



**HAL**  
open science

## On the logical foundations of moral agency

Emiliano Lorini

► **To cite this version:**

Emiliano Lorini. On the logical foundations of moral agency. 11th International Conference on Deontic Logic in Computer Science (DEON 2012), Jul 2012, Bergen, Norway. pp.108-122, <10.1007/978-3-642-31570-1\_8>. <hal-03671073>

**HAL Id: hal-03671073**

**<https://hal.science/hal-03671073v1>**

Submitted on 19 May 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

# On the logical foundations of moral agency\*

Emiliano Lorini  
Université de Toulouse, IRIT-CNRS, France

## Abstract

The aim of this work is to provide a logical analysis of moral agency. Although this concept has been extensively studied in social philosophy and in social sciences, it has been far less studied in the field of deontic logic and multi-agent systems (MASs). We discuss different aspects of moral agency such as the distinction between desires and moral values and the concept of moral agent. All these concepts are formalized in a variant of STIT logic with explicit actions.

## 1 Introduction

Although the concepts of morality and moral agency have been extensively studied in social philosophy and in social sciences, they have been far less studied in the areas of multi-agent systems (MASs) and normative multi-agent systems (NorMAS). Some works have been done on the extension of the BDI (Belief, Desire, Intention) model with normative concepts such as obligation [6, 10], but none of them have really focused on the integration of moral aspects into the architecture of a cognitive agent. Developing formal models of cognitive agents integrating a moral dimension is a promising research avenue for these two areas. Indeed, as shown by social scientists [9, 8], decisions of human agents are often affected by moral sentiments and moral concerns (*e.g.*, concerns for fairness or equity). Therefore, to take the presence of moral attitudes into account becomes extremely important when developing formal and computational models of social procedures to be applied to human societies, and when developing logical models of artificial agents which are expected to interact with human agents (*e.g.*, trading agents, recommender systems, tutoring agents, etc.). The aim of this paper is to propose a logical framework in which different aspects of moral agency can be formalized such as the distinction between desires and moral attitudes and the concept of moral agent.

The rest of the paper is organized as follows. Section 2 establishes the conceptual basis of the logical analysis of moral agency developed in the second part of the paper, whereas in Section 3 the logic DL-MA (*Dynamic Logic of Mental attitudes and joint Actions*) is presented. As we will show, DL-MA can be seen as a variant of STIT logic — the logic of “Seeing To It That” [4, 14] — with explicit actions. In the second part

---

\*The proceedings version of the paper published by Springer-Verlage (LNCS series) contains a mistake in Definition 2. The mistake has been amended in this version of the paper.

of the paper (Section 4), the logic DL-MA is used to develop a logical analysis of the different aspects of moral agency discussed in Section 2.

## 2 Moral agency: conceptual basis

Some background and clarifications of the notion of moral agency are needed in order to ground the logical analysis presented in Section 4 on a solid conceptual basis.

**Desires vs. moral values** A model of moral agency should be able to explain the two different origins of an agent's motivations. Some motivations originate from the agents's desires. A desire can be conceived as an agent's attitude consisting in an anticipatory mental representation of a pleasant state of affairs (representational dimension of desires) that motivates the agent to achieve it (motivational dimension of desires). In this perspective, the motivational dimension of an agent's desire is realized through its representational dimension. For example when an agent desires to be at the Japanese restaurant eating sushi, he imagines himself eating sushi at the Japanese restaurant and this representation gives him pleasure. This pleasant representation motivates him to go to the Japanese restaurant in order to eat sushi. Agents are motivated not only by their desires but also by their moral values. Moral values, and more generally moral attitudes (ideals, standards, etc.), originate from an agent's capability of discerning what from his point of view is (morally) *good* from what is (morally) *bad*. If an agent has a certain ideal  $\varphi$ , then he thinks that the realization of the state of affairs  $\varphi$  ought to be promoted because  $\varphi$  is *good* in itself. A similar distinction has also been made by philosophers and by social scientists. For instance, Searle [19] has recently proposed a theory of how an agent may want something without desiring it and on the problem of reasons for acting based on moral values and independent from desires. In his theory of morality [12], Harsanyi distinguishes a person's *ethical preferences* from her *subjective preferences* and argues that a moral choice is a choice that is based on ethical preferences.<sup>1</sup>

The distinction between desires and moral values allows to identify two different kinds of moral dilemmas. The first kind of moral dilemma is the one which originates from the logical conflict between two moral values. The paradigmatic example is the situation of a soldier during a war. As a member of the army, the soldier feels obliged to kill his enemies, if this is the only way to defend his country. But, as a catholic, he thinks that human life should be respected. Therefore, he feels morally obliged not to kill other people. The other kind of moral dilemma is the one which originates from the logical conflict between desires and moral values. The paradigmatic example is that of Adam and Eve in the garden of Eden. They are tempted by the desire to eat the forbidden fruit and, at the same time, they have a moral obligation not to do it.

---

<sup>1</sup>An alternative to Harsanyi's dual theory of morality is Sen's theory [20]. According to Sen, moral attitudes are kind of meta-level attitudes. In particular, moral judgments are *rankings of preference rankings*. In this sense, differently from Harsanyi's theory, in Sen's theory subjective preferences and ethical preferences are not at the same level. A comparison between these two theories of morality goes beyond the objectives of the present work.

**Moral agents** An important distinction for a logical theory of morality is *self-regarding* agent versus *moral* agent. A purely *self-regarding* agent is an agent who acts in order to maximize the satisfaction of his own desires, while a purely *moral agent* is an agent who acts in order to maximize the fulfilment of his own moral values. In other words, if an agent is purely self-regarding, the utility of this act for him coincides with the personal good the agent will obtain by performing this action — where the agent’s personal good coincides with the satisfaction of the agent’s own desires —. If an agent is purely moral, the utility of this act for him coincides with the moral good the agent will promote by performing this action — where the agent’s promotion of the moral good coincides with the accomplishment of his own moral values —. The notion of self-regarding agent should not be confused with the *rationality assumption* of classical decision and game theory: according to classical decision and game theory, individuals are rational in the sense that they maximize their utility. The notions of self-regarding agent and moral agent are not in contradiction with this assumption. We can safely say that a given agent acts to maximize his utility even even though he does not act to maximize the satisfaction of his own desires (as he also cares about the fulfilment of his own moral values). Therefore, the person’s act is moral even though the person is still acting to maximize her utility. Of course, purely self-regarding agents and purely moral agents are just extremes cases. An agent is more or less moral depending on whether the utility of a given option for him is more or less affected by his moral values. More precisely, the higher is the influence of the agent’s moral values in evaluating the utility of a given decision option, more moral the agent is. The extent to which an agent’s utility is affected by his moral values can be called *degree of moral sensitivity*.<sup>2</sup>

### 3 Logical framework

In the following sections the logic DL-MA (*Dynamic Logic of Mental attitudes and joint Actions*) is presented. DL-MA is a modal logic which supports reasoning about actions and capabilities of agents and of coalitions of agents. DL-MA also allows to describe epistemic states of agents as well as their desires and moral values. We first present the syntax and the semantics of DL-MA (Sections 3.1, 3.2 and 3.3). A complete axiomatization of the logic is given in Section 3.4. As we will show in Section 3.5, DL-MA is tightly related to STIT logic, as it can be seen as a variant of STIT with explicit actions.

#### 3.1 Syntax

Assume a countable set of atomic propositions  $Atm = \{p, q, \dots\}$ , a finite set of agents  $Agt = \{i_1, \dots, i_{|Ag|}\}$ , a finite set of atomic action types  $Act = \{a, b, \dots\}$ , and a finite set of natural numbers  $Num = \{x \in \mathbb{N} : 0 \leq x \leq \max\}$ , with  $\max \in \mathbb{N} \setminus \{0\}$ . We assume that the set  $Act$  includes the (in)action skip, *i.e.*, the action of doing nothing.

We define  $Prop$  to be the set of propositional formulas, that is, the set of all Boolean combinations of atomic propositions.

<sup>2</sup>For a similar idea in current economic models of morality see, *e.g.*, [3, 1].

For each agent  $i \in \text{Agt}$ , the set  $\text{Act}_i = \{a_i : a \in \text{Act}\}$  identifies the set of agent  $i$ 's actions.  $2^{\text{Agt}^*} = 2^{\text{Agt}} \setminus \{\emptyset\}$  is the set of non-empty sets of agents, also called *coalitions*. Elements of  $2^{\text{Agt}^*}$  are denoted by symbols  $H, J, \dots$ . Let

$$J\text{Act} = \{\{a_{i_1}, \dots, a_{i_{|Agt|}}\} : a_{i_k} \in \text{Act}_{i_k} \text{ for all } i_k \in \text{Agt}\}$$

be the set of all possible *joint actions* of all agents. One might think of  $J\text{Act}$  as the set of all possible *strategy profiles* in the game theoretic sense. Just as in game theory we suppose that at a given time point every agent performs exactly one action, and that all actions of different agents occur in parallel. Elements of  $J\text{Act}$  are denoted by symbols  $\delta, \delta', \delta'', \dots$ . We let  $\delta_i$  denote the element in  $\delta$  corresponding to agent  $i$ . Given a coalition  $H$  and a joint action  $\delta$  of all agents,  $\delta_H = \{a_i : i \in H \text{ and } a_i \in \delta\}$  is coalition  $H$ 's part in the joint action  $\delta$ . Let  $J\text{Act}_H = \{\delta_H : \delta \in J\text{Act}\}$  be the set of all possible joint actions of coalition  $H$ . Finally, let  $J\text{Act}^*$  be the set of all (possibly infinite) sequences of joint actions. Elements of  $J\text{Act}^*$  are denoted by symbols  $\epsilon, \epsilon', \epsilon'', \dots$ . For every  $\epsilon_1, \epsilon_2 \in J\text{Act}^*$ , we write  $\epsilon_1 \sqsubseteq \epsilon_2$  to mean that either  $\epsilon_1 = \epsilon_2$  or  $\epsilon_1$  is an initial subsequence of  $\epsilon_2$ , *i.e.*, there is  $\epsilon_3 \in J\text{Act}^*$  such that  $\epsilon_2 = \epsilon_1; \epsilon_3$ . The empty sequence of joint actions is denoted by *nil*. Infinite sequences of joint actions are called *histories*. Let  $\text{Hist}$  be the set of all histories. Elements of  $\text{Hist}$  are denoted by symbols  $h, h', \dots$ .

The language of DL-MA is defined by the following grammar in Backus-Naur Form:

$$\begin{aligned} \epsilon &::= \delta \mid \epsilon; \epsilon \\ \varphi &::= p \mid \text{occ}(\epsilon) \mid \text{des}_{i,h} \mid \text{idl}_{i,h} \mid \neg \varphi \mid \varphi_1 \wedge \varphi_2 \mid \llbracket \delta \rrbracket \varphi \mid \square \varphi \mid K_i \varphi \end{aligned}$$

where  $p$  ranges over  $\text{Atm}$ ,  $i$  ranges over  $\text{Agt}$ ,  $\delta$  ranges over  $J\text{Act}$ ,  $\epsilon$  ranges over  $J\text{Act}^*$  and  $h$  ranges over  $\text{Num}$ . The other Boolean constructions  $\top, \perp, \vee, \rightarrow$  and  $\leftrightarrow$  are defined from  $p, \neg$  and  $\wedge$  in the standard way.

DL-MA has special atomic formulas of three different kinds. The atomic formulas  $\text{occ}(\epsilon)$  represent information about occurrences of joint action sequences. The formula  $\text{occ}(\epsilon)$  has to be read “the joint action sequence  $\epsilon$  is going to occur”.

The other atomic formulas  $\text{des}_{i,h}$  and  $\text{idl}_{i,h}$  are used to rank the histories that an agent can imagine at a given world according to their *desirability* degree and to their *ideality* degree for the agent. *Desirability* captures the quantitative dimension of desires (*i.e.*, how much a given history promotes the satisfaction of the agent's desires), whereas *ideality* captures the quantitative dimension of moral values (*i.e.*, how much a given history promotes the fulfilment of the agent's moral values). Formula  $\text{des}_{i,h}$  has to be read “the current history has for agent  $i$  a degree of desirability equal to  $h$ ” while formula  $\text{idl}_{i,h}$  has to be read “the current history has for agent  $i$  a degree of ideality equal to  $h$ ”. The following two abbreviations express respectively that “the current world has for agent  $i$  a degree of desirability equal to or higher than  $h$ ” and “the current world has for agent  $i$  a degree of ideality equal to or higher than  $h$ ”:

$$\text{des}_{i,\geq h} \stackrel{\text{def}}{=} \bigvee_{h \leq k \leq \max} \text{des}_{i,k} \qquad \text{idl}_{i,\geq h} \stackrel{\text{def}}{=} \bigvee_{h \leq k \leq \max} \text{idl}_{i,k}$$

The logic DL-MA has three kinds of modal operators:  $\llbracket \delta \rrbracket$ ,  $\square$  and  $K_i$ .  $\square$  is the historic necessity operator of STIT logic. The formula  $\square \varphi$  has to be read “ $\varphi$  is true

in all histories passing through the current moment”. We define  $\diamond$  to be the dual of  $\Box$ , i.e.,  $\diamond\varphi \stackrel{\text{def}}{=} \neg\Box\neg\varphi$  where  $\diamond\varphi$  has to be read “ $\varphi$  is true in at least one history passing through the current moment”.

$\llbracket\delta\rrbracket$  is a dynamic operator describing the fact that if the joint action  $\delta$  is performed then it will lead to a state in which a given state of affairs holds. In particular,  $\llbracket\delta\rrbracket\varphi$  has to be read “if the joint  $\delta$  is performed, then  $\varphi$  will be true after its execution”. We define  $\langle\langle\delta\rangle\rangle$  to be the dual of  $\llbracket\delta\rrbracket$ , i.e.,  $\langle\langle\delta\rangle\rangle\varphi \stackrel{\text{def}}{=} \neg\llbracket\delta\rrbracket\neg\varphi$ , where  $\langle\langle\delta\rangle\rangle\varphi$  has to be read “the joint  $\delta$  is performed and  $\varphi$  will be true after its execution”.

Finally,  $K_i$  is a modal operator characterizing the concept of *ex ante* (or *choice-independent*) knowledge in the sense of Aumann & Dreze [2] (see also [18]). The formula  $K_i\varphi$  has to be read “agent  $i$  knows that  $\varphi$  is true independently from his current choice” or “agent  $i$  thinks that  $\varphi$  is true for any choice he could have made”. The dual of the operator  $K_i$  is denoted by  $\widehat{K}_i$ , i.e.,  $\widehat{K}_i\varphi \stackrel{\text{def}}{=} \neg K_i\neg\varphi$ . Aumann & Dreze distinguish *ex ante* knowledge from *interim* knowledge. *Ex ante* knowledge characterizes an agent’s knowledge assuming that no decision has yet been made by him, whereas *interim* knowledge characterizes an agent’s knowledge assuming that the agent has made his decision about which action to take, but might still be uncertain about the decisions of others. The concept of *interim* knowledge is expressed in DL-MA by the following operator  $K_i^*$ :

$$K_i^*\varphi \stackrel{\text{def}}{=} \bigwedge_{a_i \in \text{Act}_i} (\text{choose}(a_i) \rightarrow K_i(\text{choose}(a_i) \rightarrow \varphi))$$

where  $K_i^*\varphi$  has to be read “the agent  $i$  knows that  $\varphi$  is true given his current choice”.

The following abbreviations will be useful to axiomatize DL-MA. For any coalition  $H$ , joint action sequence  $\epsilon$  and joint action  $\delta_H$  we define:

$$\text{choose}(\epsilon, \delta_H) \stackrel{\text{def}}{=} \bigvee_{\delta \in \text{JAct}: \delta_H \subseteq \delta} \text{occ}(\epsilon; \delta) \quad \text{able}(\epsilon, \delta_H) \stackrel{\text{def}}{=} \diamond \text{choose}(\epsilon, \delta_H)$$

$\text{choose}(\epsilon, \delta_H)$  has to be read “at the end of the joint action sequence  $\epsilon$ , the agents in  $H$  will choose the joint action  $\delta_H$ ”, whereas  $\text{able}(\epsilon, \delta_H)$  has to be read “the agents in  $H$  will be able to perform the joint action  $\delta_H$ , at the end of the joint action sequence  $\epsilon$ ”. For notational convenience, when  $\epsilon = \text{nil}$ , we write  $\text{choose}(\delta_H)$  instead of  $\text{choose}(\text{nil}, \delta_H)$  and  $\text{able}(\delta_H)$  instead of  $\text{able}(\text{nil}, \delta_H)$ . Besides, we simply write  $\text{choose}(\epsilon, a_i)$  instead of  $\text{choose}(\epsilon, \{a_i\})$ ,  $\text{able}(\epsilon, a_i)$  instead of  $\text{able}(\epsilon, \{a_i\})$ ,  $\text{choose}(a_i)$  instead of  $\text{choose}(\{a_i\})$  and  $\text{able}(a_i)$  instead of  $\text{able}(\{a_i\})$ .

## 3.2 Action description

Similarly to Situation Calculus [17], in DL-MA actions are described in terms of their positive and negative effect preconditions. In particular, we introduce two functions:

$$\begin{aligned} \gamma^+ &: \text{Agt} \times \text{Act} \times \text{Atm} \longrightarrow \text{Prop} \\ \gamma^- &: \text{Agt} \times \text{Act} \times \text{Atm} \longrightarrow \text{Prop} \end{aligned}$$

mapping agents, actions and atomic propositions to propositional formulas. The formula  $\gamma^+(i, a, p)$  describes the *positive effect preconditions* of action  $a$  performed by

agent  $i$  with respect to  $p$ , whereas  $\gamma^-(i, a, p)$  describes the *negative effect preconditions* of action  $a$  performed by agent  $i$  with respect to  $p$ . Formula  $\gamma^+(i, a, p)$  represents the conditions under which agent  $i$  will make  $p$  true by performing action  $a$ , if no other agent interferes with  $i$ 's action; while  $\gamma^-(i, a, p)$  represents the conditions under which agent  $i$  will make  $p$  false by performing action  $a$ , if no other agent interferes with  $i$ 's action. We assume that “making  $p$  true” means changing the truth value of  $p$  from false to true, whereas “making  $p$  false” means changing the truth value of  $p$  from true to false. We make the following *coherence assumption*:

( $COH_\gamma$ ) for every  $i \in Agt$ ,  $a \in Act$  and  $p \in Atm$ ,  $\gamma^+(i, a, p)$  and  $\gamma^-(i, a, p)$  must be logically inconsistent.

$COH_\gamma$  ensures that actions do not have contradictory effects. As to the (in)action skip of doing nothing, we assume that for every agent  $i \in Agt$ ,  $\gamma^+(i, skip, p) = \gamma^-(i, skip, p) = \perp$ , *i.e.*, an agent cannot change the truth value of  $p$  by doing nothing.

### 3.3 Semantics

The semantics of DL-MA is a possible world semantics with accessibility relations associated with each modal operators, with functions for *desirability* and *ideality*, and with a function designating the history starting in a given world.

**Definition 1** A DL-MA model is a tuple  $M = \langle W, \mathcal{H}, \equiv, \{\mathcal{E}_i : i \in Agt\}, \kappa_{des}, \kappa_{idl}, \mathcal{V} \rangle$  where:

- $W$  is a set of states (or worlds),
- $\mathcal{H}$  is a total function  $\mathcal{H} : W \rightarrow Hist$ ,
- $\equiv$  and every  $\mathcal{E}_i$  are equivalence relations between states in  $W$ ,
- $\kappa_{des} : W \times Agt \rightarrow Num$  and  $\kappa_{idl} : W \times Agt \rightarrow Num$  are total functions mapping worlds and agents to natural numbers in  $Num$ ,
- $\mathcal{V} : W \rightarrow 2^{Atm}$  is a valuation function.

As usual  $p \in \mathcal{V}(w)$  means that proposition  $p$  is true at world  $w$ .

For every world  $w \in W$ ,  $\mathcal{H}(w)$  identifies the history starting in  $w$ . For notational convenience, for all  $\epsilon \in JAct^*$ ,  $i \in Agt$  and  $a_i \in Act_i$  we write  $\epsilon; a_i \sqsubseteq \mathcal{H}(w)$  if there is  $\delta \in JAct$  such that  $a_i \in \delta$  and  $\epsilon; \delta \sqsubseteq \mathcal{H}(w)$ .

$\equiv$ -equivalence classes are called *moments*. If  $w$  and  $v$  belong to the same moment (*i.e.*,  $w \equiv v$ ), then the history starting in  $w$  (*i.e.*,  $\mathcal{H}(w)$ ) and the history starting in  $v$  (*i.e.*,  $\mathcal{H}(v)$ ) are said to be alternative histories (*viz.*, histories starting at the same moment). The concept of moment is the one used in STIT logic [4, 14] and, more generally, in the Ockhamist theory of time [21, 23].

The equivalence relations  $\mathcal{E}_i$  are used to interpret the *ex ante* epistemic operators  $K_i$ . They can be viewed as functions from  $W$  to  $2^W$ . Therefore, we can write  $\mathcal{E}_i(w) = \{v \in W : w\mathcal{E}_i v\}$ . The set  $\mathcal{E}_i(w)$  is the agent  $i$ 's *information state* at world  $w$ : the

set of worlds that at  $w$  agent  $i$  considers epistemically possible independently from his current choice or, more shortly, agent  $i$ 's set of epistemic alternatives at  $w$ . As  $\mathcal{E}_i$  is an equivalence relation, if  $w\mathcal{E}_iv$  then agent  $i$  has the same information state at  $w$  and  $v$ .

The functions  $\kappa_{\text{des}}$  and  $\kappa_{\text{idl}}$  represent respectively desirability gradings and ideality gradings of the possible worlds and are used to interpret the atomic formulas  $\text{des}_{i,h}$  and  $\text{idl}_{i,h}$ .  $\kappa_{\text{des}}(w, i) = h$  means that according to the agent  $i$  the world  $w$  has a degree of desirability  $h$ , whereas  $\kappa_{\text{idl}}(w, i) = h$  means that according to the agent  $i$  the world  $w$  has a degree of ideality  $h$ .

We impose the following constraints on DL-MA models. For all  $w \in W$ ,  $\delta \in JAct$ ,  $\epsilon \in JAct^*$ ,  $i \in Agt$  and  $a_i \in Act_i$  we have:

- (C1) if for all  $i \in Agt$  there is  $u_i$  such that  $w \equiv u_i$  and  $\epsilon; \delta_i \sqsubseteq \mathcal{H}(u_i)$ , then there is  $u$  such that  $w \equiv u$  and  $\epsilon; \delta \sqsubseteq \mathcal{H}(u)$ ;
- (C2) if there is  $v$  such that  $w \equiv v$  and  $\epsilon; a_i \sqsubseteq \mathcal{H}(v)$  then, for all  $u$  such that  $w\mathcal{E}_iu$ , there is  $z$  such that  $u \equiv z$  and  $\epsilon; a_i \sqsubseteq \mathcal{H}(z)$ ;
- (C3) if there is  $v$  such that  $w \equiv v$  and  $\epsilon; a_i \sqsubseteq \mathcal{H}(v)$ , then there is  $u$  such that  $w\mathcal{E}_iu$  and  $\epsilon; a_i \sqsubseteq \mathcal{H}(u)$ .

According to the Constraint **C1**, if every individual action in a joint action  $\delta$  can be chosen at the end of the joint action sequence  $\epsilon$ , then the individual actions in  $\delta$  can be chosen simultaneously at the end of  $\epsilon$ . The Constraint **C1** is a variant of the assumption of *independence of agents* of STIT logic. More intuitively, this means that agents can never be deprived of choices due to the choices made by other agents. The Constraint **C2** is a basic assumption about agents' knowledge over their abilities: if an agent  $i$  can perform the action  $a$  at the end of the joint action sequence  $\epsilon$ , then he knows this. In other words, an agent has perfect knowledge about his ability to perform a given action. The Constraint **C3** characterizes the basic property of *ex ante* knowledge: if an agent  $i$  can perform the action  $a$  at the end of the joint action sequence  $\epsilon$ , then there is a history that the agent considers possible independently from his current choice in which he performs the action  $a$  at the end of the joint action sequence  $\epsilon$ . In other words, for every action that an agent is able to perform, there is a history that the agent considers possible in which he performs this action. The rules defining truth conditions of DL-MA formulas are the standard ones for Boolean constructions  $p$ ,  $\neg$  and  $\wedge$  plus:

- $M, w \models \text{occ}(\epsilon)$  iff  $\epsilon \sqsubseteq \mathcal{H}(w)$
- $M, w \models \text{des}_{i,h}$  iff  $\kappa_{\text{des}}(w, i) = h$
- $M, w \models \text{idl}_{i,h}$  iff  $\kappa_{\text{idl}}(w, i) = h$
- $M, w \models \llbracket \delta \rrbracket \varphi$  iff if  $M, w \models \text{occ}(\delta)$  then  $M^\delta, w \models \varphi$
- $M, w \models \Box \varphi$  iff  $M, v \models \varphi$  for all  $v$  such that  $w \equiv v$
- $M, w \models K_i \varphi$  iff  $M, v \models \varphi$  for all  $v$  such that  $w\mathcal{E}_iv$

where model  $M^\delta$  is defined according to Definition 2 below.

**Definition 2 (Update via joint action)** Given a DL-MA model  $M = \langle W, \mathcal{H}, \equiv, \{\mathcal{E}_i : i \in \text{Agt}\}, \kappa_{\text{des}}, \kappa_{\text{idl}}, \mathcal{V} \rangle$ , the update of  $M$  by  $\delta$  is defined to be  $M^\delta = \langle W^\delta, \mathcal{H}^\delta, \equiv^\delta, \{\mathcal{E}_i^\delta : i \in \text{Agt}\}, \kappa_{\text{des}}^\delta, \kappa_{\text{idl}}^\delta, \mathcal{V}^\delta \rangle$  where for all  $i \in \text{Agt}$  and  $h \in \text{Hist}$ :

$$\begin{aligned}
W^\delta &= \{w \in W : M, w \models \text{occ}(\delta)\} \\
\mathcal{H}^\delta(w) &= h \text{ if } \mathcal{H}(w) = \delta'; h \text{ for some } \delta' \in \text{JAct} \\
\equiv^\delta &= \equiv \cap (W^\delta \times W^\delta) \\
\mathcal{E}_i^\delta &= \mathcal{E}_i \cap (W^\delta \times W^\delta) \\
\kappa_{\text{des}}^\delta(w, i) &= \kappa_{\text{des}}(w, i) \\
\kappa_{\text{idl}}^\delta(w, i) &= \kappa_{\text{idl}}(w, i) \\
\mathcal{V}^\delta(w) &= (\mathcal{V}(w) \setminus \{p : \exists a_i \in \delta \text{ such that } M, w \models \gamma^-(i, a, p) \text{ and} \\
&\quad \nexists b_j \in \delta \text{ such that } M, w \models \gamma^+(j, b, p)\}) \cup \\
&\quad \{p : \exists a_i \in \delta \text{ such that } M, w \models \gamma^+(i, a, p) \text{ and} \\
&\quad \nexists b_j \in \delta \text{ such that } M, w \models \gamma^-(j, b, p)\}
\end{aligned}$$

The performance of a joint action  $\delta$  modifies the physical facts via the positive effect preconditions and the negative effect preconditions, defined in Section 3.2 (see the definition of  $\mathcal{V}^\delta$ ). In particular, if there is an action in the joint action  $\delta$  whose positive effect preconditions with respect to  $p$  hold and there is no other action in the joint action  $\delta$  whose negative effect preconditions with respect to  $p$  hold, then  $p$  will be true after the occurrence of  $\delta$ ; if there is an action in the joint action  $\delta$  whose negative effect preconditions with respect to  $p$  hold and there is no other action in the joint action  $\delta$  whose positive effect preconditions with respect to  $p$  hold, then  $p$  will be false after the occurrence of  $\delta$ . Besides, the occurrence of the joint action  $\delta$  makes the current history advance one step forward (see the definition of  $\mathcal{H}^\delta$ ). As to the equivalence relations  $\equiv$  and  $\mathcal{E}_i$  for historic necessity and *ex ante* knowledge, they are restricted to the set of worlds in which the joint action  $\delta$  occurs (see the definitions of  $\equiv^\delta$  and  $\mathcal{E}_i^\delta$ ). Finally, the joint action  $\delta$  does not modify the agents' desirability and ideality gradings over the histories (see the definitions of  $\kappa_{\text{des}}^\delta, \kappa_{\text{idl}}^\delta$  and  $W^\delta$ ). As stated by the following proposition, the update via a joint action preserves the constraints on DL-MA-models.

**Proposition 1** *If  $M$  is a DL-MA-model then  $M^\delta$  is a DL-MA-model too.*

PROOF. As an example, we prove that the model  $M^\delta$  satisfies the Constraint **C1** after assuming that  $M$  is a DL-MA model. Let us suppose that for all  $i \in \text{Agt}$  there is  $u_i$  such that  $w \equiv^\delta u_i$  and  $\epsilon; \delta'_i \sqsubseteq \mathcal{H}^\delta(u_i)$ . The latter implies that for all  $i \in \text{Agt}$  there is  $u_i$  such that  $w \equiv u_i$  and  $\delta; \epsilon; \delta'_i \sqsubseteq \mathcal{H}(u_i)$ . As  $M$  is a DL-MA model which satisfies the Constraint **C1**, it follows that there is  $u$  such that  $w \equiv u$  and  $\delta; \epsilon; \delta' \sqsubseteq \mathcal{H}(u)$ . The latter implies that there is  $u$  such that  $w \equiv^\delta u$  and  $\epsilon; \delta' \sqsubseteq \mathcal{H}^\delta(u)$ .

### 3.4 Axiomatization

Our logic DL-MA has so-called reduction axioms. These axioms allow to eliminate all the dynamic operators  $\llbracket \delta \rrbracket$  from formulas. That elimination provides an axiomatics.

**Proposition 2** *The following formulas are DL-MA valid for every  $i \in \text{Agt}$ ,  $a_i \in \text{Act}$ ,  $\delta, \delta' \in \text{JAct}$  and  $\epsilon, \epsilon' \in \text{JAct}^*$ :*

$$\begin{aligned}
\text{occ}(\epsilon) &\rightarrow \bigvee_{\delta \in \text{JAct}} \text{occ}(\epsilon; \delta) && \text{(OneJAct)} \\
\text{occ}(\epsilon; \delta) &\rightarrow \neg \text{occ}(\epsilon; \delta') \text{ if } \delta \neq \delta' && \text{(UniqueJAct)} \\
\text{occ}(\epsilon) &\rightarrow \text{occ}(\epsilon') \text{ if } \epsilon' \sqsubseteq \epsilon && \text{(SubSeqJAct)} \\
\bigvee_{h \in \text{Num}} \text{des}_{i,h} &&& \text{(ComplDes)} \\
\text{des}_{i,h} &\rightarrow \neg \text{des}_{i,k} \text{ if } h \neq k && \text{(UniqueDes)} \\
\bigvee_{h \in \text{Num}} \text{idl}_{i,h} &&& \text{(ComplIdl)} \\
\text{idl}_{i,h} &\rightarrow \neg \text{idl}_{i,k} \text{ if } h \neq k && \text{(UniqueIdl)} \\
\left( \bigwedge_{i \in \text{Agt}} \text{choose}(\epsilon, \delta_i) \right) &\rightarrow \text{occ}(\epsilon; \delta) && \text{(IndepAgt)} \\
\text{able}(\epsilon, a_i) &\rightarrow \text{K}_i \text{able}(\epsilon, a_i) && \text{(KnowCan)} \\
\text{able}(\epsilon, a_i) &\rightarrow \widehat{\text{K}}_i \text{choose}(\epsilon, a_i) && \text{(ExAnteKnow)}
\end{aligned}$$

**Proposition 3** *The following equivalences are DL-MA valid for all  $p \in \text{Atm}$ ,  $i \in \text{Agt}$ ,  $\delta \in \text{JAct}$ ,  $\epsilon \in \text{JAct}^*$  and  $h \in \text{Num}$ :*

$$\begin{aligned}
\llbracket \delta \rrbracket \neg \varphi &\leftrightarrow (\text{occ}(\delta) \rightarrow \neg \llbracket \delta \rrbracket \varphi) \\
\llbracket \delta \rrbracket (\varphi \wedge \psi) &\leftrightarrow (\llbracket \delta \rrbracket \varphi \wedge \llbracket \delta \rrbracket \psi) \\
\llbracket \delta \rrbracket p &\leftrightarrow (\text{occ}(\delta) \rightarrow ((\bigvee_{a_i \in \delta} \gamma^+(i, a, p) \wedge \bigwedge_{b_j \in \delta} \neg \gamma^-(j, b, p)) \vee \\
&\quad (p \wedge \bigwedge_{a_i \in \delta} \neg \gamma^-(i, a, p)) \vee (p \wedge \bigvee_{a_i \in \delta} \gamma^+(i, a, p)))) \\
\llbracket \delta \rrbracket \text{occ}(\epsilon) &\leftrightarrow (\text{occ}(\delta) \rightarrow \text{occ}(\delta; \epsilon)) \\
\llbracket \delta \rrbracket \text{des}_{i,h} &\leftrightarrow (\text{occ}(\delta) \rightarrow \text{des}_{i,h}) \\
\llbracket \delta \rrbracket \text{idl}_{i,h} &\leftrightarrow (\text{occ}(\delta) \rightarrow \text{idl}_{i,h}) \\
\llbracket \delta \rrbracket \Box \varphi &\leftrightarrow (\text{occ}(\delta) \rightarrow \Box (\text{occ}(\delta) \rightarrow \llbracket \delta \rrbracket \varphi)) \\
\llbracket \delta \rrbracket \text{K}_i \varphi &\leftrightarrow (\text{occ}(\delta) \rightarrow \text{K}_i (\text{occ}(\delta) \rightarrow \llbracket \delta \rrbracket \varphi))
\end{aligned}$$

Let L-MA be the fragment of the logic DL-MA without dynamic operators  $\llbracket \delta \rrbracket$ . that is, let the language of L-MA be the set of formulas defined by the following grammar:

$$\varphi ::= p \mid \text{occ}(\epsilon) \mid \text{des}_{i,h} \mid \text{idl}_{i,h} \mid \neg \varphi \mid \varphi_1 \wedge \varphi_2 \mid \Box \varphi \mid \text{K}_i \varphi$$

where  $p$  ranges over  $\text{Atm}$ ,  $i$  ranges over  $\text{Agt}$ ,  $\epsilon$  ranges over  $\text{JAct}^*$  and  $h$  ranges over  $\text{Num}$ . As the rule of replacement of equivalents preserves validity, the equivalences of Proposition 3 together with the rule of replacement of equivalents allow to reduce every DL-MA formula to an equivalent L-MA formula. Call *red* the mapping which iteratively applies the above equivalences from the left to the right, starting from one of the innermost modal operators. *red* pushes the dynamic operators inside the formula, and finally eliminates them when facing an atomic formula.

**Proposition 4** *Let  $\varphi$  be a formula in the language of DL-MA. Then*

1.  $red(\varphi)$  has no dynamic operators  $\llbracket \delta \rrbracket$
2.  $red(\varphi) \leftrightarrow \varphi$  is DL-MA valid
3.  $red(\varphi)$  is DL-MA valid iff  $red(\varphi)$  is L-MA valid.

PROOF.[Sketch] The first item is clear. The second item is proved using Proposition 3 (and the rule of replacement of equivalents). The last item follows from the second item and the fact that DL-MA is a conservative extension of L-MA.

**Theorem 1** *The validities of DL-MA are completely axiomatized by*

- all principles of classical propositional logic
- axiomatization of the normal modal logic S5 for the historic necessity operator  $\square$
- axiomatization of the normal modal logic S5 for each epistemic operator  $K_i$
- the schemas of Proposition 2
- the reduction axioms of Proposition 3
- the rule of replacement of equivalents

*from  $\psi_1 \leftrightarrow \psi_2$  infer  $\varphi \leftrightarrow \varphi[\psi_1/\psi_2]$*

PROOF.[Sketch] To prove soundness is just a routine exercise. The completeness proof proceeds as follows. By standard canonical model argument, it is routine to show that the axioms and rules of inference of the normal modal logic S5 for the operator  $\square$  and for every epistemic operator  $K_i$  together with the principles in Proposition 2 and all principles of classical propositional logic provide a complete axiomatization for L-MA. Now, suppose  $\varphi$  is DL-MA valid. Then  $red(\varphi)$  is valid in L-MA due to Proposition 4. By the completeness of L-MA,  $red(\varphi)$  is also provable there. DL-MA being a conservative extension of L-MA,  $red(\varphi)$  is provable in DL-MA, too. As the reduction axioms and the rule of replacement of equivalents are part of our axiomatics, the formula  $\varphi$  must also be provable in DL-MA.

We write  $\vdash \varphi$  if  $\varphi$  is a DL-MA-theorem.

### 3.5 Relationships between DL-MA and STIT

As pointed out in the introduction, DL-MA can be seen as a variant of STIT with explicit actions: while in STIT an action is identified with the result brought about by a coalition (*i.e.*, in STIT one can only express that a given coalition  $H$  sees to it that  $\varphi$ ), in DL-MA an action is identified both with the result brought about by the coalition and with the means used by the coalition to bring about the result (*i.e.*, in DL-MA one can also express that a given coalition  $H$  sees to it that  $\varphi$  by choosing the joint action

$\delta_H$ ).<sup>3</sup> In DL-MA we can express different concepts of agency which have been studied in the context of STIT theory [4, 14].

For instance, the so-called ‘Chellas’ operator  $\text{CStit}_H$  of STIT can be defined in DL-MA as follows:

$$\text{CStit}_H\varphi \stackrel{\text{def}}{=} \bigvee_{\delta_H \in \text{JAct}_H} (\text{choose}(\delta_H) \wedge \bigwedge_{\delta'_{\text{Agt} \setminus H} \in \text{JAct}_{\text{Agt} \setminus H}} \Box(\text{occ}(\delta_H \cup \delta'_{\text{Agt} \setminus H}) \rightarrow \varphi))$$

This means that the coalition  $H$  sees to it that  $\varphi$  if and only if, the agents in  $H$  choose some joint action  $\delta_H$  such that, no matter what the agents outside  $H$  choose, if the agents in  $H$  choose  $\delta_H$  then  $\varphi$  will be true. The so-called ‘deliberative’ operator  $\text{DStit}_H$  can be defined in DL-MA as follows:

$$\text{DStit}_H\varphi \stackrel{\text{def}}{=} \text{CStit}_H\varphi \wedge \neg\Box\varphi$$

DL-MA integrates a temporal dimension which allows to express that a given coalition  $H$  sees to it that  $\varphi$  is true in the next state of the system. In particular, in DL-MA one can simulate the so-called XSTIT operator  $\text{XStit}_H$  proposed by Broersen [5]:

$$\text{XStit}_H\varphi \stackrel{\text{def}}{=} \text{CStit}_HX\varphi$$

where  $X$  is the operator ‘next’ of linear temporal logic (LTL) which is defined as:

$$X\varphi \stackrel{\text{def}}{=} \bigvee_{\delta \in \text{JAct}} \langle\langle \delta \rangle\rangle\varphi$$

## 4 Moral agency: a logical formalization

In what follows the logic DL-MA is applied to the formalization of the different aspects of moral agency discussed in Section 2. Section 4.1 provides a logical formalization of desires and moral values. In Section 4.2 we define the concept of moral sensitivity as well as a concept of preference based on desires and moral values.

### 4.1 Desires and moral values

For every agent  $i \in \text{Agt}$ , we define two types of dyadic operators, one for desirability and the other for ideality:

$$\begin{aligned} \psi \leq_i^{\text{des}} \varphi &\stackrel{\text{def}}{=} \bigwedge_{h \in \text{Num}} (\widehat{\text{K}}_i(\text{des}_{i, \geq h} \wedge \psi) \rightarrow \widehat{\text{K}}_i(\text{des}_{i, \geq h} \wedge \varphi)) \\ \psi \leq_i^{\text{idl}} \varphi &\stackrel{\text{def}}{=} \bigwedge_{h \in \text{Num}} (\widehat{\text{K}}_i(\text{idl}_{i, \geq h} \wedge \psi) \rightarrow \widehat{\text{K}}_i(\text{idl}_{i, \geq h} \wedge \varphi)) \end{aligned}$$

$\psi \leq_i^{\text{des}} \varphi$  has to be read “ $\varphi$  is for agent  $i$  at least as desirable as  $\psi$ ” and  $\psi \leq_i^{\text{idl}} \varphi$  has to be read “ $\varphi$  is for agent  $i$  at least as ideal as  $\psi$ ”. Corresponding strict orderings are

<sup>3</sup>See also [13] for a variant of STIT with explicit actions.

defined in the expected way as follows:  $\psi <_i^{\text{des}} \varphi \stackrel{\text{def}}{=} (\psi \leq_i^{\text{des}} \varphi) \wedge \neg(\varphi \leq_i^{\text{des}} \psi)$  and  $\psi <_i^{\text{idl}} \varphi \stackrel{\text{def}}{=} (\psi \leq_i^{\text{idl}} \varphi) \wedge \neg(\varphi \leq_i^{\text{idl}} \psi)$ .

As the following proposition highlights the comparative statements  $\psi \leq_i^{\text{des}} \varphi$  and  $\psi \leq_i^{\text{idl}} \varphi$  might also be read as “for every epistemically possible  $\psi$ -state there is an epistemically possible  $\varphi$ -state which is at least as desirable” and “for every epistemically possible  $\psi$ -state there is an epistemically possible  $\varphi$ -state which is at least as ideal”.

**Proposition 5** *For every  $i \in \text{Agt}$  we have:*

- $M, w \models \psi \leq_i^{\text{des}} \varphi$  if and only if for all  $v \in \mathcal{E}_i(w)$ , if  $M, v \models \psi$  then there is  $u \in \mathcal{E}_i(w)$  such that  $\kappa_{\text{des}}(v, i) \leq \kappa_{\text{des}}(u, i)$  and  $M, u \models \varphi$ ,
- $M, w \models \psi \leq_i^{\text{idl}} \varphi$  if and only if for all  $v \in \mathcal{E}_i(w)$ , if  $M, v \models \psi$  then there is  $u \in \mathcal{E}_i(w)$  such that  $\kappa_{\text{idl}}(v, i) \leq \kappa_{\text{idl}}(u, i)$  and  $M, u \models \varphi$ .

As pointed out by [16], there is no consensus in the literature on how preferential statements between formulas should be defined. In [22] other kinds of preference comparisons between formulas (*i.e.*, between sets of states) are defined. For instance, the previous  $\forall\exists$ -reading of preference statements should be distinguished from a  $\forall\forall$ -reading (“for every  $\psi$ -state and for every  $\varphi$ -state the  $\varphi$ -state is at least as desirable/ideal as the  $\psi$ -state”) and a  $\exists\exists$ -reading (“there are a  $\psi$ -state and a  $\varphi$ -state such that the  $\varphi$ -state is at least as desirable/ideal as the  $\psi$ -state”). A logical analysis of such alternative readings of preference statements is postponed to future work.

Proposition 6 captures some central properties of the operators  $\leq_i^{\text{des}}$  and  $\leq_i^{\text{idl}}$ .

**Proposition 6** *Let  $x \in \{\text{des}, \text{idl}\}$ . Then:*

$$\vdash \psi \leq_i^x \psi \quad (6a)$$

$$\vdash ((\varphi_1 \leq_i^x \varphi_2) \wedge (\varphi_2 \leq_i^x \varphi_3)) \rightarrow (\varphi_1 \leq_i^x \varphi_3) \quad (6b)$$

$$\vdash (\varphi_1 \leq_i^x \varphi_2) \vee (\varphi_2 \leq_i^x \varphi_1) \quad (6c)$$

$$\vdash \perp \leq_i^x \top \quad (6d)$$

$$\text{if } \vdash \varphi \rightarrow (\psi_1 \vee \dots \vee \psi_s) \text{ then } \vdash (\varphi \leq_i^x \psi_1) \vee \dots \vee (\varphi \leq_i^x \psi_s) \quad (6e)$$

**PROOF.** We prove theorem (6c) with  $x = \text{des}$  as an example.  $\neg(\varphi_1 \leq_i^x \varphi_2)$  is equivalent to  $\bigvee_{h \in \text{Num}} (\widehat{K}_i(\text{des}_{i, \geq h} \wedge \varphi_1) \wedge K_i(\text{des}_{i, \geq h} \rightarrow \neg\varphi_2))$  which in turn implies  $\bigwedge_{k \in \text{Num}} (\widehat{K}_i(\text{des}_{i, \geq k} \wedge \varphi_2) \rightarrow \widehat{K}_i(\text{des}_{i, \geq k} \wedge \varphi_1))$ . The latter is equivalent to  $(\varphi_2 \leq_i^x \varphi_1)$ . We have proved that  $\neg(\varphi_1 \leq_i^x \varphi_2) \rightarrow (\varphi_2 \leq_i^x \varphi_1)$  which is equivalent to  $(\varphi_1 \leq_i^x \varphi_2) \vee (\varphi_2 \leq_i^x \varphi_1)$ .

Theorems (6a)-(6c) highlight that  $\leq_i^{\text{des}}$  and  $\leq_i^{\text{idl}}$  are total preorders. Theorems (6b)-(6d) are the three fundamental principles of Lewis’s conditional logic [15].

## 4.2 Moral sensitivity and preference based on desires and moral values

We extend the logic DL-MA with special constructions of the form  $\text{moralSensit}(i, m)$  which has to be read “agent  $i$ ’s degree of moral sensitivity is equal to  $m$ ” with  $m \in$

*Num*. We call DL-MA<sup>+</sup> the resulting logic. A DL-MA<sup>+</sup> model is a tuple  $\langle M, \mathcal{S} \rangle$  where  $M$  is a DL-MA model and  $\mathcal{S}$  is a total function:

$$\mathcal{S} : W \times \text{Agt} \longrightarrow \text{Num}$$

capturing the moral sensitivity of an agent at a given state. We assume that an agent is aware of his current degree of moral sensitivity, that is, for every  $i \in \text{Agt}$  and  $w \in W$  we suppose that:

**(C4)** if  $\mathcal{S}(w, i) = m$  then, for all  $v$  such that  $w \mathcal{E}_i v$ ,  $\mathcal{S}(v, i) = m$ .

Constructions  $\text{moralSensit}(i, m)$  are interpreted by means of the function  $\mathcal{S}$  as follows:  
 $M, w \models \text{moralSensit}(i, m)$  if and only if  $\mathcal{S}(w, i) = m$

It is straightforward to adapt the proof of Theorem 1 in order to prove that the logic DL-MA<sup>+</sup> is completely axiomatized by the axioms and rules of inference of the logic DL-MA plus the following axiom schemas:

$$\bigvee_{m \in \text{Num}} \text{moralSensit}(i, m) \quad \textbf{(ComplMoral)}$$

$$\text{moralSensit}(i, m) \rightarrow \neg \text{moralSensit}(i, l) \text{ if } m \neq l \quad \textbf{(UniqueMoral)}$$

$$\text{moralSensit}(i, m) \rightarrow K_i \text{moralSensit}(i, m) \quad \textbf{(KnowMoral)}$$

We use the degree of moral sensitivity as a parameter for calculating the utility of a given history for an agent.

**Definition 3 (Utility)** *Given a DL-MA<sup>+</sup> model  $M = \langle W, \mathcal{H}, \equiv, \{\mathcal{E}_i : i \in \text{Agt}\}, \kappa_{\text{des}}, \kappa_{\text{idl}}, \mathcal{V}, \mathcal{S} \rangle$ , the utility for agent  $i$  of the history starting in the world  $w$ , denoted by  $\kappa_{\text{util}}(w, i)$ , is defined as follows:*

$$\kappa_{\text{util}}(w, i) = \mathcal{S}(w, i) \times \kappa_{\text{idl}}(w, i) + (\max - \mathcal{S}(w, i)) \times \kappa_{\text{des}}(w, i)$$

Moreover, we define

$UScale = \{x : \exists h_1, h_2, h_3 \in \text{Num} \text{ such that } x = h_1 \times h_2 + (\max - h_1) \times h_3\}$   
to be the agents' utility scale.<sup>4</sup>

According to Definition 3, the utility of a history is a function of both the degree of desirability and the degree of the ideality of the history. Degree of moral sensitivity captures the extent to which the utility of a given history is affected by moral values: the higher is the agent's moral sensitivity and the higher is the influence of the degree of ideality in determining the utility of the history; the lower is the agent's moral sensitivity and the higher is the influence of the degree of desirability in determining the utility of the history.

The next step in the analysis is to define the concept of preferred history, as a history that an agent envisages and that has a maximal degree of utility for him. More formally:

**Definition 4 (Preferred histories)** *Given a DL-MA<sup>+</sup> model  $M = \langle W, \mathcal{H}, \equiv, \{\mathcal{E}_i : i \in \text{Agt}\}, \kappa_{\text{des}}, \kappa_{\text{idl}}, \mathcal{V}, \mathcal{S} \rangle$ , agent  $i$ 's set of preferred histories at  $w$ , denoted by  $\mathcal{P}_i(w)$ , is defined as follows:*

$$\mathcal{P}_i(w) = \underset{v \in W : w \mathcal{E}_i v}{\text{argmax}} \kappa_{\text{util}}(v, i)$$

<sup>4</sup>Note that  $UScale$  is finite because  $\text{Num}$  is finite.

We use the notion of preferred history to define the concept of preference. We say that an agent prefers  $\varphi$  to be true if and only if,  $\varphi$  is true in all histories that he prefers. More generally, we assume that each agent envisages a given set of possible histories and he evaluates which are the best ones among them. This leads to a *realistic* notion of preference in the sense that an agent's set of preferred histories is a subset of the set of histories that the agent envisages.

**Definition 5 (Preference,  $\text{Pref}_i$ )** *At world  $w$  agent  $i$  prefers  $\varphi$  to be true, i.e.,  $M, w \models \text{Pref}_i\varphi$ , if and only if  $M, v \models \varphi$  for all  $v \in \mathcal{P}_i(w)$ .*

As the following proposition highlights, the concept of preference semantically defined in Definition 5 is syntactically expressible in the logic DL-MA<sup>+</sup>.

**Proposition 7** *For all  $i \in \text{Agt}$  we have  $M, w \models \text{Pref}_i\varphi$  if and only if*

$$M, w \models \bigvee_{x \in \text{UScale}} (\widehat{K}_i \text{util}_{i,x} \wedge \bigwedge_{y \in \text{UScale}: y > x} K_i \neg \text{util}_{i,y} \wedge K_i (\text{util}_{i,x} \rightarrow \varphi))$$

where

$$\text{util}_{i,x} \stackrel{\text{def}}{=} \bigvee_{k,l,m \in \text{Num}: x=m \times l + (\max - m) \times k} (\text{moralSensit}(i,m) \wedge \text{des}_{i,k} \wedge \text{idl}_{i,l})$$

Our operator of preference corresponds to Cohen & Levesque's goal operator [7]. However, differently from DL-MA, their logic does not explain the relationships between goals and desires, and between goals and moral values.

Proposition 8 captures some basic properties of our preference operator.

**Proposition 8** *For every  $i \in \text{Agt}$  we have:*

$$\vdash (\text{Pref}_i\varphi \wedge \text{Pref}_i\psi) \rightarrow \text{Pref}_i(\varphi \wedge \psi) \quad (8a)$$

$$\vdash \neg(\text{Pref}_i\varphi \wedge \text{Pref}_i\neg\varphi) \quad (8b)$$

$$\text{if } \vdash \varphi \text{ then } \vdash \text{Pref}_i\varphi \quad (8c)$$

$$\vdash \text{Pref}_i\varphi \rightarrow K_i \text{Pref}_i\varphi \quad (8d)$$

$$\vdash \neg \text{Pref}_i\varphi \rightarrow K_i \neg \text{Pref}_i\varphi \quad (8e)$$

$$\vdash K_i\varphi \rightarrow \text{Pref}_i\varphi \quad (8f)$$

$$\vdash ((\neg\varphi <_i^{\text{des}} \varphi) \wedge \text{moralSensit}(i,0)) \rightarrow \text{Pref}_i\varphi \quad (8g)$$

$$\vdash ((\neg\varphi <_i^{\text{idl}} \varphi) \wedge \text{moralSensit}(i,\max)) \rightarrow \text{Pref}_i\varphi \quad (8h)$$

Theorems (8a) and (8b) together with the rule of inference (8c) highlights that  $\text{Pref}_i$  is a KD normal modal operator. Theorems (8d) and (8e) are properties of positive and negative introspection for preferences. Theorem (8f) is the syntactic counterpart of the *realism* principle for preferences discussed above. A similar property is satisfied by Cohen & Levesque's goal operator. Finally, theorems (8g) and (8h) highlight the logical relationship between desires and preferences, and moral values and preferences. According to (8g), if an agent  $i$  has a minimal degree of moral sensitivity and considers  $\varphi$  strictly more desirable than  $\neg\varphi$ , then  $i$  prefers  $\varphi$  to be true. According to (8h), if an agent  $i$  has a maximal degree of moral sensitivity and considers  $\varphi$  strictly more ideal than  $\neg\varphi$ , then  $i$  prefers  $\varphi$  to be true.

The last aspect of moral agency we consider is the relationship between preferences and choices. To this aim we need to introduce a notion of preference-based rationality.

We say that a given agent  $i$  is rational, denoted by  $\text{Rat}_i$ , if and only if for every action  $a$ , if  $i$  prefers to do  $a$  and knows that is able to do it, then he decides to do  $a$ ; and if he prefers not to do  $a$ , then he does not decide to do  $a$ . In other words, an agent is rational if his decision to perform a given action is completely determined by his own preferences:

$$\text{Rat}_i \stackrel{\text{def}}{=} \bigwedge_{a \in \text{Act}} (((\text{Pref}_i \text{choose}(a_i) \wedge \text{K}_i \text{able}(a_i)) \rightarrow \text{choose}(a_i)) \wedge (\text{Pref}_i \neg \text{choose}(a_i) \rightarrow \neg \text{choose}(a_i)))$$

The following two theorems, which follow from theorems (8g) and (8h) in Proposition 8, explain how desires and moral values motivate an agent to perform a given action:

$$\begin{aligned} & \vdash ((\neg \text{choose}(a_i) <_i^{\text{des}} \text{choose}(a_i)) \wedge \text{moralSensit}(i, 0) \wedge \text{Rat}_i \wedge \text{K}_i \text{able}(a_i)) \rightarrow \\ & \text{choose}(a_i) \\ & \vdash ((\neg \text{choose}(a_i) <_i^{\text{idl}} \text{choose}(a_i)) \wedge \text{moralSensit}(i, \text{max}) \wedge \text{Rat}_i \wedge \text{K}_i \text{able}(a_i)) \rightarrow \\ & \text{choose}(a_i) \end{aligned}$$

According to the preceding two theorems, if agent  $i$  is rational, has a minimal/maximal degree of moral sensitivity, considers the situation in which he performs action  $a$  strictly more desirable/ideal than the situation in which he does not perform  $a$  and knows that he is able to do  $a$ , then  $i$  decides to do  $a$ .

## 5 Conclusion

We have devised a logic which supports reasoning about actions and capabilities of agents and coalitions, epistemic states of agents as well as their desires and moral values. We have used it to provide a formal analysis of the concept of moral agency.

Directions of future work are manifold. For instance, there are important aspects of moral agency that have not been addressed in this work and that we intend to study in the future. One of them is the concept of moral emotion [11]. Moral emotions such as guilt, moral pride and reproach are emotions which are based either on the fulfillment or on the violation of an agent's moral values by the agent himself or by another agent. Another issue we plan to investigate, and which has been briefly mentioned in Section 2, is the relationships between an agent's moral values and external norms (*e.g.*, obligations, prohibitions, etc.). As to the logical part, we have provided a complete axiomatization of the logic DL-MA. Future work will be devoted to prove its decidability.

## 6 Acknowledgements

This research has been supported by the French ANR project EmoTES “Emotions in strategic interaction: theory, experiments, logical and computational studies”, contract No. 11-EMCO-004-01.

## References

- [1] I. Alger and J. W. Weibull. Homo moralis: preference evolution under incomplete information and assortative matching. Technical report, Toulouse School of Economics (TSE), 2012.
- [2] R.J. Aumann and J.H. Dreze. Rational expectations in games. *American Economic Review*, 98(1):72–86, 2008.
- [3] P. Battigalli and M. Dufwenberg. Guilt in games. *The American Economic Review*, 97(2):170–176, 2007.
- [4] N. Belnap, M. Perloff, and M. Xu. *Facing the future: agents and choices in our indeterminist world*. Oxford University Press, New York, 2001.
- [5] J. Broersen. A logical analysis of the interaction between ‘obligation-to-do’ and ‘knowingly doing’. In *Proc. of the Ninth International Conference on Deontic Logic in Computer Science (DEON’08)*, volume 5076 of *LNCS*, pages 140–154. Springer-Verlag, 2008.
- [6] J. Broersen, M. Dastani, J. Hulstijn, and L. van der Torre. Goal generation in the BOID architecture. *Cognitive Science Quarterly*, 2(3-4):428–447, 2002.
- [7] P. R. Cohen and H. J. Levesque. Intention is choice with commitment. *Artificial Intelligence*, 42:213–261, 1990.
- [8] E. Fehr and K. M. Schmidt. Theories of fairness and reciprocity: Evidence and economic applications. In *Advances in Economics and Econometrics*. Cambridge University Press, 2003.
- [9] H. Gintis, S. Bowles, R. Boyd, and E. Fehr, editors. *Moral sentiments and material interests*. MIT Press, Cambridge, 2005.
- [10] G. Governatori and A. Rotolo. BIO logical agents: Norms, beliefs, intentions in defeasible logic. *Journal of Autonomous Agents and Multi Agent Systems*, 17(1):36–69, 2008.
- [11] J. Haidt. The moral emotions. In R. J. Davidson, K. R. Scherer, and H. H. Goldsmith, editors, *Handbook of Affective Sciences*, pages 852–870. Oxford University Press, 2003.
- [12] J. Harsanyi. Morality and the theory of rational behaviour. In A. K. Sen and B. Williams, editors, *Utilitarianism and beyond*. Cambridge University Press, Cambridge, 1982.
- [13] A. Herzig and E. Lorini. A dynamic logic of agency I: STIT, abilities and powers. *Journal of Logic, Language and Information*, 19(1):89–121, 2010.
- [14] J. F. Horty. *Agency and Deontic Logic*. Oxford University Press, Oxford, 2001.
- [15] D. Lewis. *Counterfactuals*. Harvard University Press, 1973.

- [16] F. Liu. *Changing for the better: Preference dynamics and agent diversity*. PhD thesis, University of Amsterdam, The Netherlands, 2008.
- [17] R. Reiter. *Knowledge in action: logical foundations for specifying and implementing dynamical systems*. MIT Press, Cambridge, 2001.
- [18] O. Roy. Epistemic logic and the foundations of decision and game theory. *Journal of the Indian Council of Philosophical Research*, 27(2):283–314, 2010.
- [19] J. Searle. *Rationality in Action*. MIT Press, Cambridge, 2001.
- [20] A. K. Sen. Rational fools: a critique of the behavioral foundations of economic theory. *Philosophy and public affairs*, 6:317–344, 1977.
- [21] R. Thomason. Combinations of tense and modality. In D. Gabbay and F. Guentner, editors, *Handbook of Philosophical Logic*. Reidel, Dordrecht, 1984.
- [22] J. van Benthem, P. Girard, and O. Roy. Everything else being equal: a modal logic for *ceteris paribus* preferences. *Journal of Philosophical Logic*, 38(1):83–125, 2009.
- [23] A. Zanardo. Branching-time logic with quantification over branches: The point of view of modal logic. *Journal of Symbolic Logic*, 61(1):143–166, 1996.