^C$  and their states  $s$ , and thus  $\models^C \neg \mathbf{Goal}_i \varphi$ .
6. This item is easily shown by selecting an appropriate contingency  $\varphi$  and an arbitrary tautology  $\psi$ , such that for certain  $M$  and  $s$  holds that  $M, s \models^C \mathbf{Goal}_i \varphi$ . For then  $M, s \models^C (\varphi \rightarrow \psi) \wedge \mathbf{Goal}_i \varphi$  while — by the previous item —  $M, s \not\models^C \mathbf{Goal}_i \psi$ .
7. Item 7 is proved similarly to item 6.

⊠

7.15. PROPOSITION. For all  $i \in \mathbf{A}$ ,  $\alpha, \beta \in \mathbf{Ac}$  and  $\varphi \in \mathbf{L}$  we have:

1.  $\models^C \mathbf{Intend}_i(\alpha, \varphi) \rightarrow \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \top$
2.  $\models^C \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \top \leftrightarrow \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \mathbf{Committed}_i \alpha$
3.  $\models^C \mathbf{Committed}_i \alpha \rightarrow \neg \mathbf{A}_i \text{commit\_to } \beta$
4.  $\models^C [\text{do}_i(\text{commit\_to } \alpha)] \neg \mathbf{A}_i \text{commit\_to } \beta$
5.  $\models^C \mathbf{Committed}_i \alpha \leftrightarrow \langle \text{do}_i(\text{uncommit } \alpha) \rangle \neg \mathbf{Committed}_i \alpha$
6.  $\models^C \mathbf{Intend}_i(\alpha, \varphi) \rightarrow \neg \mathbf{A}_i \text{uncommit } \alpha$
7.  $\models^C (\mathbf{C}_i \varphi \leftrightarrow \mathbf{K}_i \mathbf{C}_i \varphi) \rightarrow (\mathbf{A}_i \text{uncommit } \alpha \leftrightarrow \mathbf{K}_i \mathbf{A}_i \text{uncommit } \alpha)$
8.  $\models^C \mathbf{Committed}_i \alpha \wedge \neg \mathbf{Can}_i(\alpha, \top) \rightarrow \mathbf{Can}_i(\text{uncommit } \alpha, \neg \mathbf{Committed}_i \alpha)$

PROOF: We show the second, third, fourth, seventh and eight item; the other ones follow directly from the respective definitions. Let  $M \in \mathbf{M}^C$  with state  $s$ , and  $i \in \mathbf{A}$ ,  $\alpha, \beta \in \mathbf{Ac}$  be arbitrary.

2. Let  $M, s \models^C \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \top$  and let  $M', s = \mathbf{r}^C(i, \text{commit\_to } \alpha)(M, s)$ . We have to show that  $M', s \models^C \mathbf{Committed}_i \alpha$ , i.e. we have to show that  $\forall s' \in [s]_{R'(i)} \exists \alpha_1 \in \text{CR}_{M'}^C(i, \alpha, s') \exists \alpha_2 \in \text{Agenda}'(i, s') \exists \alpha'_2 \in \text{CR}_M^C(i, \alpha_2, s') (\text{Prefix}(\alpha_1, \alpha'_2))$ . An inspection of Definition 7.8 shows that for all  $s' \in [s]_{R(i)} = [s]_{R'(i)}$  holds that  $\text{Agenda}'(i, s')$  contains  $\alpha$ . This implies that  $M', s \models^C \mathbf{Committed}_i \alpha$ . Thus  $M, s \models^C \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \mathbf{Committed}_i \alpha$ , which suffices to conclude that item 2 holds.



3. If  $M, s \models^C \mathbf{Committed}_i \alpha$  then, by Definition 7.11, we have that  $\text{Agenda}(i, s) \neq \emptyset$ . Hence, by Definition 7.8,  $M, s \models^C \neg \mathbf{A}_i \text{commit\_to } \beta$ .
4. If  $\mathbf{r}^C(i, \text{commit\_to } \alpha)(M, s) = \emptyset$  then  $M, s \models^C [\text{do}_i(\text{commit\_to } \alpha)] \neg \mathbf{A}_i \text{commit\_to } \beta$  is trivially true. Else  $M, s \models^C \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \mathbf{Committed}_i \alpha$  by item 2 of this proposition, and, by item 3, this implies  $M, s \models^C \langle \text{do}_i(\text{commit\_to } \alpha) \rangle \neg \mathbf{A}_i \text{commit\_to } \beta$ , which suffices to conclude item 4.
7. Suppose  $M, s \models^C \mathbf{C}_i \varphi \leftrightarrow \mathbf{K}_i \mathbf{C}_i \varphi$  and  $M, s \models^C \mathbf{A}_i \text{uncommit } \alpha$ . This implies that  $M, s \models^C \neg \mathbf{Intend}_i(\alpha, \varphi)$ , for all  $\varphi \in C(i, s)$ . That is,  $M, s \models^C \neg \mathbf{Can}_i(\alpha, \varphi) \vee \neg \mathbf{K}_i \mathbf{Goal}_i \varphi$  for all  $\varphi \in C(i, s)$ . But by the introspective properties of knowledge the latter implies that  $M, s \models^C \mathbf{K}_i \neg \mathbf{Can}_i(\alpha, \varphi) \vee \mathbf{K}_i \neg \mathbf{K}_i \mathbf{Goal}_i \varphi$ , for all  $\varphi \in C(i, s)$ . Hence  $M, s \models^C \mathbf{K}_i(\neg \mathbf{Can}_i(\alpha, \varphi) \vee \neg \mathbf{K}_i \mathbf{Goal}_i \varphi)$ , for all  $\varphi \in C(i, s)$ , and thus for all  $s' \in [s]_{R(i)}$  it holds that  $M, s' \models^C \neg \mathbf{Intend}_i(\alpha, \varphi)$  for all  $\varphi \in C(i, s)$ , i.e. for all  $\varphi$  such that  $M, s \models^C \mathbf{C}_i \varphi$ . By the assumption  $M, s \models^C \mathbf{C}_i \varphi \leftrightarrow \mathbf{K}_i \mathbf{C}_i \varphi$ , we obtain that  $M, s' \models^C \neg \mathbf{Intend}_i(\alpha, \varphi)$  for all  $\varphi$  such that  $M, s \models^C \mathbf{K}_i \mathbf{C}_i \varphi$ , i.e. for all  $\varphi \in C(i, s')$  for all  $s' \in [s]_{R(i)}$ . Consequently,  $M, s' \models^C \mathbf{A}_i \text{uncommit } \alpha$ . Thus  $M, s \models^C \mathbf{K}_i \mathbf{A}_i \text{uncommit } \alpha$ , which suffices to conclude that item 7 indeed holds.
8. Suppose  $M, s \models^C \mathbf{Committed}_i \alpha \wedge \neg \mathbf{Can}_i(\alpha, \top)$ . Then  $M, s \models^C \neg \mathbf{Can}_i(\alpha, \varphi)$  for all  $\varphi \in L$ , and thus  $M, s \models^C \neg \mathbf{Intend}_i(\alpha, \varphi)$  for all  $\varphi \in C(i, s)$ . Then, by definition of  $\mathbf{c}^C$ ,  $M, s \models^C \mathbf{A}_i \text{uncommit } \alpha$ , and, by the previous item,  $M, s \models^C \mathbf{K}_i \mathbf{A}_i \text{uncommit } \alpha$ . Also,  $M, s \models^C \mathbf{Committed}_i \alpha$  implies  $M, s \models^C \mathbf{K}_i \mathbf{Committed}_i \alpha$  by Proposition 7.16(1), and, by item 5 of this proposition,  $M, s \models^C \mathbf{K}_i \langle \text{do}_i(\text{uncommit } \alpha) \rangle \neg \mathbf{Committed}_i \alpha$ . Thus  $M, s \models^C \mathbf{K}_i \langle \text{do}_i(\text{uncommit } \alpha) \rangle \neg \mathbf{Committed}_i \alpha \wedge \mathbf{K}_i \mathbf{A}_i \text{uncommit } \alpha$ . This implies that  $M, s \models^C \mathbf{Can}_i(\text{uncommit } \alpha, \neg \mathbf{Committed}_i \alpha)$ , which suffices to conclude item 8.

⊠

7.16. PROPOSITION. *For all  $i \in A$ ,  $\alpha, \alpha_1, \alpha_2 \in Ac$  and all  $\varphi \in L$  we have:*

1.  $\models^C \mathbf{Committed}_i \alpha \rightarrow \mathbf{K}_i \mathbf{Committed}_i \alpha$
2.  $\models^C \mathbf{Committed}_i(\alpha_1; \alpha_2) \rightarrow \mathbf{Committed}_i \alpha_1 \wedge \mathbf{K}_i[\text{do}_i(\alpha_1)] \mathbf{Committed}_i \alpha_2$
3.  $\models^C \mathbf{Committed}_i \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \wedge \mathbf{K}_i \varphi \rightarrow \mathbf{Committed}_i(\text{confirm } \varphi; \alpha_1)$
4.  $\models^C \mathbf{Committed}_i \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \wedge \mathbf{K}_i \neg \varphi \rightarrow \mathbf{Committed}_i(\text{confirm } \neg \varphi; \alpha_2)$
5.  $\models^C \mathbf{Committed}_i \text{while } \varphi \text{ do } \alpha \text{ od} \wedge \mathbf{K}_i \varphi \rightarrow \mathbf{Committed}_i((\text{confirm } \varphi; \alpha); \text{while } \varphi \text{ do } \alpha \text{ od})$

PROOF: We successively show all items. Let  $M \in \mathbf{M}^C$  with state  $s$  and  $\varphi \in L$ ,  $\alpha, \alpha_1, \alpha_2 \in Ac$  be arbitrary.

1. Assume that  $M, s \models^C \mathbf{Committed}_i \alpha$ . By Definition 7.11 it then follows that  $\forall s' \in [s]_{R(i)} \exists \alpha_1 \in \text{CR}_M^C(i, \alpha, s') \exists \alpha_2 \in \text{Agenda}(i, s') \exists \alpha'_2 \in \text{CR}_M^C(i, \alpha_2, s') (\text{Prefix}(\alpha_1, \alpha'_2))$ . Since  $[s]_{R(i)}$  is an equivalence class we have that  $\forall s'' \in [s]_{R(i)} \forall s' \in [s'']_{R(i)} \exists \alpha_1 \in \text{CR}_M^C(i, \alpha, s') \exists \alpha_2 \in \text{Agenda}(i, s') \exists \alpha'_2 \in \text{CR}_M^C(i, \alpha_2, s') (\text{Prefix}(\alpha_1, \alpha'_2))$ , which implies  $M, s'' \models^C \mathbf{Committed}_i \alpha$  for all  $s'' \in [s]_{R(i)}$ , and thus  $M, s \models^C \mathbf{K}_i \mathbf{Committed}_i \alpha$ .
2. Let  $M, s \models^C \mathbf{Committed}_i \alpha_1; \alpha_2$ , i.e. for all  $s' \in [s]_{R(i)}$  some  $\beta_1 \in \text{CR}_M^C(i, \alpha_1; \alpha_2, s')$ ,  $\beta_2 \in \text{Agenda}(i, s')$  and  $\beta'_2 \in \text{CR}_M^C(i, \beta_2, s')$  exist such that  $\text{Prefix}(\beta_1, \beta'_2)$  holds. But then also for all  $s' \in [s]_{R(i)}$  some  $\gamma_1 \in \text{CR}_M^C(i, \alpha_1, s')$ ,  $\gamma_2 \in \text{Agenda}(i, s')$  and  $\gamma'_2 \in \text{CR}_M^C(i, \gamma_2, s')$  exist such that  $\text{Prefix}(\gamma_1, \gamma'_2)$  holds (viz.  $\gamma_2 = \beta_2$  and  $\gamma'_2 = \beta'_2$  satisfy the requirement, since  $\text{Prefix}(\beta_1, \gamma_1)$ ). Thus  $M, s \models^C \mathbf{Committed}_i \alpha_1$ .

Furthermore, consider some state  $s' \in [s]_{R(i)}$  and some state  $s'' \in \mathbf{r}^C(i, \alpha_1)(M, s')$ . By the above we obtain that there exist some  $\beta_1 \in \text{CR}_M^C(i, \alpha_1; \alpha_2, s')$ ,  $\beta_2 \in \text{Agenda}(i, s')$  and

- $\beta_2' \in \text{CR}_M^C(i, \beta_2, s')$  exist such that  $\text{Prefix}(\beta_1, \beta_2')$  holds. By Proposition 7.3 we have that  $\text{CR}_M^C(i, \beta_1, s') = \text{CR}_M^C(i, \alpha_1, s'); \text{CR}_M^C(i, \alpha_2, s'')$ . Suppose  $\text{CR}_M^C(i, \alpha_1, s') = \{b_1; b_2; \dots; b_m\}$ . Since  $\text{CR}_M^C(i, \beta_2, s') = \text{CR}_M^C(i, \alpha_1; \alpha_2, s'); \gamma = \text{CR}_M^C(i, \alpha_1, s'); \text{CR}_M^C(i, \alpha_2, s''); \gamma = b_1; b_2; \dots; b_m; \text{CR}_M^C(i, \alpha_2, s''); \gamma$ , for some  $\gamma \in \text{Acseq}$ , we have by Proposition 7.5 that  $\langle s', \beta_2 \rangle \rightarrow_{b_1} \dots \rightarrow_{b_m} \langle s'', \delta_2 \rangle$ , with  $\delta_2$  such that  $\text{CR}_M^C(i, \delta_2, s'') = \text{CR}_M^C(i, \alpha_2, s''); \gamma$ . This implies that  $\delta_2 \in \text{Agenda}(i, s''')$  for all  $s''' \in [s'']_{R(i)}$ . Thus we have that, for all  $s''' \in [s'']_{R(i)}$ , there exist  $\delta_1 \in \text{CR}_M^C(i, \alpha_2, s''')$ ,  $\delta_2 \in \text{Agenda}(i, s''')$  and  $\delta_2' \in \text{CR}_M^C(i, \delta_2, s''')$  such that  $\text{Prefix}(\delta_1, \delta_2')$ .
3. Assume that  $M, s \models^C \mathbf{Committed}_i \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \wedge \mathbf{K}_i \varphi$ . By definition of  $\text{CR}_M^C$  and CS we have  $\text{CR}_M^C(i, \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi}, s') = \text{CR}_M^C(i, \text{confirm } \varphi; \alpha_1, s')$  for all  $s' \in [s]_{R(i)}$ . Hence it follows that  $\exists \beta_1 \in \text{CR}_M^C(i, \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi}, s') \exists \beta_2 \in \text{Agenda}(i, s') \exists \beta_2' \in \text{CR}_M^C(i, \beta_2, s') (\text{Prefix}(\beta_1, \beta_2'))$  implies  $\exists \beta_1 \in \text{CR}_M^C(i, (\text{confirm } \varphi; \alpha_1), s') \exists \beta_2 \in \text{Agenda}(i, s') \exists \beta_2' \in \text{CR}_M^C(i, \beta_2, s') (\text{Prefix}(\beta_1, \beta_2'))$  for all  $s' \in [s]_{R(i)}$ . Then it is indeed the case that from  $M, s \models^C \mathbf{Committed}_i \text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi}$  it follows that  $M, s \models^C \mathbf{Committed}_i (\text{confirm } \varphi; \alpha_1)$ , which suffices to conclude this item.
  4. This item is completely analogous to the previous one.
  5. From the definition of  $\text{CR}_M^C$  and CS it follows that in the case that  $M, s \models^C \varphi$ ,  $\text{CR}_M^C(i, \text{while } \varphi \text{ do } \alpha \text{ od}, s) = \text{CR}_M^C(i, (\text{confirm } \varphi; \alpha); \text{while } \varphi \text{ do } \alpha \text{ od}, s)$ . By a similar argument as the one given in the proof of item 3 one concludes that  $M, s \models^C \mathbf{Committed}_i \text{while } \varphi \text{ do } \alpha \text{ od} \wedge \mathbf{K}_i \varphi \rightarrow \mathbf{Committed}_i ((\text{confirm } \varphi; \alpha); \text{while } \varphi \text{ do } \alpha \text{ od})$ , which concludes item 5.

□

## References

- [1] L. Åqvist. Deontic logic. In D.M. Gabbay and F. Guenther, editors, *Handbook of Philosophical Logic*, volume 2, chapter 11, pages 605–714. D. Reidel, Dordrecht, 1984.
- [2] H.-N. Castañeda. The paradoxes of deontic logic: the simplest solution to all of them in one fell swoop. In Risto Hilpinen, editor, *New Studies in Deontic Logic*, pages 37–85. Reidel, Dordrecht, 1981.
- [3] C. Castelfranchi, D. D'Aloisi, and F. Giacomelli. A framework for dealing with belief-goal dynamics. In M. Gori and G. Soda, editors, *Topics in Artificial Intelligence*, volume 992 of *Lecture Notes in Computer Science (subseries LNAI)*, pages 237–242. Springer-Verlag, 1995.
- [4] L. Cavedon, L. Padgham, A. Rao, and E. Sonenberg. Revisiting rationality for agents with intentions. In X. Yao, editor, *Bridging the Gap: Proceedings of the Eight Australian Joint Conference on Artificial Intelligence*, pages 131–138. World Scientific, 1995.
- [5] P.R. Cohen and H.J. Levesque. Intention is choice with commitment. *Artificial Intelligence*, 42:213–261, 1990.
- [6] F. Dignum and B. van Linder. Modelling rational agents in a dynamic environment: Putting Humpty Dumpty together again. In J.L. Fiadeiro and P.-Y. Schobbens, editors, *Proceedings of the 2nd Workshop of the ModelAge Project*, pages 81–91, 1996.
- [7] F. Dignum and B. van Linder. Modelling social agents in a dynamic environment: Making agents talk. Submitted, 1996.

- [8] F. Dignum, J.-J.Ch. Meyer, R.J. Wieringa, and R. Kuiper. A modal approach to intentions, commitments and obligations: Intention plus commitment yields obligation. In M.A. Brown and J. Carmo, editors, *Deontic Logic, Agency and Normative Systems*, Springer Workshops in Computing, pages 80–97. Springer-Verlag, 1996.
- [9] R. Fagin and J.Y. Halpern. Belief, awareness and limited reasoning. *Artificial Intelligence*, 34:39–76, 1988.
- [10] G. Gazdar, G. Pullum, R. Carpenter, E. Klein, T. Hukari, and R. Levine. Category structures. *Computational Linguistics*, 14:1–19, 1988.
- [11] J.Y. Halpern and Y. Moses. A guide to completeness and complexity for modal logics of knowledge and belief. *Artificial Intelligence*, 54:319–379, 1992.
- [12] D. Harel. Dynamic logic. In D.M. Gabbay and F. Guenther, editors, *Handbook of Philosophical Logic*, volume 2, chapter 10, pages 497–604. D. Reidel, Dordrecht, 1984.
- [13] J. Hintikka. *Knowledge and Belief*. Cornell University Press, Ithaca, NY, 1962.
- [14] W. van der Hoek, B. van Linder, and J.-J. Ch. Meyer. A logic of capabilities. Technical Report IR-330, Vrije Universiteit Amsterdam, July 1993.
- [15] W. van der Hoek, B. van Linder, and J.-J. Ch. Meyer. A logic of capabilities. In A. Nerode and Yu. V. Matiyasevich, editors, *Proceedings of the Third International Symposium on the Logical Foundations of Computer Science (LFCS'94)*, volume 813 of *Lecture Notes in Computer Science*, pages 366–378. Springer-Verlag, 1994.
- [16] W. van der Hoek, B. van Linder, and J.-J. Ch. Meyer. Unravelling nondeterminism: On having the ability to choose (extended abstract). In P. Jorrand and V. Sgurev, editors, *Proceedings of the Sixth International Conference on Artificial Intelligence: Methodology, Systems, Applications (AIMSA'94)*, pages 163–172. World Scientific, 1994.
- [17] A. Kenny. *Will, Freedom and Power*. Basil Blackwell, Oxford, 1975.
- [18] Y. Lespérance, H. Levesque, F. Lin, D. Marcu, R. Reiter, and R. Scherl. Foundations of a logical approach to agent programming. In M. Wooldridge, J.P. Müller, and M. Tambe, editors, *Intelligent Agents Volume II – Agent Theories, Architectures, and Languages*, volume 1037 of *Lecture Notes in Computer Science (subseries LNAI)*, pages 331–347. Springer-Verlag, 1996.
- [19] H. Levesque. A logic of implicit and explicit belief. In *Proceedings of the Fourth National Conference on Artificial Intelligence (AAAI'84)*, pages 198–202. The AAAI Press/The MIT Press, 1984.
- [20] F. Lin and R. Reiter. State constraints revisited. *Journal of Logic and Computation, Special Issue on Actions and Processes*, 1994.
- [21] B. van Linder, W. van der Hoek, and J.-J. Ch. Meyer. Tests as epistemic updates. In A.G. Cohn, editor, *Proceedings of the 11th European Conference on Artificial Intelligence (ECAI'94)*, pages 331–335. John Wiley & Sons, 1994.

- [22] B. van Linder, W. van der Hoek, and J.-J. Ch. Meyer. Actions that make you change your mind. In A. Laux and H. Wansing, editors, *Knowledge and Belief in Philosophy and Artificial Intelligence*, pages 103–146. Akademie Verlag, 1995.
- [23] B. van Linder, W. van der Hoek, and J.-J. Ch. Meyer. The dynamics of default reasoning (extended abstract). In C. Froidevaux and J. Kohlas, editors, *Symbolic and Quantitative Approaches to Reasoning and Uncertainty*, volume 946 of *Lecture Notes in Computer Science (subseries LNAI)*, pages 277–284. Springer-Verlag, 1995. Full version to appear in *Data & Knowledge Engineering*.
- [24] B. van Linder, W. van der Hoek, and J.-J. Ch. Meyer. Seeing is believing (and so are hearing and jumping). In M. Gori and G. Soda, editors, *Topics in Artificial Intelligence (Proc. AIIA'95)*, volume 992 of *Lecture Notes in Computer Science (subseries LNAI)*, pages 402–423. Springer-Verlag, 1995. Full version to appear in *Journal of Logic, Language and Information*.
- [25] J.-J. Ch. Meyer and R.J. Wieringa. Deontic logic: A concise overview. In J.-J. Ch. Meyer and R.J. Wieringa, editors, *Deontic Logic in Computer Science*, chapter 1, pages 3–16. John Wiley & Sons, 1993.
- [26] R.C. Moore. Reasoning about knowledge and action. Technical Report 191, SRI International, 1980.
- [27] G. Plotkin. A structural approach to operational semantics. Technical Report DAIME FN-19, Aarhus University, 1981.
- [28] V. Rantala. Impossible worlds semantics and logical omniscience. *Acta Philosophica Fennica*, 35:106–115, 1982.
- [29] A.S. Rao and M.P. Georgeff. Asymmetry thesis and side-effect problems in linear time and branching time intention logics. In J. Mylopoulos and R. Reiter, editors, *Proceedings of the Twelfth International Joint Conference on Artificial Intelligence (IJCAI'91)*, pages 498–504. Morgan Kaufmann, 1991.
- [30] A.S. Rao and M.P. Georgeff. Modeling rational agents within a BDI-architecture. In J. Allen, R. Fikes, and E. Sandewall, editors, *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning (KR'91)*, pages 473–484. Morgan Kaufmann, 1991.
- [31] J. Raz, editor. *Practical Reasoning*. Oxford Readings in Philosophy. Oxford University Press, 1978.
- [32] R. Reiter. Proving properties of states in the situation calculus. *Artificial Intelligence*, 64:337–351, 1993.
- [33] A. Ross. Imperatives and logic. *Theoria*, 7:53–71, 1941.
- [34] E. Sandewall and Y. Shoham. Non-monotonic temporal reasoning. In C.J. Hogger D.M. Gabbay and J.A. Robinson, editors, *Handbook of Logic in Artificial Intelligence and Logic Programming*, volume 4, pages 439–498. Clarendon Press, Oxford, 1995.
- [35] Y. Shoham. Agent-oriented programming. *Artificial Intelligence*, 60:51–92, 1993.

- [36] Y. Shoham. Agent-oriented programming. *Artificial Intelligence*, 60(1):51–92, 1993.
- [37] M.P. Singh. *Multiagent Systems*, volume 799 of *LNAI*. Springer-Verlag, Berlin-Heidelberg, 1994.
- [38] E. Spaan. *Complexity of Modal Logics*. PhD thesis, Universiteit van Amsterdam, 1993.
- [39] S.R. Thomas. *PLACA, an Agent-Oriented Programming Language*. PhD thesis, Stanford University, 1993.
- [40] M. Wooldridge and N. R. Jennings. Intelligent agents: Theory and practice. *The Knowledge Engineering Review*, 10(2):115–152, 1995.
- [41] G.H. von Wright. *Norm and Action*. Routledge & Kegan Paul, London, 1963.
- [42] G.H. von Wright. On so-called practical inference. In J. Raz, editor, *Practical Reasoning*, chapter III, pages 46–62. Oxford University Press, 1978.