



# A Qualitative Theory of Cognitive Attitudes and their Change

Emiliano Lorini

## ► To cite this version:

Emiliano Lorini. A Qualitative Theory of Cognitive Attitudes and their Change. Theory and Practice of Logic Programming, 2021, 21 (4), pp.428-458. 10.1017/S1471068421000053 . hal-03453908

**HAL Id: hal-03453908**

**<https://hal.science/hal-03453908>**

Submitted on 30 Nov 2021

**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.

# A Qualitative Theory of Cognitive Attitudes and their Change

EMILIANO LORINI  
IRIT-CNRS, Toulouse University, France

## Abstract

We present a general logical framework for reasoning about agents' cognitive attitudes of both epistemic type and motivational type. We show that it allows us to express a variety of relevant concepts for qualitative decision theory including the concepts of knowledge, belief, strong belief, conditional belief, desire, conditional desire, strong desire and preference. We also present two extensions of the logic, one by the notion of choice and the other by dynamic operators for belief change and desire change, and we apply the former to the analysis of single-stage games under incomplete information. We provide sound and complete axiomatizations for the basic logic and for its two extensions.

The paper is “under consideration in Theory and Practice of Logic Programming (TPLP)”.

## 1 Introduction

Since the seminal work of Hintikka on epistemic logic [28], of Von Wright on the logic of preference [55, 56] and of Cohen & Levesque on the logic of intention [19], many formal logics for reasoning about cognitive attitudes of agents such as knowledge and belief [24], preference [32, 48], desire [23], intention [44, 30] and their combination [38, 54] have been proposed. Generally speaking, these logics are nothing but formal models of rational agency relying on the idea that an agent endowed with cognitive attitudes makes decisions on the basis of what she believes and of what she desires or prefers.

The idea of describing rational agents in terms of their epistemic and motivational attitudes is something that these logics share with classical decision theory and game theory. Classical decision theory and game theory provide a quantitative account of individual and strategic decision-making by assuming that agents' beliefs and desires can be respectively modeled by subjective probabilities and utilities. Qualitative approaches to individual and strategic decision-making have been proposed in AI [16, 22] to characterize criteria that a rational agent should adopt for making decisions when she cannot build a probability distribution over the set of possible events and her preference over the set of possible outcomes cannot be expressed by a utility function but only by a qualitative ordering over the outcomes. For example, going beyond expected utility maximization, qualitative criteria such as the maxmin principle (choose the action that

will minimize potential loss) and the maxmax principle (choose the action that will maximize potential gain) have been studied and axiomatically characterized [18, 17].

The aim of this paper is to present an expressive logical framework for representing both the static and the dynamic aspects of a rich variety of agents' cognitive attitudes in a multi-agent setting. In agreement with philosophical theories [41, 43, 29, 34], our logic allows us to distinguish two general categories of cognitive attitudes: *epistemic* attitudes, including belief and knowledge, and *motivational* ones, including desire and preference. Moreover, in agreement with rational choice theory, it allows us to capture a notion of choice which depends on what an agent believes and prefers.<sup>1</sup>

The example depicted in Figure 1 brings to the fore the epistemic and motivational attitudes that are involved in everyday situations whereby artificial agents are supposed to interact. There are two autonomous agents meeting at a crossroad: agent 1 and agent 2. The two agents could be either two mobile robots or two autonomous vehicles. Each agent can decide either to stop or to continue. If an agent stops, then it will lose time. If both agents decide to continue, they will collide and, consequently, each of them will lose time. Therefore, for an agent not to lose time, it has to continue, while the other agent decides to stop.

In this situation, each agent is identified with the set of cognitive attitudes it endorses. For instance, it is reasonable to suppose that the two agents know that in the situation they face necessarily some of them will lose time and that if one of them loses time by letting the other pass, there will be no collision. On the motivational side, it is reasonable to suppose that each agent is strongly motivated by two desires, namely, the desire not to lose time and the desire to avoid a collision.

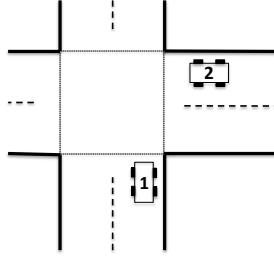


Figure 1: Crossroad game

On the dynamic side, we consider two basic forms of cognitive attitude change, namely, belief change and desire change. While belief change has been extensively studied in the area of belief revision [1, 15, 20, 46, 42, 13] and dynamic epistemic logic (DEL) [47, 10, 51, 4], desire change is far less studied and understood. We will study two basic forms of cognitive attitude revision, namely, *radical* attitude revision and *conservative* attitude revision. While the distinction between radical and conservative belief revision has been drawn before (see, e.g., [47]), the distinction between radical

<sup>1</sup> Rational choice theory (RCT) is a umbrella term for a family of theories prescribing that an agent should choose the course of action that, according to her beliefs, leads to the most desirable (or most preferred) consequences. In other words, RCT relies on the general assumption that agents make optimal choices in the light of her beliefs, desires and preferences. See [40] for more details on RCT.

and conservative desire revision is new. Radical belief revision by an input  $\varphi$  makes all states at which  $\varphi$  is true more plausible than all states at which  $\varphi$  is false, whereas conservative belief revision by  $\varphi$  simply promotes the most plausible states in which  $\varphi$  is true to the highest plausibility rank, but apart from that, it keeps the old plausibility ordering. For example, suppose in the crossroad game of Figure 1, agent 1 and agent 2 can communicate. Agent 1 informs agent 2 that “if they both lose time, then there will no collision” and agent 2 trusts what agent 1 says. Then, by performing a conservative belief revision, agent 2 will promote the most plausible situations in which the formula announced by 1 is true to the highest plausibility rank. As a consequence, agent 2 will start to believe what 1 has just said.

Symmetrically, radical desire revision by  $\varphi$  makes all states at which  $\varphi$  is true more desirable than all states at which  $\varphi$  is false, whereas conservative desire revision by  $\varphi$  simply demotes the least desirable states in which  $\varphi$  is false to the lowest desirability rank, but apart from that, it keeps the old desirability ordering. For example, suppose in the crossroad game agent 1 has just learnt that agent 2 is an ambulance which has to transport a patient to the hospital as quickly as possible. Consequently, 1 starts to be altruistically motivated by the fact that 2 does not lose time. Thus, by performing a radical desire revision, agent 1 will start to consider all situations in which 2 does not lose time more desirable than the situations in which it does. This radical desire revision operation leads agent 1 to strongly desire that agent 2 does not lose time.

The paper is organized as follows. In Section 2, we present the semantics and syntax of our logic, called Dynamic Logic of Cognitive Attitudes (DLCA). At the semantic level, it exploits two orderings that capture, respectively, an agent’s comparative plausibility and comparative desirability over states. At the syntactic level, it uses program constructs of dynamic logic (sequential composition, non-deterministic choice, intersection, complement, converse and test) to build complex cognitive attitudes from simple ones. Following [39, 25], it also exploits nominals in order to axiomatize intersection and complement of programs. In Section 3, we illustrate the expressive power of our logic by using it to formalize a variety of cognitive attitudes of agents including knowledge, belief, strong belief, conditional belief, desire, strong desire, conditional desire and preference. We instantiate some of these concepts in the crossroad game depicted in Figure 1. In Section 4, we present a sound and complete axiomatization for our logic. In Section 5 we present the first extension of our logic by the notion of choice and apply it to the analysis of single-stage games under incomplete information. Section 6 presents the second extension of our logic by dynamic operators for belief and desire change. In Section 7 we conclude. Formal proofs are given in a technical annex at the end of the paper.<sup>2</sup>

---

<sup>2</sup>This paper is an extended and improved version of [35]. The JELIA’19 paper did not include the two extensions of Section 5 and Section 6, or the detailed proof of the completeness theorem for the logic DLCA. Also, the logical analysis of the cognitive attitudes in Section 3 has been extended: (i) we included the notion of conditional desire which was not considered in the JELIA’19 paper, and (ii) we added new logical validities which describe interesting properties of cognitive attitudes.

## 2 Dynamic Logic of Cognitive Attitudes

Let  $Atm$  be a countable infinite set of atomic propositions, let  $Nom$  be a countable infinite set of nominals disjoint from  $Atm$  and let  $Agt$  be a finite set of agents.

**Definition 1 (Multi-agent cognitive model)** A multi-agent cognitive model (MCM) is a tuple  $M = (W, (\preceq_{i,P})_{i \in Agt}, (\preceq_{i,D})_{i \in Agt}, (\equiv_i)_{i \in Agt}, V)$  where:

- $W$  is a set of worlds or states;
- for every  $i \in Agt$ ,  $\preceq_{i,P}$  and  $\preceq_{i,D}$  are preorders on  $W$  and  $\equiv_i$  is an equivalence relation on  $W$  such that for all  $\tau \in \{P, D\}$  and for all  $w, v \in W$ :

$$(C1) \preceq_{i,\tau} \subseteq \equiv_i,$$

$$(C2) \text{ if } w \equiv_i v \text{ then } w \preceq_{i,\tau} v \text{ or } v \preceq_{i,\tau} w;$$

- $V : W \rightarrow 2^{Atm \cup Nom}$  is a valuation function such that for all  $w, v \in W$ :

$$(C3) V_{Nom}(w) \neq \emptyset,$$

$$(C4) \text{ if } V_{Nom}(w) \cap V_{Nom}(v) \neq \emptyset \text{ then } w = v;$$

$$\text{where } V_{Nom}(w) = Nom \cap V(w).$$

$w \preceq_{i,P} v$  means that, according to agent  $i$ ,  $v$  is at least as plausible as  $w$ , whereas  $w \preceq_{i,D} v$  means that, according to agent  $i$ ,  $v$  is at least as desirable as  $w$ . Finally,  $w \equiv_i v$  means that  $w$  and  $v$  are indistinguishable for agent  $i$ . For every  $w \in W$ ,  $\equiv_i(w)$  is also called agent  $i$ 's information set at state  $w$ . According to Constraint C1, an agent can only compare the plausibility (resp. desirability) of two states in her information set. According to Constraint C2, the plausibility (resp. desirability) of two states in an agent's information set are always comparable. Constraints C3 and C4 capture the two basic properties of nominals: every state is associated with at least one nominal and there are no different states associated with the same nominal.

Note that there is no connection between binary relations  $\preceq_{i,P}$  and  $\preceq_{i,D}$ . In accord with classical decision and game theory in which an agent's subjective probability and utility function do not interact, we adopt a normative view of epistemic and motivational attitudes according to which an agent's epistemic plausibility and desirability are assumed to be independent.<sup>3</sup> Therefore, we do not consider cognitive biases typical of human reasoning such as wishful thinking, as the tendency to form beliefs according to what is desired in the absence of a clear evidence against it [37]. Nonetheless, as we will show in Section 3.3, the primitive relations  $\preceq_{i,P}$  and  $\preceq_{i,D}$  can be combined to obtain a notion of realistic preference which is essential for elucidating the connection between an agent's beliefs and desires and her choices.

We introduce the following modal language  $\mathcal{L}_{DLCA}(Atm, Nom, Agt)$ , or simply  $\mathcal{L}_{DLCA}$ , for the Dynamic Logic of Cognitive Attitudes DLCA:

<sup>3</sup>The normative view is usually opposed to the descriptive view. The normative view is aimed at describing the reasoning and decision-making of ideal agents conforming to standards of rationality, while the descriptive view is concerned with psychologically realistic cognitive agents who systematically violate standards of rationality and exhibit different types of cognitive bias.

$$\begin{aligned}
\pi &::= \equiv_i | \preceq_{i,P} | \preceq_{i,D} | \preceq_{i,P}^\sim | \preceq_{i,D}^\sim | \pi; \pi' | \pi \cup \pi' | \pi \cap \pi' | \neg \pi | \varphi? \\
\varphi &::= p | x | \neg \varphi | \varphi \wedge \varphi' | [\pi] \varphi
\end{aligned}$$

where  $p$  ranges over  $Atm$ ,  $x$  ranges over  $Nom$  and  $i$  ranges over  $Agt$ . The other Boolean constructions  $\top$ ,  $\perp$ ,  $\vee$ ,  $\rightarrow$  and  $\leftrightarrow$  are defined from  $p$ ,  $\neg$  and  $\wedge$  in the standard way. The propositional language built from the set of atomic propositions  $Atm$  is noted  $\mathcal{L}_{PL}(Atm)$ . Note that the sets  $Atm$ ,  $Nom$  and  $Agt$  define the signature of the language  $\mathcal{L}_{DLCA}$ . They are not part of the model since every atomic proposition  $p$ , nominal  $x$  and modal formula  $[\pi]\varphi$  should be interpretable relative to any MCM.

Elements  $\pi$  are called *cognitive programs* or, more shortly, *programs*. The set of all programs is noted  $\mathcal{P}(Atm, Nom, Agt)$ , or simply,  $\mathcal{P}$ .

Cognitive programs correspond to the basic constructions of Propositional Dynamic Logic (PDL) [26]: atomic programs of type  $\equiv_i$ ,  $\preceq_{i,P}$ ,  $\preceq_{i,D}$ ,  $\preceq_{i,P}^\sim$  and  $\preceq_{i,D}^\sim$ , sequential composition ( $;$ ), non-deterministic choice ( $\cup$ ), intersection ( $\cap$ ), converse ( $\neg$ ) and test ( $?$ ). A given cognitive program  $\pi$  corresponds to a specific configuration of the agents' cognitive states including their epistemic states and their motivational states.

The formula  $[\pi]\varphi$  has to be read “ $\varphi$  is true, according to the cognitive program  $\pi$ ”. As usual, we define  $\langle \pi \rangle$  to be the dual operator of  $[\pi]$ , that is,  $\langle \pi \rangle \varphi =_{def} \neg[\pi]\neg\varphi$ .

The atomic program  $\equiv_i$  represents the standard S5, partition-based and fully introspective notion of knowledge [24, 5].  $[\equiv_i]\varphi$  has to be read “ $\varphi$  is true according to what agent  $i$  knows” or more simply “agent  $i$  knows that  $\varphi$  is true”, which just means that “ $\varphi$  is true in all worlds that agent  $i$  envisages”.

The atomic programs  $\preceq_{i,P}$  and  $\preceq_{i,D}$  capture, respectively, agent  $i$ 's plausibility ordering and agent  $i$ 's desirability ordering over facts. In particular,  $[\preceq_{i,P}]\varphi$  has to be read “ $\varphi$  is true at all states that, according to agent  $i$ , are at least as plausible as the current one”, while  $[\preceq_{i,D}]\varphi$  has to be read “ $\varphi$  is true at all states that, according to agent  $i$ , are at least as desirable as the current one”. The atomic programs  $\preceq_{i,P}^\sim$  and  $\preceq_{i,D}^\sim$  are the complements of the atomic programs  $\preceq_{i,P}$  and  $\preceq_{i,D}$ , respectively. In particular,  $[\preceq_{i,P}^\sim]\varphi$  has to be read “ $\varphi$  is true at all states that, according to agent  $i$ , are *not* at least as plausible as the current one”, while  $[\preceq_{i,D}^\sim]\varphi$  has to be read “ $\varphi$  is true at all states that, according to agent  $i$ , are *not* at least as desirable as the current one”. The program constructs  $;$ ,  $\cup$ ,  $\cap$ ,  $\neg$  and  $?$  are used to define complex cognitive programs from the atomic cognitive programs. For example, the formula  $[\preceq_{i,P} \cup \preceq_{i,D}]\varphi$  has to be read “ $\varphi$  is true at all states that, according to agent  $i$ , are either at least as plausible *or* at least as desirable as the current one”, whereas the formula  $[\preceq_{i,P} \cap \preceq_{i,D}]\varphi$  has to be read “ $\varphi$  is true at all states that, according to agent  $i$ , are at least as plausible *and* at least as desirable as the current one”.

The following definition provides truth conditions for formulas in  $\mathcal{L}_{DLCA}$ :

**Definition 2 (Truth conditions)** *Let  $M = (W, (\preceq_{i,P})_{i \in Agt}, (\preceq_{i,D})_{i \in Agt}, (\equiv_i)_{i \in Agt}, V)$  be*

a MCM and let  $w \in W$ . Then:

$$\begin{aligned}
M, w \models p &\iff p \in V(w), \\
M, w \models x &\iff x \in V(w), \\
M, w \models \neg\phi &\iff M, w \not\models \phi, \\
M, w \models \phi \wedge \psi &\iff M, w \models \phi \text{ and } M, w \models \psi, \\
M, w \models [\pi]\phi &\iff \forall v \in W : \text{if } wR_\pi v \text{ then } M, v \models \phi,
\end{aligned}$$

where the binary relation  $R_\pi$  on  $W$  is inductively defined as follows, with  $\tau \in \{P, D\}$ :

$$\begin{aligned}
wR_{\equiv_i} v &\text{ iff } w \equiv_i v, \\
wR_{\succeq_{i,\tau}} v &\text{ iff } w \succeq_{i,\tau} v, \\
wR_{\succeq_{i,\tau}^\sim} v &\text{ iff } w \equiv_i v \text{ and } w \not\prec_{i,\tau} v, \\
wR_{\pi;\pi'} v &\text{ iff } \exists u \in W : wR_\pi u \text{ and } uR_{\pi'} v, \\
wR_{\pi \cup \pi'} v &\text{ iff } wR_\pi v \text{ or } wR_{\pi'} v, \\
wR_{\pi \cap \pi'} v &\text{ iff } wR_\pi v \text{ and } wR_{\pi'} v, \\
wR_{-\pi} v &\text{ iff } vR_\pi w, \\
wR_{\phi?} v &\text{ iff } w = v \text{ and } M, w \models \phi.
\end{aligned}$$

For notational convenience, we use  $wR_\pi v$  and  $(w, v) \in R_\pi$  as interchangeable notations.

We can build a variety of cognitive programs capturing different types of plausibility and desirability relations between possible worlds. For instance, for every  $\tau \in \{P, D\}$ , we can define:

$$\begin{aligned}
\succeq_{i,\tau} &=_{\text{def}} - \prec_{i,\tau}, \\
\succ_{i,\tau} &=_{\text{def}} \succeq_{i,\tau} \cap \succeq_{i,\tau}^\sim, \\
\succeq_{i,\tau}^\sim &=_{\text{def}} - \prec_{i,\tau}^\sim, \\
\prec_{i,\tau} &=_{\text{def}} \succeq_{i,\tau} \cap \succeq_{i,\tau}^\sim, \\
\approx_{i,\tau} &=_{\text{def}} \succeq_{i,\tau} \cap \succeq_{i,\tau}^\sim.
\end{aligned}$$

The five definitions denote respectively “at most as plausible (resp. desirable) as”, “less plausible (resp. desirable) than”, “not at most as plausible (resp. desirable) as”, “more plausible (resp. desirable) than” and “equally plausible (resp. desirable) as”.

For every formula  $\phi$  in  $\mathcal{L}_{\text{DLCA}}$  we say that  $\phi$  is valid, noted  $\models_{\text{MCM}} \phi$ , if and only if for every multi-agent cognitive model  $M$  and world  $w$  in  $M$ , we have  $M, w \models \phi$ . Conversely, we say that  $\phi$  is satisfiable if  $\neg\phi$  is not valid.

For a given multi-agent cognitive model  $M = (W, (\succeq_{i,P})_{i \in \text{Agt}}, (\succeq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, N, V)$ , we define  $\|\phi\|_M = \{v \in W : M, v \models \phi\}$  to be the truth set of  $\phi$  in  $M$ . Moreover, for every  $w \in W$  and for every  $i \in \text{Agt}$ , we define  $\|\phi\|_{i,w,M} = \{v \in W : M, v \models \phi \text{ and } w \equiv_i v\}$  to be the truth set of  $\phi$  from  $i$ 's point of view at state  $w$  in  $M$ .

### 3 Formalization of Cognitive Attitudes

In this section, we show how the logic DLCA can be used to model the variety of cognitive attitudes of agents that we have briefly discussed in the introduction.

#### 3.1 Epistemic Attitudes

We start with the family of epistemic attitudes by defining a standard notion of belief. We say that an agent believes that  $\phi$  if and only if  $\phi$  is true at all states that the agent considers maximally plausible.

**Definition 3 (Belief)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that agent  $i$  believes that  $\phi$  at  $w$ , noted  $M, w \models B_i \phi$ , if and only if  $\text{Best}_{i,P}(w) \subseteq \|\phi\|_M$  where  $\text{Best}_{i,P}(w) = \{v \in W : w \equiv_i v \text{ and } \forall u \in W, \text{ if } w \equiv_i u \text{ then } u \preceq_{i,P} v\}$ .

As the following proposition highlights, the previous notion of belief is expressible in the logic DLCA by means of the cognitive program  $\equiv_i; [\prec_{i,P}] \perp ?$ .

**Proposition 1** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models B_i \phi \text{ iff } M, w \models [\equiv_i; [\prec_{i,P}] \perp ?] \phi.$$

It is worth noting that the set  $\text{Best}_{i,P}(w)$  in Definition 3 might be empty, since it is not necessarily the case that the relation  $\preceq_{i,P}$  is conversely well-founded.<sup>4</sup> As a consequence, the belief operator  $B_i$  does not necessarily satisfy Axiom D, i.e., the formula  $B_i \phi \wedge B_i \neg \phi$  is satisfiable in the logic DLCA.

In the literature on epistemic logic [11], mere belief of Definition 3 is usually distinguished from strong belief. Specifically, we say that an agent strongly believes that  $\phi$  if and only if, according to agent  $i$ , all  $\phi$ -worlds are strictly more plausible than all  $\neg \phi$ -worlds.

**Definition 4 (Strong belief)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that agent  $i$  strongly believes that  $\phi$  at  $w$ , noted  $M, w \models SB_i \phi$ , if and only if  $\forall v \in \|\phi\|_{i,w,M}$  and  $\forall u \in \|\neg \phi\|_{i,w,M} : u \prec_{i,P} v$ .

As the following proposition highlights, the previous notion of strong belief is expressible in the logic DLCA by means of the cognitive program  $\equiv_i; \phi?; \preceq_{i,P}$ .

**Proposition 2** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models SB_i \phi \text{ iff } M, w \models [\equiv_i; \phi?; \preceq_{i,P}] \phi.$$

<sup>4</sup>This means that there could be a world  $v$  such that  $w \equiv_i v$  and there is a  $\preceq_{i,P}$ -infinite ascending chain from  $v$ .



Strong belief that  $\phi$  implies belief that  $\phi$ , if the agent envisages at least one state in which  $\phi$  is true. This property is expressed by the following validity:

$$\models_{MCM} (\text{SB}_i \phi \wedge \langle \equiv_i \rangle \phi) \rightarrow \text{B}_i \phi \quad (1)$$

Conditional belief is another notion which has been studied by epistemic logicians given its important role in belief dynamics [47]. We say that an agent believes that  $\phi$  conditional on  $\psi$ , or she would believe that  $\phi$  if she learnt that  $\psi$ , if and only if, according to the agent, all most plausible  $\psi$ -worlds are also  $\phi$ -worlds.

**Definition 5 (Conditional belief)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that agent  $i$  would believe that  $\phi$  if she learnt that  $\psi$  at  $w$ , noted  $M, w \models \text{B}_i(\psi, \phi)$ , if and only if  $\text{Best}_{i,P}(\psi, w) \subseteq \|\phi\|_M$ , where  $\text{Best}_{i,P}(\psi, w) = \{v \in \|\psi\|_{i,w,M} : \forall u \in \|\psi\|_{i,w,M}, u \preceq_{i,P} v\}$ .

Note that  $\text{Best}_{i,P}(\top, w) = \text{Best}_{i,P}(w)$ .

As for belief and strong belief, we have a specific cognitive program  $\equiv_i; (\psi \wedge [\prec_{i,P}] \neg \psi)?$  corresponding to the belief that  $\phi$  conditional on  $\psi$ , so that the latter can be represented in the language of the logic DLCA.

**Proposition 3** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models \text{B}_i(\psi, \phi) \text{ iff } M, w \models [\equiv_i; (\psi \wedge [\prec_{i,P}] \neg \psi)?] \phi.$$

### 3.2 Motivational Attitudes I: Desires

The first kind of motivational attitude we consider is desire. Following [23], we say that an agent desires that  $\phi$  if and only if all states that the agent envisages at which  $\phi$  is true are not minimally desirable for her. In other words, desiring that  $\phi$  consists in having some degree of attraction for all situations in which  $\phi$  is true, since minimally desirable states are those to which the agent is not attracted at all.

**Definition 6 (Desire)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that agent  $i$  desires that  $\phi$  at  $w$ , noted  $M, w \models \text{D}_i \phi$ , if and only if  $\text{Worst}_{i,D}(w) \cap \|\phi\|_M = \emptyset$ , where  $\text{Worst}_{i,D}(w) = \{v \in W : w \equiv_i v \text{ and } \forall u \in W, \text{ if } w \equiv_i u \text{ then } v \preceq_{i,D} u\}$ .

As the following proposition highlights, the previous notion of desire is characterized by the cognitive program  $\equiv_i; [\succ_{i,D}] \perp ?$ .

**Proposition 4** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models \text{D}_i \phi \text{ iff } M, w \models [\equiv_i; [\succ_{i,D}] \perp ?] \neg \phi.$$

Similarly to the set  $Best_{i,P}(w)$  in Definition 3, the set  $Worst_{i,D}(w)$  in Definition 6 might be empty, since it is not necessarily the case that the relation  $\preceq_{i,D}$  is well-founded.<sup>5</sup> As a consequence, desires are not necessarily consistent and an agent may desire the tautology, i.e., the formulas  $D_i\varphi \wedge D_i\neg\varphi$  and  $D_i\top$  are satisfiable in the logic DLCA. As emphasized by [23], this notion of desire satisfies the following property:

$$\models_{MCM} D_i\varphi \rightarrow D_i(\varphi \wedge \psi) \quad (2)$$

Indeed, if an agent has some degree of attraction for all situations in which  $\varphi$  is true then, clearly, it should have some degree of attraction for all situations in which  $\varphi \wedge \psi$  is true, since all  $\varphi \wedge \psi$ -situations are also  $\varphi$ -situations.

Note that there is no counterpart of this property for belief, as the formula  $B_i\varphi \rightarrow B_i(\varphi \wedge \psi)$  is clearly not valid.<sup>6</sup>

It is a property that the notion of desire shares with the *open reading* of the concept of permission studied in the area of deontic logic (see, e.g., [3, 31]).<sup>7</sup> One way of blocking this inference is by strengthening the notion of desire. We say that an agent strongly desires that  $\varphi$  if and only if, according to agent  $i$ , all  $\varphi$ -worlds are strictly more desirable than all  $\neg\varphi$ -worlds.

**Definition 7 (Strong desire)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that agent  $i$  strongly desires that  $\varphi$  at  $w$ , noted  $M, w \models SD_i\varphi$ , if and only if  $\forall v \in \|\varphi\|_{i,w,M}$  and  $\forall u \in \|\neg\varphi\|_{i,w,M} : u \prec_{i,D} v$ .

As for desire, there exists a cognitive program which characterizes strong desire, namely, the program  $\equiv_i; \varphi?; \preceq_{i,D}$ .

**Proposition 5** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models SD_i\varphi \text{ iff } M, w \models [\equiv_i; \varphi?; \preceq_{i,D}] \varphi.$$

We have that strong desire implies desire, when the agent envisages at least one state in which  $\varphi$  is false:

$$\models_{MCM} (SD_i\varphi \wedge \langle \equiv_i \rangle \neg\varphi) \rightarrow D_i\varphi \quad (3)$$

Unlike desire, it is not necessarily the case that strongly desiring that  $\varphi$  implies strongly desiring that  $\varphi \wedge \psi$ , i.e.,  $SD_i\varphi \wedge \neg SD_i(\varphi \wedge \psi)$  is satisfiable in the logic DLCA. Indeed, strongly desiring that  $\varphi$  is compatible with envisaging a situation in which  $\varphi \wedge \psi$  holds and another situation in which  $\varphi \wedge \neg\psi$  holds such that the first situation is less desirable than the second.

<sup>5</sup>This means that there could be a world  $v$  such that  $w \equiv_i v$  and there is a  $\preceq_{i,D}$ -infinite descending chain from  $v$ .

<sup>6</sup>See [23] for more details about the differences between the notion of belief and the notion of desire.

<sup>7</sup> According to deontic logicians, there are at least two candidate readings of the statement “ $\varphi$  is permitted”: (i) every instance of  $\varphi$  is OK according to the normative regulation, and (ii) at least one instance of  $\varphi$  (but possibly not all) is OK according to the normative regulation. The former is the so-called *open reading* of permission.

The last motivational attitude we consider is conditional desire which parallels the notion of conditional belief of Definition 5. We say that an agent desires that  $\varphi$  conditional on  $\psi$ , or she would desire that  $\varphi$  if she started to desire that  $\psi$ , if and only if, according to agent  $i$ , there is no least desirable  $\neg\psi$ -world which is also a  $\varphi$ -world. The idea behind this notion is the following. If the agent started to desire that  $\psi$ , all  $\psi$ -worlds would start to have some degree of attraction for her and the least desirable  $\neg\psi$ -worlds would become the minimally desirable worlds. Therefore, the fact that there is no least desirable  $\neg\psi$ -world which is also a  $\varphi$ -world guarantees that, if the agent started to desire that  $\psi$ , no  $\varphi$ -world would be included in the set of minimally desirable worlds for the agent. The latter means that, if the agent started to desire that  $\psi$ , all  $\varphi$ -worlds would have some degree of attraction for her and she would desire that  $\varphi$ .

**Definition 8 (Conditional desire)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that agent  $i$  would desire that  $\varphi$  if she started to desire that  $\psi$  at  $w$ , noted  $M, w \models D_i(\psi, \varphi)$ , if and only if  $\text{Worst}_{i,D}(\neg\psi, w) \cap \|\varphi\|_M = \emptyset$ , with  $\text{Worst}_{i,D}(\neg\psi, w) = \{v \in \|\neg\psi\|_{i,w,M} : \forall u \in \|\neg\psi\|_{i,w,M}, v \preceq_{i,D} u\}$ .

As for the other cognitive attitudes, there is a specific cognitive program which characterizes conditional desire.

**Proposition 6** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models D_i(\psi, \varphi) \text{ iff } M, w \models [\equiv_i; (\neg\psi \wedge [\succ_{i,D}] \psi)?] \neg\varphi.$$

In Section 2, we emphasized that the relations  $\preceq_{i,P}$  and  $\preceq_{i,D}$  do not interact since our logic is aimed at modeling ideal rational agents with no wishful thinking and, more generally, with no cognitive biases. We conclude this section by showing how the assumption of independence between epistemic plausibility and desirability could be relaxed and, consequently, how wishful thinking could be modeled in our framework.

A wishful thinker is nothing but an agent who systematically believes what she strongly desires in the absence of a reason to believe the contrary. Such a connection between the agent's beliefs and desires is captured by the following “wishful thinking” (WT) constraint on MCMs:

$$\forall w \in W : \text{Best}_{i,P}(w) \subseteq \text{Best}_{i,D}(w) \text{ or } \text{Best}_{i,P}(w) \subseteq \text{Worst}_{i,D}(w),$$

where  $\text{Best}_{i,P}(w)$  and  $\text{Worst}_{i,D}(w)$  are defined as in Definitions 3 and 6, and  $\text{Best}_{i,D}(w) = \{v \in W : w \equiv_i v \text{ and } \forall u \in W, \text{ if } w \equiv_i u \text{ then } u \preceq_{i,D} v\}$ . It is routine to verify that if the MCM  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  satisfies the previous constraint WT, then the following holds for every  $w \in W$ :

$$M, w \models (\text{SD}_i \varphi \wedge \neg \text{B}_i \neg \varphi) \rightarrow \text{B}_i \varphi.$$

We leave for future work an in-depth analysis of the variant of our logic in which wishful thinking is enabled.

### 3.3 Motivational Attitudes II: Preferences

We consider two views about comparative statements between formulas of the form “agent  $i$  prefers  $\phi$  to  $\psi$ ” or “the state of affairs  $\phi$  is for agent  $i$  at least as good as the state of affairs  $\psi$ ”. According to the optimistic view, when assessing whether  $\phi$  is at least as good as  $\psi$ , an agent focuses on the best  $\phi$ -situations in comparison with the best  $\psi$ -situations. Specifically, an “optimistic” agent  $i$  prefers  $\phi$  to  $\psi$  if and only if, for every  $\psi$ -situation envisaged by  $i$  there exists a  $\phi$ -situation envisaged by  $i$  such that the latter is at least as desirable as the former. According to the pessimistic view, she focuses on the worst  $\phi$ -situations in comparison with the worst  $\psi$ -situations. Specifically, a “pessimistic” agent  $i$  prefers  $\phi$  to  $\psi$  if and only if, for every  $\phi$ -situation envisaged by  $i$  there exists a  $\psi$ -situation envisaged by  $i$  such that the former is at least as desirable as the latter.

Let us first define a dyadic operator for preference according to the optimistic view.

**Definition 9 (Preference: optimistic view)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that, according to agent  $i$ 's optimistic assessment,  $\phi$  is at least as good as  $\psi$  at  $w$ , noted  $M, w \models P_i^{\text{Opt}}(\psi \preceq \phi)$ , if and only if  $\forall u \in \|\psi\|_{i,w,M}, \exists v \in \|\phi\|_{i,w,M} : u \preceq_{i,D} v$ .

As the following proposition highlights, it is expressible in the language  $\mathcal{L}_{\text{DLCA}}$ .

**Proposition 7** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models P_i^{\text{Opt}}(\psi \preceq \phi) \text{ iff } M, w \models [\equiv_i; \psi?] \langle \preceq_{i,D} \rangle \phi.$$

Let us now define preference according to the pessimistic view.

**Definition 10 (Preference: pessimistic view)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that, according to agent  $i$ 's pessimistic assessment,  $\phi$  is at least as good as  $\psi$  at  $w$ , noted  $M, w \models P_i^{\text{Pess}}(\psi \preceq \phi)$ , if and only if  $\forall v \in \|\phi\|_{i,w,M}, \exists u \in \|\psi\|_{i,w,M} : u \preceq_{i,D} v$ .

As for the optimistic view, the pessimistic view is also expressible in the language  $\mathcal{L}_{\text{DLCA}}$ .

**Proposition 8** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models P_i^{\text{Pess}}(\psi \preceq \phi) \text{ iff } M, w \models [\equiv_i; \phi?] \langle \succeq_{i,D} \rangle \psi.$$

Thanks to the totality of the relation  $\preceq_{i,D}$  (Constraint C2 in Definition 1), dyadic preference over formulas is total too. This fact is illustrated by the following validity. For every  $x \in \{\text{Opt}, \text{Pess}\}$ :

$$\models_{\text{MCM}} P_i^x(\psi \preceq \phi) \vee P_i^x(\phi \preceq \psi) \quad (4)$$

To see this suppose  $M, w \models \neg P_i^{\text{Opt}}(\psi \preceq \phi)$  for an arbitrary model  $M$  and world  $w$  in  $M$ . Because of Constraint C2 in Definition 1, the latter implies that  $\exists u \in \|\psi\|_{i,w,M}, \forall v \in$

$\|\varphi\|_{i,w,M} : v \prec_{i,D} u$ . Therefore,  $\forall v \in \|\varphi\|_{i,w,M}, \exists u \in \|\psi\|_{i,w,M} : v \preceq_{i,D} u$  which is equivalent to  $M, w \models \text{RP}_i^{\text{Opt}}(\varphi \preceq \psi)$ . The case  $x = \text{Pess}$  can be proved in an analogous way.

The previous notion of (optimistic and pessimistic) preference does not depend on what the agent believes. This means that, in order to assess whether  $\varphi$  is at least as good as  $\psi$ , an agent also takes into account worlds that are implausible (or, more generally, not maximally plausible). *Realistic* preference requires that an agent compares two formulas  $\varphi$  and  $\psi$  only with respect to the set of most plausible states. This idea has been discussed in the area of qualitative decision theory by different authors [16, 18, 17].

The following definition introduces *realistic* preference according to the optimistic view.

**Definition 11 (Realistic preference: optimistic view)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that, according to agent  $i$ 's optimistic assessment,  $\varphi$  is realistically at least as good as  $\psi$  at  $w$ , noted  $M, w \models \text{RP}_i^{\text{Opt}}(\psi \preceq \varphi)$ , if and only if  $\forall u \in \text{Best}_{i,P}(w) \cap \|\psi\|_{i,w,M}, \exists v \in \text{Best}_{i,P}(w) \cap \|\varphi\|_{i,w,M} : u \preceq_{i,D} v$ .

The idea is that an “optimistic” agent  $i$  considers  $\varphi$  *realistically* at least as good as  $\psi$  if and only if, for every  $\psi$ -situation in agent  $i$ 's belief set there exists a  $\varphi$ -situation in agent  $i$ 's belief set such that the latter is at least as good as the former.

The previous notion as well is expressible in the language  $\mathcal{L}_{\text{DLCA}}$ .

**Proposition 9** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models \text{RP}_i^{\text{Opt}}(\psi \preceq \varphi) \text{ iff } M, w \models [\equiv_i; [\prec_{i,P}] \perp ?; \psi?] \langle \preceq_{i,D} \cap (\equiv_i; [\prec_{i,P}] \perp ?) \rangle \varphi.$$

The following definition introduces *realistic* preference according to the pessimistic view.

**Definition 12 (Realistic preference: pessimistic view)** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . We say that, according to agent  $i$ 's pessimistic assessment,  $\varphi$  is realistically at least as good as  $\psi$  at  $w$ , noted  $M, w \models \text{RP}_i^{\text{Pess}}(\psi \preceq \varphi)$ , if and only if  $\forall v \in \text{Best}_{i,P}(w) \cap \|\varphi\|_{i,w,M}, \exists u \in \text{Best}_{i,P}(w) \cap \|\psi\|_{i,w,M} : u \preceq_{i,D} v$ .

The idea is that a “pessimistic” agent  $i$  considers  $\varphi$  *realistically* at least as good as  $\psi$  if and only if, for every  $\varphi$ -situation in agent  $i$ 's belief set there exists a  $\psi$ -situation in agent  $i$ 's belief set such that the former is at least as good as the latter.

It is also expressible in the language  $\mathcal{L}_{\text{DLCA}}$ .

**Proposition 10** Let  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and let  $w \in W$ . Then, we have

$$M, w \models \text{RP}_i^{\text{Pess}}(\psi \preceq \varphi) \text{ iff } M, w \models [\equiv_i; [\prec_{i,P}] \perp ?; \varphi?] \langle \succeq_{i,D} \cap (\equiv_i; [\prec_{i,P}] \perp ?) \rangle \psi.$$

Like dyadic preference over formulas, realistic dyadic preference over formulas is total. In fact, for every for  $x \in \{Opt, Pess\}$ , we have:

$$\models_{MCM} RP_i^x(\psi \preceq \varphi) \vee RP_i^x(\varphi \preceq \psi) \quad (5)$$

The following abbreviations define *strict* variants of dyadic preference operators:

$$\begin{aligned} P_i^{Opt}(\psi \prec \varphi) &=_{def} \neg P_i^{Opt}(\varphi \preceq \psi) \\ P_i^{Pess}(\psi \prec \varphi) &=_{def} \neg P_i^{Pess}(\varphi \preceq \psi) \\ RP_i^{Opt}(\psi \prec \varphi) &=_{def} \neg RP_i^{Opt}(\varphi \preceq \psi) \\ RP_i^{Pess}(\psi \prec \varphi) &=_{def} \neg RP_i^{Pess}(\varphi \preceq \psi) \end{aligned}$$

$P_i^{Opt}(\psi \prec \varphi)$  (resp.  $P_i^{Pess}(\psi \prec \varphi)$ ) has to be read “according to  $i$ ’s optimistic (resp. pessimistic) assessment,  $\varphi$  is better than  $\psi$ ”.  $RP_i^{Opt}(\psi \prec \varphi)$  (resp.  $RP_i^{Pess}(\psi \prec \varphi)$ ) has to be read “according to agent  $i$ ’s optimistic (resp. pessimistic) assessment,  $\varphi$  is realistically better than  $\psi$ ”.

We conclude this section by defining two notions of monadic preference and corresponding two notions of *realistic* monadic preference, respectively noted  $P_i^{Opt}\varphi$ ,  $P_i^{Pess}\varphi$ ,  $RP_i^{Opt}\varphi$  and  $RP_i^{Pess}\varphi$ :

$$\begin{aligned} P_i^{Opt}\varphi &=_{def} P_i^{Opt}(\neg\varphi \prec \varphi) \\ P_i^{Pess}\varphi &=_{def} P_i^{Pess}(\neg\varphi \prec \varphi) \\ RP_i^{Opt}\varphi &=_{def} RP_i^{Opt}(\neg\varphi \prec \varphi) \\ RP_i^{Pess}\varphi &=_{def} RP_i^{Pess}(\neg\varphi \prec \varphi) \end{aligned}$$

An optimistic (resp. pessimistic) agent has a preference for  $\varphi$ , noted  $P_i^{Opt}\varphi$  (resp.  $P_i^{Pess}\varphi$ ), if and only if, according to her optimistic (resp. pessimistic) assessment,  $\varphi$  is better than  $\neg\varphi$ . An optimistic (resp. pessimistic) agent has a realistic preference for  $\varphi$ , noted  $RP_i^{Opt}\varphi$  (resp.  $RP_i^{Pess}\varphi$ ), if and only if, according to her optimistic (resp. pessimistic) assessment,  $\varphi$  is realistically better than  $\neg\varphi$ .

The following validity illustrates the relationship between the notion of desire defined in Definition 6 and the previous notion of pessimistic monadic preference:

$$\models_{MCM} \neg D_i \top \rightarrow (D_i \varphi \leftrightarrow P_i^{Pess}\varphi) \quad (6)$$

This means that if there exists at least a minimally desirable state for agent  $i$  (condition  $\neg D_i \top$ ), then  $i$  desires that  $\varphi$  if and only if, according to her pessimistic assessment,  $\varphi$  is better than  $\neg\varphi$ .

### 3.4 Example

In the previous sections, we have defined a variety of cognitive attitudes of epistemic and motivational type. Let us illustrate them with the help of the crossroad scenario sketched in the introduction. For simplicity, we assume that  $Ag_t = \{1, 2\}$  and that the

set of atomic propositions  $Atm$  includes the following elements with their corresponding meaning:  $co$  (“agent 1 and agent 2 collide”),  $lo_1$  (“agent 1 loses time”) and  $lo_2$  (“agent 2 loses time”).

We are going to make different hypotheses about the agents’ cognitive attitudes and present a number of conclusions that can be drawn from them. Our initial hypothesis concerns the agents’ knowledge:

$$\varphi_1 =_{def} \bigwedge_{i \in \{1,2\}} [\equiv_i] \left( ((lo_1 \wedge \neg lo_2) \rightarrow \neg co) \wedge ((\neg lo_1 \wedge lo_2) \rightarrow \neg co) \wedge \neg(\neg lo_1 \wedge \neg lo_2) \right).$$

According to hypothesis  $\varphi_1$ , agents 1 and 2 know (i) that there will be no collision if one of them loses time by letting the other pass, and (ii) that necessarily one of them will lose time (since if they both pass, there will be a collision so that they will both lose time).

Our second hypothesis concerns what the agents merely envisage:

$$\varphi_2 =_{def} \bigwedge_{i \in \{1,2\}} \left( \langle \equiv_i \rangle co \wedge \langle \equiv_i \rangle (lo_1 \wedge \neg lo_2) \wedge \langle \equiv_i \rangle (\neg lo_1 \wedge lo_2) \wedge \langle \equiv_i \rangle (lo_1 \wedge lo_2 \wedge \neg co) \right).$$

According to hypothesis  $\varphi_2$ , agents 1 and 2 envisage four possible situations: (i) the situations in which they collide, (ii) the two situations in which one of them loses its time while the other does not, and (iii) the situation in which they both lose time because of a collision.

We conclude with the following hypothesis about the agents’ motivations, according to which each agent strongly desires not to collide and strongly desires not to lose time:

$$\varphi_3 =_{def} \bigwedge_{i \in \{1,2\}} (SD_i \neg lo_i \wedge SD_i \neg co).$$

As the following validities highlight, the previous hypotheses lead to different conclusions about the agents’ epistemic and motivational attitudes:

$$\models_{MCM} \varphi_1 \rightarrow \bigwedge_{i \in \{1,2\}} \left( [\equiv_i] (co \rightarrow (lo_1 \wedge lo_2)) \wedge B_i(\neg lo_1, lo_2) \wedge B_i(\neg lo_2, lo_1) \right) \quad (7)$$

$$\models_{MCM} (\varphi_2 \wedge \varphi_3) \rightarrow \bigwedge_{i \in \{1,2\}} (SD_i(\neg lo_i \wedge \neg co) \wedge D_i \neg lo_i \wedge D_i \neg co) \quad (8)$$

$$\models_{MCM} (\varphi_1 \wedge \varphi_2 \wedge \varphi_3) \rightarrow \bigwedge_{i \in \{1,2\}} (D_i(lo_1 \wedge \neg lo_2) \wedge D_i(\neg lo_1 \wedge lo_2)) \quad (9)$$

$$\models_{MCM} (\varphi_2 \wedge \varphi_3) \rightarrow \bigwedge_{i \in \{1,2\}} (\neg SD_i(lo_1 \wedge \neg lo_2) \wedge \neg SD_i(\neg lo_1 \wedge lo_2)) \quad (10)$$

The single hypothesis  $\varphi_1$  leads to the conclusion (i) that the agents know that a collision implies that they both lose time, and (ii) that they believe that an agent loses time

conditional on the fact that the other does not. Thanks to the set of hypotheses  $\{\varphi_2, \varphi_3\}$ , we can conclude (i) that each agent strongly desires not to lose time and to avoid a collision, and (ii) each agent has both the desire not to lose time and the desire to avoid a collision. Finally, thanks to the set of hypotheses  $\{\varphi_1, \varphi_2, \varphi_3\}$ , we can conclude that each agent finds desirable the situations in which only one of them loses time by letting the other pass. As the last validity indicates, such situations are merely desirable for the agent but not strongly desirable.

## 4 Axiomatization

In this section, we provide a sound and complete axiomatization for the Dynamic Logic of Cognitive Attitudes (DLCA). The first step consists in precisely defining this logic which includes several axioms and rule of necessitation for the modalities  $[\pi]$  as well as one non-standard rule of inference for nominals.

**Definition 13 (Logic DLCA)** *We define DLCA to be the extension of classical propositional logic given by the following axioms and rules with  $\tau \in \{P, D\}$ :*

$([\pi]\varphi \wedge [\pi](\varphi \rightarrow \psi)) \rightarrow [\pi]\psi$	(K $_{\pi}$ )
$[\equiv_i]\varphi \rightarrow \varphi$	(T $_{\equiv_i}$ )
$[\equiv_i]\varphi \rightarrow [\equiv_i][\equiv_i]\varphi$	(4 $_{\equiv_i}$ )
$\neg[\equiv_i]\varphi \rightarrow [\equiv_i]\neg[\equiv_i]\varphi$	(5 $_{\equiv_i}$ )
$[\preceq_{i,\tau}]\varphi \rightarrow \varphi$	(T $_{\preceq_{i,\tau}}$ )
$[\preceq_{i,\tau}]\varphi \rightarrow [\preceq_{i,\tau}][\preceq_{i,\tau}]\varphi$	(4 $_{\preceq_{i,\tau}}$ )
$[\equiv_i]\varphi \rightarrow [\preceq_{i,\tau}]\varphi$	(Inc $_{\preceq_{i,\tau}, \equiv_i}$ )
$(\langle \equiv_i \rangle \varphi \wedge \langle \equiv_i \rangle \psi) \rightarrow (\langle \equiv_i \rangle (\varphi \wedge \langle \preceq_{i,\tau} \rangle \psi) \vee \langle \equiv_i \rangle (\psi \wedge \langle \preceq_{i,\tau} \rangle \varphi))$	(Conn $_{\preceq_{i,\tau}, \equiv_i}$ )
$[\pi; \pi']\varphi \leftrightarrow [\pi][\pi']\varphi$	(Red $_{\cdot}$ )
$[\pi \cup \pi']\varphi \leftrightarrow ([\pi]\varphi \wedge [\pi']\varphi)$	(Red $_{\cup}$ )
$([\pi]\varphi \wedge [\pi']\psi) \rightarrow [\pi \cap \pi'](\varphi \wedge \psi)$	(Add1 $_{\cap}$ )
$(\langle \pi \rangle x \wedge \langle \pi' \rangle x) \rightarrow \langle \pi \cap \pi' \rangle x$	(Add2 $_{\cap}$ )
$\varphi \rightarrow [\pi]\langle -\pi \rangle \varphi$	(Conv1 $_{-}$ )
$\varphi \rightarrow [-\pi]\langle \pi \rangle \varphi$	(Conv2 $_{-}$ )
$([\preceq_{i,\tau}]\varphi \wedge [\preceq_{i,\tau}^{\sim}]\varphi) \leftrightarrow [\equiv_i]\varphi$	(Comp1 $_{\sim}$ )
$\langle \preceq_{i,\tau} \rangle x \rightarrow [\preceq_{i,\tau}^{\sim}]\neg x$	(Comp2 $_{\sim}$ )
$[?\varphi]\psi \rightarrow (\varphi \rightarrow \psi)$	(Red $_{?}$ )
$\langle \pi \rangle (x \wedge \varphi) \rightarrow [\pi'](x \rightarrow \varphi)$	(Most $_x$ )
$\frac{\varphi}{[\pi]\varphi}$	(Nec $_{\pi}$ )
$\frac{[\pi]\neg x \text{ for all } x \in \text{Nom}}{[\pi]\perp}$	(Cov)



Note that the primitive operators  $[\preceq_{i,P}]$  and  $[\preceq_{i,D}]$  are S4 (or KT4), while  $[\equiv_i]$  is S5. The only interaction principles between these three operators are the “inclusion” Axiom **Inc** $_{\preceq_{i,\tau},\equiv_i}$  and the “connectedness” Axiom **Conn** $_{\preceq_{i,\tau},\equiv_i}$ . Operators  $[\preceq_{i,P}]$  and  $[\preceq_{i,D}]$  do not interact since, as we have emphasized in Section 2, epistemic plausibility and desirability are assumed to be independent notions.

For every  $\varphi \in \mathcal{L}_{\text{DLCA}}$ , we write  $\vdash \varphi$  to denote the fact that  $\varphi$  is a theorem of DLCA, i.e., there exists an at most countably infinite sequence  $\psi_0, \psi_1, \dots$  such that  $\psi_0 = \varphi$  and for all  $k \geq 0$ ,  $\psi_k$  is an instance of some axiom or  $\psi_k$  can be obtained from some later members of the sequence by an application of some inference rule.

The rest of this section is devoted to prove that the logic DLCA is sound and complete for the class of multi-agent cognitive models.

Soundness, namely checking that the axioms are valid and the rules of inferences preserve validity, is a routine exercise. Notice that the admissibility of the rule of inference **Cov** is guaranteed by the fact that the set of nominals *Nom* is infinite.

As for completeness, the proof is organized in several steps. We use techniques from dynamic logic and modal logic with names [39, 25].

In the rest of this section, we denote sets of formulas from  $\mathcal{L}_{\text{DLCA}}$  by  $\Sigma, \Sigma', \dots$ . Let  $\varphi \in \mathcal{L}_{\text{DLCA}}$  and  $\Sigma \subseteq \mathcal{L}_{\text{DLCA}}$ , we define:

$$\Sigma + \varphi = \{\psi \in \mathcal{L}_{\text{DLCA}} : \varphi \rightarrow \psi \in \Sigma\}.$$

Let us start by defining the concepts of theory and maximal consistent theory.

**Definition 14 (Theory)** *A set of formulas  $\Sigma$  is said to be a theory if it contains all theorems of DLCA and is closed under modus ponens and rule **Cov**. It is said to be a consistent theory if it is a theory and  $\perp \notin \Sigma$ . It is said to be a maximal consistent theory (MCT) if it is a consistent theory and, for each consistent theory  $\Sigma'$ , we have that if  $\Sigma \subseteq \Sigma'$  then  $\Sigma = \Sigma'$ .*

We have the following property for theories.

**Proposition 11** *Let  $\Sigma$  be a theory and let  $\varphi \in \mathcal{L}_{\text{DLCA}}$ . Then,  $\Sigma + \varphi$  is a theory. Moreover, if  $\Sigma$  is consistent then either  $\Sigma + \varphi$  is consistent or  $\Sigma + \neg\varphi$  is consistent.*

**PROOF.** Let us first prove that if  $\Sigma$  is a theory then  $\Sigma + \varphi$  is a theory as well. Suppose  $\Sigma$  is a theory. Then,  $\Sigma + \varphi$  clearly contains all theorems of DLCA. Moreover, suppose  $\psi \rightarrow \psi', \psi \in \Sigma + \varphi$ . Thus, by definition of  $\Sigma + \varphi$ , we have  $\varphi \rightarrow \psi, \varphi \rightarrow (\psi \rightarrow \psi') \in \Sigma$ . Since  $\Sigma$  is closed under modus ponens and contains all theorems of DLCA, the latter implies  $(\varphi \rightarrow \psi) \wedge (\varphi \rightarrow (\psi \rightarrow \psi')) \in \Sigma$ . Consequently, since  $\Sigma$  is closed under modus ponens,  $\varphi \rightarrow \psi' \in \Sigma$ . Hence,  $\psi' \in \Sigma + \varphi$ . This means that  $\Sigma + \varphi$  is closed under modus ponens. Finally, let us show that  $\Sigma + \varphi$  is closed under **Cov**. Suppose  $[\pi] \neg x \in \Sigma + \varphi$  for all  $x$ . Thus, by definition of  $\Sigma + \varphi$ ,  $\varphi \rightarrow [\pi] \neg x \in \Sigma$  for all  $x$ . Since  $\Sigma$  is a theory, the latter implies that  $[\varphi; \pi] \neg x \in \Sigma$  for all  $x$ . Thus, since  $\Sigma$  is a theory,  $[\varphi; \pi] \perp \in \Sigma$  and, consequently,  $\varphi \rightarrow [\pi] \perp \in \Sigma$ . It follows that  $[\pi] \perp \in \Sigma + \varphi$ .

Let us show that if  $\Sigma$  is consistent then either  $\Sigma + \varphi$  is consistent or  $\Sigma + \neg\varphi$  is consistent. Suppose the antecedent is true while the consequent is false. Then,  $\varphi \rightarrow$

$\perp \in \Sigma$  and  $\neg\phi \rightarrow \perp \in \Sigma$ . Since  $\Sigma$  is a theory, we have  $(\phi \rightarrow \perp) \wedge (\neg\phi \rightarrow \perp) \in \Sigma$ . Thus,  $\perp \in \Sigma$  which is in contradiction with the fact that  $\Sigma$  is consistent. ■

The following proposition highlights some standard properties of MCTs.

**Proposition 12** *Let  $\Sigma$  be a MCT. Then, for all  $\phi, \psi \in \mathcal{L}_{\text{DLCA}}$ :*

- $\phi \in \Sigma$  or  $\neg\phi \in \Sigma$ ,
- $\phi \vee \psi \in \Sigma$  iff  $\phi \in \Sigma$  or  $\psi \in \Sigma$ .

PROOF. We only prove the first item by reductio ad absurdum. Suppose  $\Sigma$  is a MCT,  $\phi \notin \Sigma$  and  $\neg\phi \notin \Sigma$ . We clearly have  $\Sigma \subseteq \Sigma + \phi$  and  $\Sigma \subseteq \Sigma + \neg\phi$ . Moreover,  $\phi \in \Sigma + \phi$  and  $\neg\phi \in \Sigma + \neg\phi$ . Thus,  $\Sigma \subset \Sigma + \phi$  and  $\Sigma \subset \Sigma + \neg\phi$ . By Proposition 11,  $\Sigma + \phi$  and  $\Sigma + \neg\phi$  are theories. Moreover, either  $\Sigma + \phi$  is consistent or  $\Sigma + \neg\phi$  is consistent. This contradicts the fact that  $\Sigma$  is a MCT. ■

The following variant of the Lindenbaum's lemma is proved in the same way as [39, Lemma 4.15].

**Lemma 1** *Let  $\Sigma$  be a consistent theory and let  $\phi \notin \Sigma$ . Then, there exists a MCT  $\Sigma^+$  such that  $\Sigma \subseteq \Sigma^+$  and  $\phi \notin \Sigma^+$ .*

The following lemma highlights a fundamental property of MCTs.

**Lemma 2** *Let  $\Sigma$  be a MCT. Then, there exists  $x \in \text{Nom}$  such  $x \in \Sigma$ .*

PROOF. We prove the lemma by reductio ad absurdum. Let  $\Sigma$  be a MCT. Moreover, suppose that, for all  $x \in \text{Nom}$ ,  $x \notin \Sigma$ . By Proposition 12, it follows that, for all  $x \in \text{Nom}$ ,  $\neg x \in \Sigma$ .

By Axiom **Red**<sub>?</sub>, we have  $\neg x \leftrightarrow [?\top]\neg x \in \Sigma$  for all  $x \in \text{Nom}$ . Thus, for all  $x \in \text{Nom}$ ,  $[?\top]\neg x \in \Sigma$ . Hence, since  $\Sigma$  is closed under **Cov**,  $[?\top]\perp \in \Sigma$ . By Axiom **Red**<sub>?</sub>, the latter is equivalent to  $\perp \in \Sigma$ . The latter is contradiction with the fact that  $\Sigma$  is a MCT. ■

Let us now define the canonical model for our logic.

**Definition 15 (Canonical model)** *The canonical model is the tuple  $M^c = (W^c, (\preceq_{i,P}^c)_{i \in \text{Agt}}, (\preceq_{i,D}^c)_{i \in \text{Agt}}, (\equiv_i^c)_{i \in \text{Agt}}, V^c)$  such that:*

- $W^c$  is the set of all MCTs,
- for all  $i \in \text{Agt}$ , for all  $\tau \in \{P, D\}$ , for all  $w, v \in W^c$ ,  $w \preceq_{i,\tau}^c v$  iff, for all  $\phi \in \mathcal{L}_{\text{DLCA}}$ , if  $[\preceq_{i,\tau}]\phi \in w$  then  $\phi \in v$ ,
- for all  $i \in \text{Agt}$ , for all  $w, v \in W^c$ ,  $w \equiv_i^c v$  iff, for all  $\phi \in \mathcal{L}_{\text{DLCA}}$ , if  $[\equiv_i]\phi \in w$  then  $\phi \in v$ ,
- for all  $w \in W^c$ ,  $V^c(w) = (\text{Atm} \cup \text{Nom}) \cap w$ .

Let us now define the canonical relations for the complex programs  $\pi$ .

**Definition 16 (Canonical relation)** Let  $M^c = (W^c, (\preceq_{i,P}^c)_{i \in \text{Agt}}, (\preceq_{i,D}^c)_{i \in \text{Agt}}, (\equiv_i^c)_{i \in \text{Agt}}, V^c)$  be the canonical model. Then, for all  $\pi \in \mathcal{P}$  and for all  $w, v \in W^c$ :

$$wR_\pi^c v \text{ iff, for all } \varphi \in \mathcal{L}_{\text{DLCA}}, \text{ if } [\pi]\varphi \in w \text{ then } \varphi \in v.$$

The following Lemma 3 highlights one fundamental property of the canonical model.

**Lemma 3** Let  $M^c = (W^c, (\preceq_{i,P}^c)_{i \in \text{Agt}}, (\preceq_{i,D}^c)_{i \in \text{Agt}}, (\equiv_i^c)_{i \in \text{Agt}}, V^c)$  be the canonical model. Then, for all  $\Sigma, \Sigma' \in W^c$ , for all  $\pi \in \mathcal{P}$  and for all  $x \in \text{Nom}$ , if  $x \in \Sigma, x \in \Sigma'$  and  $\Sigma R_\pi^c \Sigma'$  then  $\Sigma = \Sigma'$ .

PROOF. Let us first prove that (i) if  $x \in \Sigma$  and  $\varphi \in \Sigma$  then  $[\pi](x \rightarrow \varphi) \in \Sigma$ . Suppose  $x, \varphi \in \Sigma$ . Thus,  $x \wedge \varphi \in \Sigma$  since  $\Sigma$  is a MCT. Moreover,  $(x \wedge \varphi) \rightarrow [\pi](x \rightarrow \varphi) \in \Sigma$ , because of Axiom **Most<sub>x</sub>**. Hence,  $[\pi](x \rightarrow \varphi) \in \Sigma$ .

Now let us prove by absurdum that (ii) if  $x \in \Sigma, \Sigma'$  and  $\Sigma R_\pi^c \Sigma'$  then  $\Sigma = \Sigma'$ . Suppose  $x \in \Sigma, \Sigma', \Sigma R_\pi^c \Sigma'$  and  $\Sigma \neq \Sigma'$ . The latter implies that there exists  $\varphi$  such that  $\varphi \in \Sigma$  and  $\varphi \notin \Sigma'$ . By item (i) above, it follows that  $[\pi](x \rightarrow \varphi) \in \Sigma$ . Since  $\Sigma R_\pi^c \Sigma'$ , the latter implies that  $x \rightarrow \varphi \in \Sigma'$ . Since  $x \in \Sigma'$ , it follows that  $\varphi \in \Sigma'$  which leads to a contradiction. ■

The next step consists in proving the following existence lemma.

**Lemma 4** Let  $M^c = (W^c, (\preceq_{i,P}^c)_{i \in \text{Agt}}, (\preceq_{i,D}^c)_{i \in \text{Agt}}, (\equiv_i^c)_{i \in \text{Agt}}, V^c)$  be the canonical model, let  $w \in W^c$  and let  $\langle \pi \rangle \varphi \in \mathcal{L}_{\text{DLCA}}$ . Then, if  $\langle \pi \rangle \varphi \in w$  then there exists  $v \in W^c$  such that  $wR_\pi^c v$  and  $\varphi \in v$ .

PROOF. Suppose  $w$  is a MCT and  $\langle \pi \rangle \varphi \in w$ . It follows that  $[\pi]w = \{\psi : [\pi]\psi \in w\}$  is a consistent theory. Indeed, it is easy to check that  $[\pi]w$  contains all theorems of DLCA, is closed under modus ponens and rule **Cov**. Let us prove that it is consistent by reductio ad absurdum. Suppose  $\perp \in [\pi]w$ . Thus,  $[\pi]\perp \in w$ . Hence,  $[\pi]\neg\varphi \in w$ . Since  $\langle \pi \rangle \varphi \in w$ ,  $\perp \in w$ . The latter contradicts the fact that  $w$  is a MCT. Let us distinguish two cases.

Case 1:  $\varphi \in [\pi]w$ . Thus,  $\neg\varphi \notin [\pi]w$  since  $w$  is consistent. Thus, by Lemma 1, there exists MCT  $v$  such that  $[\pi]w \subseteq v$ ,  $\varphi \in v$  and  $\neg\varphi \notin v$ . By definition of  $R_\pi^c$ ,  $wR_\pi^c v$ .

Case 2:  $\varphi \notin [\pi]w$ . By Proposition 11,  $[\pi]w + \varphi$  is a theory since  $[\pi]w$  is a theory.  $[\pi]w + \varphi$  is consistent. Suppose it is not. Thus,  $\varphi \rightarrow \perp \in [\pi]w$  and, consequently,  $\neg\varphi \in [\pi]w$ . Hence,  $[\pi]\neg\varphi \in w$ . It follows that  $\perp \in w$ , since  $\langle \pi \rangle \varphi \in w$ . But this contradicts the fact that  $w$  is a MCT. Thus,  $[\pi]w + \varphi$  is a consistent theory. Moreover,  $\varphi \in [\pi]w + \varphi$ ,  $\neg\varphi \notin [\pi]w + \varphi$  and  $[\pi]w \subseteq [\pi]w + \varphi$ . By Lemma 1, there exists MCT  $v$  such that  $[\pi]w \subseteq v$ ,  $\varphi \in v$  and  $\neg\varphi \notin v$ . By definition of  $R_\pi^c$ ,  $wR_\pi^c v$ . ■

The following truth lemma is proved in the usual way by induction on the structure of  $\varphi$  thanks to Lemma 4.

**Lemma 5** Let  $M^c = (W^c, (\preceq_{i,P}^c)_{i \in \text{Agt}}, (\preceq_{i,D}^c)_{i \in \text{Agt}}, (\equiv_i^c)_{i \in \text{Agt}}, V^c)$  be the canonical model, let  $w \in W^c$  and let  $\varphi \in \mathcal{L}_{\text{DLCA}}$ . Then,  $M^c, w \models \varphi$  iff  $\varphi \in w$ .

PROOF. The proof is by induction on the structure of  $\phi$ . We only prove the case in which  $\phi$  is a modal formula  $[\pi]\psi$ . As for the right-to-left direction we have:

$$\begin{aligned} [\pi]\psi \in w & \text{ only if } \forall v \in R_\pi^c(w) : \psi \in v \text{ (by definition of } R_\pi^c) \\ & \text{iff } \forall v \in R_\pi^c(w) : M^c, v \models \psi \text{ (by induction hypothesis)} \\ & \text{iff } M^c, w \models [\pi]\psi \end{aligned}$$

As for the left-to-right direction, we prove that if  $[\pi]\psi \notin w$  then  $M^c, w \not\models [\pi]\psi$  that, given the property of MCSs, is equivalent to proving that if  $\langle \pi \rangle \psi \in w$  then  $M^c, w \models \langle \pi \rangle \psi$ . Suppose  $\langle \pi \rangle \psi \in w$ . Then, by Lemma 4, there exists  $v \in W^c$  such that  $wR_\pi^c v$  and  $\psi \in v$ . Hence, by induction hypothesis, there exists  $v \in W^c$  such that  $wR_\pi^c v$  and  $M^c, v \models \psi$ . The latter is equivalent to  $M^c, w \models \langle \pi \rangle \psi$ . ■

The pre-final stage of the proof consists in introducing an alternative semantics for the language  $\mathcal{L}_{DLCA}$  which turns out to be equivalent to the original semantics based on MCMs.

**Definition 17 (Quasi multi-agent cognitive model)** *A quasi multi-agent cognitive model (quasi-MCM) is a tuple  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  where  $W$ ,  $\preceq_{i,P}$ ,  $\preceq_{i,D}$ ,  $\equiv_i$  and  $V$  are as in Definition 1 except that Constraint C4 is replaced by the following weaker constraint. For all  $w, v \in W$ :*

**(C4\*)** *if  $V_{\text{Nom}}(w) \cap V_{\text{Nom}}(v) \neq \emptyset$  and  $wR_\pi v$  for some  $\pi \in \mathcal{P}$  then  $w = v$ .*

By the generated submodel property, it is easy to show that the semantics in terms of MCMs and the semantics in terms of quasi-MCMs are equivalent with respect to the language  $\mathcal{L}_{DLCA}$ .

**Proposition 13** *Let  $\phi \in \mathcal{L}_{DLCA}$ . Then,  $\phi$  is valid relative to the class of MCMs if and only if  $\phi$  is valid relative to the class of quasi-MCMs.*

The following theorem highlights that the canonical model is indeed a structure of the right type.

**Lemma 6** *The canonical model  $M^c$  is a quasi-MCM.*

PROOF. The fact that  $M^c$  satisfies Constraints C3 and C4\* follows from Lemma 2 and Lemma 3. To prove that  $\equiv_i$  is an equivalence relation that  $\preceq_{i,D}^c$  and  $\preceq_{i,P}^c$  are preorders and that  $M^c$  satisfies Constraints C1 and C2 is just a routine exercise. Indeed, Axioms  $\mathbf{T}_{\equiv_i}$ ,  $\mathbf{4}_{\equiv_i}$ ,  $\mathbf{5}_{\equiv_i}$ ,  $\mathbf{T}_{\preceq_{i,\tau}}$ ,  $\mathbf{4}_{\preceq_{i,\tau}}$ ,  $\mathbf{Inc}_{\preceq_{i,\tau}, \equiv_i}$  and  $\mathbf{Conn}_{\preceq_{i,\tau}, \equiv_i}$  are canonical for these semantic conditions.

To conclude, we need to prove that the following six conditions hold, for  $i \in \text{Agt}$

and  $\tau \in \{P, D\}$ :

$$\begin{aligned}
(w, v) \in R_{\succeq_{i, \tau}}^c & \text{ iff } (w, v) \in R_{\equiv_i}^c \text{ and } (w, v) \notin R_{\prec_{i, \tau}}^c \\
(w, v) \in R_{\pi, \pi'}^c & \text{ iff } \exists u \in W^c : (w, u) \in R_{\pi}^c \text{ and } (u, v) \in R_{\pi'}^c \\
(w, v) \in R_{\pi \cup \pi'}^c & \text{ iff } (w, v) \in R_{\pi}^c \text{ or } (w, v) \in R_{\pi'}^c \\
(w, v) \in R_{\pi \cap \pi'}^c & \text{ iff } (w, v) \in R_{\pi}^c \text{ and } (w, v) \in R_{\pi'}^c \\
(w, v) \in R_{-\pi}^c & \text{ iff } (v, w) \in R_{\pi}^c \\
wR_{\phi?}^c v & \text{ iff } w = v \text{ and } M^c, w \models \phi
\end{aligned}$$

We only prove the second and fourth conditions which are the most difficult ones to prove.

Let us start with the proof of the second condition. The right-to-left direction is standard. We only prove the left-to-right direction. Suppose  $(w, v) \in R_{\pi, \pi'}^c$ . Let  $[\pi]w = \{\psi : [\pi]\psi \in w\}$ . Moreover, let  $\langle \pi' \rangle v = \{\langle \pi' \rangle \psi : \psi \in v\}$ . Finally, let  $\langle \pi' \rangle \psi_1, \langle \pi' \rangle \psi_2, \dots$  be an enumeration of the elements of  $\langle \pi' \rangle v$ . We define  $\Sigma^1 = [\pi]w + \langle \pi' \rangle \psi_1$  and, for all  $k > 1$ ,  $\Sigma^k = \Sigma^{k-1} + \langle \pi' \rangle \psi_k$ . By Lemma 11 and the fact that  $[\pi]w$  is a theory, it can be shown that every  $\Sigma^k$  is a theory. Moreover, by induction on  $k$ , it can be shown that every  $\Sigma^k$  is consistent. Since  $\Sigma^{k-1} \subseteq \Sigma^k$  for all  $k > 1$ , it follows that  $\Sigma = \bigcup_{k>1} \Sigma^k$  is a consistent theory. By Lemma 1 and the definition of  $\Sigma$ , there exists  $u \in W^c$  such that  $\Sigma \subseteq u$ ,  $(w, u) \in R_{\pi}^c$  and  $(u, v) \in R_{\pi'}^c$ .

Let us now prove the fourth condition. Suppose  $(w, v) \in R_{\pi \cap \pi'}^c$ . By Definition 16 and Proposition 12, it follows that, for all  $\phi$ , if  $\phi \in v$  then  $\langle \pi \cap \pi' \rangle \phi \in w$ . The latter implies that for all  $\phi$ , if  $\phi \in v$  then  $\langle \pi \cap \pi' \rangle (\phi \vee \perp) \in w$  since  $\vdash \langle \pi \cap \pi' \rangle \phi \rightarrow \langle \pi \cap \pi' \rangle (\phi \vee \perp)$ . By Axiom **K** $_{\pi}$ , it follows that, for all  $\phi$ , if  $\phi \in v$  then  $\langle \pi \rangle \phi \vee \langle \pi' \rangle \perp \in w$ . Thus, for all  $\phi$ , if  $\phi \in v$  then  $\langle \pi \rangle \phi \in w$ , since  $\vdash (\langle \pi \rangle \phi \vee \langle \pi' \rangle \perp) \rightarrow \langle \pi \rangle \phi$ . In a similar way, we can prove that, for all  $\phi$ , if  $\phi \in v$  then  $\langle \pi' \rangle \phi \in w$ . By Definition 16 and Proposition 12, it follows that  $(w, v) \in R_{\pi}^c$  and  $(w, v) \in R_{\pi'}^c$ .

Now suppose  $(w, v) \in R_{\pi}^c$  and  $(w, v) \in R_{\pi'}^c$ . Thus, by Definition 16 and Proposition 12, (i) for all  $\phi$ , if  $\phi \in v$  then  $\langle \pi \rangle \phi \in w$  and  $\langle \pi' \rangle \phi \in w$ . By Proposition 12 and Lemma 2, we have that (ii) there exists  $x \in \text{Nom}$  such that, for all  $\phi$ ,  $\phi \in v$  iff  $x \wedge \phi \in v$ . Item (i) and item (ii) together imply that (iii) there exists  $x \in \text{Nom}$  such that, for all  $\phi$ , if  $\phi \in v$  then  $\langle \pi \rangle (x \wedge \phi) \in w$  and  $\langle \pi' \rangle (x \wedge \phi) \in w$ . We are going to prove the following theorem:

$$\vdash (\langle \pi \rangle (x \wedge \phi) \wedge \langle \pi' \rangle (x \wedge \phi)) \rightarrow \langle \pi \cap \pi' \rangle (x \wedge \phi)$$

By Axiom **K** $_{\pi}$ ,  $\langle \pi \rangle (x \wedge \phi) \wedge \langle \pi' \rangle (x \wedge \phi)$  implies  $\langle \pi \rangle x \wedge \langle \pi' \rangle x$ . By Axiom **Add2** $_{\cap}$ , the latter implies  $\langle \pi \cap \pi' \rangle x$ . Moreover, by Axiom **Inc** $_{\prec_{i, \tau}, \equiv_i}$  and Axiom **Most** $_x$ ,  $\langle \pi \rangle (x \wedge \phi)$  implies  $[\equiv_{\emptyset}](x \rightarrow \phi)$ . By Axiom **Inc** $_{\prec_{i, \tau}, \equiv_i}$ , the latter implies  $[\pi \cap \pi'](x \rightarrow \phi)$ . By Axiom **K** $_{\pi}$ ,  $[\pi \cap \pi'](x \rightarrow \phi)$  and  $\langle \pi \cap \pi' \rangle x$  together imply  $\langle \pi \cap \pi' \rangle (x \wedge \phi)$ . Thus,  $\langle \pi \rangle (x \wedge \phi) \wedge \langle \pi' \rangle (x \wedge \phi)$  implies  $\langle \pi \cap \pi' \rangle (x \wedge \phi)$ .

From previous item (iii) and the previous theorem it follows that there exists  $x \in \text{Nom}$  such that, for all  $\phi$ , if  $\phi \in v$  then  $\langle \pi \cap \pi' \rangle (x \wedge \phi)$ . The latter implies that, for all  $\phi$ , if  $\phi \in v$  then  $\langle \pi \cap \pi' \rangle \phi$ . The latter implies that  $(w, v) \in R_{\pi \cap \pi'}^c$ . ■

Let us conclude the proof by supposing  $\not\models \neg \phi$ . Therefore, by Lemma 1 and the fact that the set of DLCA-theorems is a consistent theory, there exists a MCT  $w$  such that

$\neg\phi \notin w$ . Thus, by Proposition 12, we can find a MCT  $w$  such that  $\phi \in w$ . By Lemma 5, the latter implies  $M^c, w \models \phi$  for some  $w \in W^c$ . Since, by Lemma 6,  $M^c$  is a quasi-MCM, it follows that  $\phi$  is satisfiable relative to the class of quasi-MCMs. Therefore, by Proposition 13,  $\phi$  is satisfiable relative to the class of MCMs.

We can finally state the main result of this section.

**Theorem 1** *The logic DLCA is sound and complete for the class of multi-agent cognitive models.*

## 5 Application to Game Theory

In this section, we apply our logical framework to the analysis of single-stage games under incomplete information in which agents only play once (i.e., interaction is non-repeated) and may not know some relevant characteristic of others including their preferences, choices and beliefs.

Let  $Act$  be a set of action names with elements noted  $a, b, \dots$ . Let a joint action be a function  $\delta : Agt \longrightarrow Act$  and the set of joint actions be denoted by  $JAct$ .

For every coalition  $C \in 2^{Agt}$  and for every  $\delta \in JAct$ , let  $\delta_C$  be the  $C$ -restriction of  $\delta$ , that is, the function  $\delta_C : C \longrightarrow Act$  such that  $\delta_C(i) = \delta(i)$  for all  $i \in C$ . For notational convenience, we write  $-i$  instead of  $Agt \setminus \{i\}$ , with  $i \in Agt$ .

In order to model strategic interaction in our setting, we extend MCMs of Definition 1 by agents' choices. We call MCM with choices the resulting models.

**Definition 18 (Multi-agent cognitive model with choices)** *A multi-agent cognitive model with choices (MCMC) is a tuple  $M = (W, (\preceq_{i,P})_{i \in Agt}, (\preceq_{i,D})_{i \in Agt}, (\equiv_i)_{i \in Agt}, (C_i)_{i \in Agt}, V)$ , where  $M = (W, (\preceq_{i,P})_{i \in Agt}, (\preceq_{i,D})_{i \in Agt}, (\equiv_i)_{i \in Agt}, V)$  is a MCM and every  $C_i$  is a choice function  $C_i : W \longrightarrow Act$ , which satisfies the following constraint, for each  $i \in Agt$  and  $\delta \in JAct$ :*

**(C5)** *if  $\forall j \in Agt, \exists w_j \in W$  such that  $w \equiv_i w_j$  and  $C_j(w_j) = \delta(j)$ , then  $\exists v \in W$  such that  $w \equiv_i v$  and,  $\forall j \in Agt, C_j(v) = \delta(j)$ .*

For every  $w \in W$ ,  $C_i(w)$  denotes agent  $i$ 's actual choice at  $w$ . If  $w \equiv_i v$  and  $C_j(v) = a$ , then  $a$  is a potential choice of agent  $j$  from agent  $i$ 's perspective.

According to Constraint C5, agents' choices are subjectively independent, in the sense that every agent  $i$  knows that an agent cannot be deprived of her choices due to the choices made by the others. In other words, suppose that, from agent  $i$ 's perspective,  $\delta(j)$  is a potential choice of  $j$  for every agent  $j$ . Then, from agent  $i$ 's perspective, there should be a state at which the agents choose the joint action  $\delta$ . It is a subjective version of the property of choice independence formulated in the “seeing to it that” (STIT) framework [12, 33, 7].

At the syntactic level, we extend the language  $\mathcal{L}_{DLCA}$  by special constants for choices of type  $\text{play}(i, a)$ , with  $i \in Agt$  and  $a \in Act$ , denoting the fact that “agent  $i$  plays (or chooses) action  $a$ ”. The resulting language is noted  $\mathcal{L}_{DLCAG}$ , where DLCAG stands for “Dynamic Logic of Cognitive Attitudes in Games” and a constant  $\text{play}(i, a)$  is interpreted relative to a MCMC  $M$  and a world  $w$  in  $M$ , as follows:

$$M, w \models \text{play}(i, a) \iff C_i(w) = a.$$

Let  $\delta \in JAct$  and  $C \in 2^{Agt}$ . We define:

$$\text{play}(\delta_C) =_{\text{def}} \bigwedge_{i \in C} \text{play}(i, \delta_C(i)).$$

For every formula  $\varphi$  in  $\mathcal{L}_{\text{DLCAG}}$  we say that  $\varphi$  is valid, noted  $\models_{\text{MCMC}} \varphi$ , if and only if for every multi-agent cognitive model with choices  $M$  and world  $w$  in  $M$ , we have  $M, w \models \varphi$ .

**Definition 19 (Logic DLCAG)** We define DLCAG to be the extension of logic DLCA given by the following axioms:

$$\text{play}(i, a) \rightarrow \neg \text{play}(i, b) \text{ if } a \neq b \quad (\text{MostAct})$$

$$\bigvee_{a \in Act} \text{play}(i, a) \quad (\text{LeastAct})$$

$$\left( \bigwedge_{j \in Agt} \langle \equiv_i \rangle \text{play}(j, \delta(j)) \right) \rightarrow \langle \equiv_i \rangle \text{play}(\delta_{Agt}) \quad (\text{SIC})$$

Axiom **MostAct** means that an agent chooses at most one action from  $Act$  while, according to Axiom **LeastAct**, an agent chooses at least one action from  $Act$ . Axiom **SIC** is the syntactic counterpart of *subjective choice independence* expressed by Constraint C5.

We can adapt the techniques used for proving Theorem 1 in order to prove the following Theorem 2.

**Theorem 2** *The logic DLCAG is sound and complete for the class of multi-agent cognitive models with choices.*

PROOF. Verifying that the logic DLCAG is sound for the class of MCMCs is a routine exercise. As for completeness, the proof is just a straightforward adaptation of the proof of completeness of the logic DLCA. First, we need to define corresponding notions of theory and maximal consistent theory for the logic DLCAG which are akin to the ones for the logic DLCA and use them to define the canonical model and the canonical relation for DLCAG. The canonical model for DLCAG is defined to be a tuple  $M^c = (W^c, (\preceq_{i,P}^c)_{i \in Agt}, (\preceq_{i,D}^c)_{i \in Agt}, (\equiv_i^c)_{i \in Agt}, (C_i^c)_{i \in Agt}, V^c)$  where  $W^c$  is the set of all MCTs for DLCAG,  $\preceq_{i,P}^c$ ,  $\preceq_{i,D}^c$ ,  $\equiv_i^c$  and  $V^c$  are defined as in the definition of the canonical model for DLCA (Definition 15), and  $C_i^c : W^c \rightarrow 2^{Act}$  such that, for all  $a \in Act$  and  $w \in W^c$ ,  $a \in C_i^c(w)$  if and only if  $\text{play}(i, a) \in w$ . The canonical relation for DLCAG is defined in the same way as the canonical relation for DLCA (Definition 16).

It is immediate to adapt the proof of the existence and truth lemma for DLCA (Lemma 4 and Lemma 5) to prove corresponding existence and truth lemma for DLCAG.

Secondly, we need to define the notion of quasi multi-agent cognitive model with choices (quasi-MCMC) which is analogous to the definition of quasi-MCM (Definition 17). In particular, a quasi-MCMC is defined to be a tuple  $M = (W, (\preceq_{i,P})_{i \in Agt}, (\preceq_{i,D})_{i \in Agt}, (\equiv_i)_{i \in Agt}, (C_i)_{i \in Agt}, V)$  where  $W$ ,  $\preceq_{i,P}$ ,  $\preceq_{i,D}$ ,  $\equiv_i$ ,  $C_i$  and  $V$  are as in Definition 18 except that Constraint C4 is replaced by the weaker Constraint C4\* of Definition 17. As

for MCMs, by the generated submodel property, it is easy to show that the semantics in terms of MCMs and the semantics in terms of quasi-MCMs are equivalent with respect to the language  $\mathcal{L}_{\text{DLCAG}}$ .

The only property that has to be checked carefully is whether the canonical model for DLCAG is indeed a quasi-MCMC. To this aim, we need to extend the proof of Lemma 6 in order to verify that the canonical model for DLCAG satisfies Constraint C5 of Definition 18 and that, for all  $i \in \text{Act}$  and for all  $w \in W^c$ ,  $C_i^c(w)$  is a singleton. Suppose  $a, b \in C_i^c(w)$  for  $a \neq b$ . The latter means that  $\text{play}(i, a), \text{play}(i, b) \in w$ . We have  $\text{play}(i, a) \rightarrow \neg \text{play}(i, b) \in w$ , because of Axiom **MostAct**. Thus,  $\neg \text{play}(i, b) \in w$ . Hence,  $\perp \in w$  which contradicts the fact that  $w$  is a consistent theory. Consequently, the set  $C_i^c(w)$  has at most one element. Now, let us prove that  $C_i^c(w)$  has at least one element. Because of Axiom **LeastAct**, we have  $\bigvee_{a \in \text{Act}} \text{play}(i, a) \in w$ . Thus, there exists  $a \in \text{Act}$  such that  $\text{play}(i, a) \in w$ . Hence,  $C_i^c(w)$  is non-empty. Now, let us prove that the canonical model for DLCAG satisfies Constraint C5. Suppose  $\forall j \in \text{Agt}, \exists w_j \in W^c$  such that  $w \equiv_i^c w_j$  and  $C_j^c(w_j) = \delta(j)$ . The latter means that  $\forall j \in \text{Agt}, \exists w_j \in W^c$  such that  $w \equiv_i^c w_j$  and  $\text{play}(j, \delta(j)) \in w_j$ . Thus, we have that,  $\forall j \in \text{Agt}, \langle \equiv_i \rangle \text{play}(j, \delta(j)) \in w$ . Hence,  $\bigwedge_{j \in \text{Agt}} \langle \equiv_i \rangle \text{play}(j, \delta(j)) \in w$ . By Axiom **SIC**,  $\left( \bigwedge_{j \in \text{Agt}} \langle \equiv_i \rangle \text{play}(j, \delta(j)) \right) \rightarrow \langle \equiv_i \rangle \text{play}(\delta_{\text{Agt}}) \in w$ . Consequently,  $\langle \equiv_i \rangle \text{play}(\delta_{\text{Agt}}) \in w$ . By the existence lemma for DLCAG, the latter implies that  $\exists v \in W^c$  such that  $w \equiv_i^c v$  and  $\text{play}(\delta_{\text{Agt}}) \in v$ . Thus,  $\exists v \in W^c$  such that  $w \equiv_i^c v$  and  $C_j(v) = \delta(j)$  for every  $j \in \text{Agt}$ . ■

With the support of the language  $\mathcal{L}_{\text{DLCAG}}$ , we can define a variety of notions from the theory of games under incomplete information. The first notion we consider is best response, both from the perspective of an optimistic agent and from the perspective of a pessimistic one:

$$\begin{aligned} \text{BR}_i^{\text{Opt}}(a, \delta_{-i}) &=_{\text{def}} \bigwedge_{b \in \text{Act}} \text{RP}_i^{\text{Opt}} \left( (\text{play}(i, b) \wedge \text{play}(\delta_{-i})) \preceq (\text{play}(i, a) \wedge \text{play}(\delta_{-i})) \right), \\ \text{BR}_i^{\text{Pess}}(a, \delta_{-i}) &=_{\text{def}} \bigwedge_{b \in \text{Act}} \text{RP}_i^{\text{Pess}} \left( (\text{play}(i, b) \wedge \text{play}(\delta_{-i})) \preceq (\text{play}(i, a) \wedge \text{play}(\delta_{-i})) \right). \end{aligned}$$

We say that playing action  $a$  is for agent  $i$  an optimistic (resp. pessimistic) best response to the others' joint action  $\delta_{-i}$ , noted  $\text{BR}_i^{\text{Opt}}(a, \delta_{-i})$  (resp.  $\text{BR}_i^{\text{Pess}}(a, \delta_{-i})$ ) if and only if for every action  $b$ , according to agent  $i$ 's optimistic (resp. pessimistic) assessment, playing  $a$  while the others play  $\delta_{-i}$  is realistically at least as good as playing  $b$  while the others play  $\delta_{-i}$ .

As for best response, we can define two types of *subjective* Nash equilibrium, one for optimistic agents and the other for pessimistic ones. Our notion of subjective Nash equilibrium corresponds to a *qualitative* variant of the notion of Bayesian Nash equilibrium (BNE): a similar qualitative variant of BNE is studied by [2] in the context of possibility theory. The joint action  $\delta$  is said to be a subjective optimistic (resp. pessimistic) Nash equilibrium, noted  $\text{NE}^{\text{Opt}}(\delta)$  (resp.  $\text{NE}^{\text{Pess}}(\delta)$ ), if no agent  $i$  wants to unilaterally deviate from the chosen strategy  $\delta(i)$ , under that the assumption that  $i$  is



optimistic (resp. pessimistic):

$$\begin{aligned} \text{NE}^{Opt}(\delta) &=_{\text{def}} \bigwedge_{i \in \text{Agt}} \text{BR}_i^{Opt}(\delta(i), \delta_{-i}), \\ \text{NE}^{Pess}(\delta) &=_{\text{def}} \bigwedge_{i \in \text{Agt}} \text{BR}_i^{Pess}(\delta(i), \delta_{-i}). \end{aligned}$$

Note that assuming the finiteness of the set of agents  $\text{Agt}$  is essential for defining Nash equilibrium, since our language is finitary and does not allow universal quantification over infinite sets.

Given the distinction between optimistic and pessimistic agent, two notions of rationality can be defined. Agent  $i$  is said to be optimistic (resp. pessimistic) rational, noted  $\text{Rat}_i^{Opt}$  (resp.  $\text{Rat}_i^{Pess}$ ), if she cannot choose an action that, according to her optimistic (resp. pessimistic) assessment, is better not to choose than to choose:

$$\begin{aligned} \text{Rat}_i^{Opt} &=_{\text{def}} \bigwedge_{a \in \text{Act}} \left( \text{play}(i, a) \rightarrow \text{RP}_i^{Opt}(\neg \text{play}(i, a) \preceq \text{play}(i, a)) \right), \\ \text{Rat}_i^{Pess} &=_{\text{def}} \bigwedge_{a \in \text{Act}} \left( \text{play}(i, a) \rightarrow \text{RP}_i^{Pess}(\neg \text{play}(i, a) \preceq \text{play}(i, a)) \right). \end{aligned}$$

As the following proposition indicates, the action chosen by an optimistic (resp. pessimistic) rational agent is, according to the agent's optimistic (resp. pessimistic) assessment, at least as good as the other actions she may choose.

**Proposition 14** *Let  $i \in \text{Agt}$  and  $x \in \{Opt, Pess\}$ . Then,*

$$\models_{MCMC} (\text{Rat}_i^x \wedge \text{play}(i, a)) \rightarrow \bigwedge_{b \in \text{Act}} \text{RP}_i^x(\text{play}(i, b) \preceq \text{play}(i, a)) \quad (11)$$

**PROOF.** Let us prove the case  $x = Opt$ . Let  $M$  be a MCMC and let  $w$  be a world in  $M$ . Suppose  $M, w \models \text{Rat}_i^{Opt}$  and  $M, w \models \text{play}(i, a)$ . Thus,  $M, w \models \text{RP}_i^{Opt}(\neg \text{play}(i, a) \preceq \text{play}(i, a))$ . The latter means that  $\forall u \in \text{Best}_{i,P}(w) \cap \|\neg \text{play}(i, a)\|_{i,w,M}, \exists v \in \text{Best}_{i,P}(w) \cap \|\text{play}(i, a)\|_{i,w,M} : u \preceq_{i,D} v$ . Since  $\models_{MCMC} \text{play}(i, a) \rightarrow \neg \text{play}(i, b)$  if  $a \neq b$ , the latter implies  $\forall b \in \text{Act}, \forall u \in \text{Best}_{i,P}(w) \cap \|\text{play}(i, b)\|_{i,w,M}, \exists v \in \text{Best}_{i,P}(w) \cap \|\text{play}(i, a)\|_{i,w,M} : u \preceq_{i,D} v$ . Thus,  $\bigwedge_{b \in \text{Act}} \text{RP}_i^{Opt}(\text{play}(i, b) \preceq \text{play}(i, a))$ . The case  $x = Pess$  can be proved in an analogous way. ■

The following proposition elucidates the connection between the notions of belief, rationality and Nash equilibrium: if all agents are optimistic (resp. pessimistic) rational and have a correct belief about the others' actual choices, then the joint action they choose is a subjective optimistic (resp. pessimistic) Nash equilibrium.<sup>8</sup>

**Proposition 15** *Let  $x \in \{opt, pess\}$  and  $\delta \in \text{JAct}$ . Then:*

$$\models_{MCMC} \left( \text{play}(\delta) \wedge \bigwedge_{i \in \text{Agt}} (\text{Rat}_i^x \wedge \text{B}_i \text{play}(\delta_{-i})) \right) \rightarrow \text{NE}_i^x(\delta) \quad (12)$$

<sup>8</sup> A similar epistemic characterization of Nash equilibrium is provided by Aumann & Brandenburger (A&B) [6] in the context of games with complete information. See also [45] for a similar result using a probabilistic approach.

PROOF. Let  $M$  be a MCMC and let  $w$  be a world in  $M$ . Suppose  $M, w \models \text{play}(i, \delta(i))$ ,  $M, w \models \text{Rat}_i^x$  and  $M, w \models \text{B}_i \text{play}(\delta_{-i})$ , for all  $i \in \text{Agt}$ . By Proposition 14, it follows that  $\bigwedge_{b \in \text{Act}} \text{RP}_i^x(\text{play}(i, b) \preceq \text{play}(i, \delta(i)))$  and  $M, w \models \text{B}_i \text{play}(\delta_{-i})$ , for all  $i \in \text{Agt}$ . From the latter, we can conclude that  $M, w \models \text{BR}_i^x(\delta(i), \delta_{-i})$ , for all  $i \in \text{Agt}$ . Thus,  $M, w \models \text{NE}_i^x(\delta)$ . ■

We conclude this section by illustrating the game-theoretic concepts involved in the crossroad game described in Section 3.4. It is a game under incomplete information since an agent does not necessarily know the other agent's beliefs and desires. It is single-stage since that interaction is non-repeated and agents are supposed to choose simultaneously.

**Example (cont.)** Let us suppose that the set of actions that agents 1 and 2 can choose is  $\text{Act} = \{C, S\}$ , where  $C$  is the action “to continue” and  $S$  is the action “to stop”. The following hypotheses capture the agents' knowledge and beliefs about actions and their effects:

$$\begin{aligned} \varphi_4 =_{\text{def}} & \bigwedge_{i \in \{1, 2\}} [\equiv_i] \left( ((\text{play}(1, C) \wedge \text{play}(2, C)) \rightarrow \text{co}) \wedge \right. \\ & ((\text{play}(1, C) \wedge \text{play}(2, S)) \rightarrow (\neg lo_1 \wedge lo_2)) \wedge \\ & ((\text{play}(1, S) \wedge \text{play}(2, C)) \rightarrow (lo_1 \wedge \neg lo_2)) \wedge \\ & \left. ((\text{play}(1, S) \wedge \text{play}(2, S)) \rightarrow (lo_1 \wedge lo_2 \wedge \neg \text{co})) \right), \\ \varphi_5 =_{\text{def}} & \left( (\widehat{\text{B}}_1 \text{play}(1 \mapsto C, 2 \mapsto C) \leftrightarrow \widehat{\text{B}}_1 \text{play}(1 \mapsto S, 2 \mapsto C)) \wedge \right. \\ & (\widehat{\text{B}}_1 \text{play}(1 \mapsto C, 2 \mapsto S) \leftrightarrow \widehat{\text{B}}_1 \text{play}(1 \mapsto S, 2 \mapsto S)) \wedge \\ & (\widehat{\text{B}}_2 \text{play}(1 \mapsto C, 2 \mapsto C) \leftrightarrow \widehat{\text{B}}_2 \text{play}(1 \mapsto C, 2 \mapsto S)) \wedge \\ & \left. (\widehat{\text{B}}_2 \text{play}(1 \mapsto S, 2 \mapsto C) \leftrightarrow \widehat{\text{B}}_2 \text{play}(1 \mapsto S, 2 \mapsto S)) \right), \end{aligned}$$

where  $\widehat{\text{B}}_i \varphi =_{\text{def}} \neg \text{B}_i \neg \varphi$ . According to the hypothesis  $\varphi_4$ , the agents know that (i) if they both continue, they will collide, (ii) if one of them continues while the other stops, then the first will lose its time while the second will not, and (iii) if they both stop, each of them will lose its time but there will be no collision. According to the hypothesis  $\varphi_5$ , the fact that an agent considers possible that the other will decide to continue (resp. to stop) does not depend on the agent's choice. This hypothesis is justified by the assumption that an agent's beliefs are ex ante, i.e., relative to the instant before an agent makes its choice.

As the following validity indicates, the previous hypotheses  $\varphi_4$  and  $\varphi_5$  together with the hypotheses  $\varphi_1$  and  $\varphi_3$  stated in Section 3.4 lead to the conclusion that (i) an agent's action of continuing is both an optimistic and a pessimistic best response to the other agent's action of stopping, and an agent's action of stopping is both an optimistic and a pessimistic best response to the other agent's action of continuing. For every

$x \in \{opt, pess\}$ , we have:

$$\models_{MCMC} (\varphi_1 \wedge \varphi_3 \wedge \varphi_4 \wedge \varphi_5) \rightarrow (\text{BR}_1^x(S, 2 \mapsto C) \wedge \text{BR}_1^x(C, 2 \mapsto S) \wedge \text{BR}_2^x(S, 1 \mapsto C) \wedge \text{BR}_2^x(C, 1 \mapsto S)) \quad (13)$$

## 6 Dynamic Extension

The logics we presented so far merely provide a static picture of the cognitive attitudes and choices of multiple agents in interactive situations. Following the tradition of dynamic epistemic logic (DEL) [50], in this section we move from a static to a dynamic perspective and extend the language  $\mathcal{L}_{DLCA}$  by a variety of dynamic operators for cognitive attitude change. We consider two types of cognitive attitude change, namely, radical attitude and conservative attitude change. Radical attitude change, both in its epistemic and in its motivational form, satisfies a strong form of success postulate. Particularly, if an agent forms the belief that  $\varphi$ , as a consequence of a radical belief revision by  $\varphi$ , then she should also form the strong belief that  $\varphi$ . Analogously, if an agent forms the desire that  $\varphi$ , as a consequence of a radical desire revision by  $\varphi$ , then she should also form the strong desire that  $\varphi$ . On the contrary, after a conservative belief (resp. desire) revision by  $\varphi$  is performed, an agent may form the belief (resp. desire) that  $\varphi$  without forming the strong belief (resp. strong desire) that  $\varphi$ . While radical and conservative belief revision have been studied before in the literature on DEL [47, 10], we are the first to apply DEL techniques to the analysis of desire revision and to oppose belief revision to desire revision in the DEL setting.<sup>9</sup>

In the rest of this section, we first define the semantics of radical belief revision and desire revision operators (Section 6.1). Then, we turn to conservative attitude change and define the semantics of conservative belief revision and desire revision operators (Section 6.2). Finally, we provide an axiomatics for the dynamic extension of our logic DLCA (Section 6.3).

### 6.1 Radical Attitude Revision

Radical attitude revision operators are of the form  $[\uparrow_{i,\tau} \varphi]$ , with  $\tau \in \{P, D\}$ . They describe the consequences of a radical revision operation. In particular, the formula  $[\uparrow_{i,P} \varphi] \psi$  is meant to stand for “ $\psi$  holds, after agent  $i$  has radically revised her beliefs with  $\varphi$ ”, whereas  $[\uparrow_{i,D} \varphi] \psi$  is meant to stand for “ $\psi$  holds, after agent  $i$  has radically revised her desires with  $\varphi$ ”. We assume that radical revision operations are public, i.e., if an agent radically revises her beliefs (resp. desires) with  $\varphi$ , then this is common knowledge among all agents. This assumption could be easily relaxed by using action models as introduced in [8, 9], which would allow us to model private and semi-private attitude change operations. Radical revision operators are interpreted relative to

<sup>9</sup>Research in the DEL area has rather concentrated on preference change [49, 52], leaving desire change unexplored.

a MCM  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  and a world  $w$  in  $W$ , as follows:

$$M, w \models [\uparrow_{i,\tau} \varphi] \psi \iff M^{\uparrow_{i,\tau} \varphi}, w \models \psi,$$

where

$$M^{\uparrow_{i,P} \varphi} = (W, (\preceq_{i,P}^{\uparrow_{i,P} \varphi})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V),$$

$$M^{\uparrow_{i,D} \varphi} = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D}^{\uparrow_{i,D} \varphi})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V),$$

$$\preceq_{i,\tau}^{\uparrow_{i,\tau} \varphi} = \left\{ (w, v) \in W \times W : ((M, w \models \varphi \text{ iff } M, v \models \varphi) \text{ and } w \preceq_{i,\tau} v) \text{ or } (M, w \models \neg \varphi, M, v \models \varphi \text{ and } w \equiv_i v) \right\},$$

and  $\preceq_{j,\tau}^{\uparrow_{i,\tau} \varphi} = \preceq_{j,\tau}$  for all  $j \in \text{Agt}$  such that  $i \neq j$ .

Radical belief and desire revision are completely symmetric from the point of view of the plausibility and desirability ordering. Agent  $i$ 's radical belief revision with  $\varphi$  transforms agent  $i$ 's plausibility ordering  $\preceq_{i,P}$  into the new plausibility ordering  $\preceq_{i,P}^{\uparrow_{i,P} \varphi}$ . In particular, it makes all  $\varphi$ -worlds in  $i$ 's information set more plausible than all  $\neg\varphi$ -worlds and, within those two zones, it keeps the old plausibility ordering. Analogously, agent  $i$ 's radical desire revision with  $\varphi$  transforms agent  $i$ 's desirability ordering  $\preceq_{i,D}$  into the new desirability ordering  $\preceq_{i,D}^{\uparrow_{i,D} \varphi}$ . It makes all  $\varphi$ -worlds in  $i$ 's information set more desirable than all  $\neg\varphi$ -worlds and, within those two zones, it keeps the old desirability ordering.

As emphasized above, radical revision satisfies a strong form of success principle which is formally expressed by the following two validities. Let  $\varphi \in \mathcal{L}_{\text{PL}}(\text{Atm})$ . Then,

$$\models_{\text{MCM}} \langle \equiv_i \rangle \varphi \rightarrow [\uparrow_{i,P} \varphi] (\text{B}_i \varphi \wedge \text{SB}_i \varphi) \quad (14)$$

$$\models_{\text{MCM}} \langle \equiv_i \rangle \neg \varphi \rightarrow [\uparrow_{i,D} \varphi] (\text{D}_i \varphi \wedge \text{SD}_i \varphi) \quad (15)$$

This means that (i) if  $\varphi$  is compatible with an agent's knowledge then, after she has radically revised her beliefs with  $\varphi$ , the agent will both believe that  $\varphi$  and strongly believe that  $\varphi$ , and (ii) if  $\neg\varphi$  is compatible with an agent's knowledge then, after she has radically revised her desires with  $\varphi$ , the agent will both desire that  $\varphi$  and strongly desire that  $\varphi$ .

The two validities highlight that belief and desire behave in a slightly different way under radical revision, despite the fact that the plausibility and desirability ordering are modified in the same way.

We have the following additional validities, for  $\varphi \in \mathcal{L}_{\text{PL}}(\text{Atm})$ :

$$\models_{\text{MCM}} [\uparrow_{i,P} \varphi] (\text{B}_i \varphi \rightarrow \text{SB}_i \varphi) \quad (16)$$

$$\models_{\text{MCM}} [\uparrow_{i,D} \varphi] (\text{D}_i \varphi \rightarrow \text{SD}_i \varphi) \quad (17)$$

This means that the formation of a belief (resp. desire) through radical belief (resp. desire) revision necessarily entails the formation of a strong belief (resp. strong desire) with the same content.

**Example (cont.)** Let us go back to the crossroad game in order to illustrate the radical desire revision mechanism. Suppose agent 1 performs a radical desire revision operation with  $\neg lo_2$ , since it learns that agent 2 is an ambulance which has to lose no time at the crossroad. By the previous validity (15), we can prove that, under the hypothesis  $\varphi_2$  stated in Section 3.4, 1 will both desire and strongly desire that  $\neg lo_2$ , after the radical desire revision operation with  $\neg lo_2$ :

$$\models_{MCM} \varphi_2 \rightarrow [\uparrow_{1,D} \neg lo_2](D_1 \neg lo_2 \wedge SD_1 \neg lo_2) \quad (18)$$

Moreover, under the set of hypotheses  $\{\varphi_1, \varphi_2, \varphi_3\}$ , after the radical desire revision operation with  $\neg lo_2$ , 1 will not strongly desire anymore not to lose time, but it will merely desire it:

$$\models_{MCM} (\varphi_1 \wedge \varphi_2 \wedge \varphi_3) \rightarrow [\uparrow_{1,D} \neg lo_2](D_1 \neg lo_1 \wedge \neg SD_1 \neg lo_1) \quad (19)$$

As the following proposition indicates, we have reduction axioms which allow us to eliminate radical attitude revision operators from a formula.

**Proposition 16** *The following equivalences are valid:*

$$\begin{aligned} [\uparrow_{i,\tau} \varphi] p &\leftrightarrow p \\ [\uparrow_{i,\tau} \varphi] x &\leftrightarrow x \\ [\uparrow_{i,\tau} \varphi] \neg \psi &\leftrightarrow \neg [\uparrow_{i,\tau} \varphi] \psi \\ [\uparrow_{i,\tau} \varphi] (\psi_1 \wedge \psi_2) &\leftrightarrow ([\uparrow_{i,\tau} \varphi] \psi_1 \wedge [\uparrow_{i,\tau} \varphi] \psi_2) \\ [\uparrow_{i,\tau} \varphi] [\pi] \psi &\leftrightarrow [F^{\uparrow_{i,\tau} \varphi}(\pi)] [\uparrow_{i,\tau} \varphi] \psi \end{aligned}$$

where for all  $j \in \text{Agt}$  and for all  $\tau, \tau' \in \{P, D\}$ :

$$\begin{aligned} F^{\uparrow_{i,\tau} \varphi}(\equiv_j) &= \equiv_j \\ F^{\uparrow_{i,\tau} \varphi}(\preceq_{i,\tau}) &= (\varphi?; \preceq_{i,\tau}; \varphi?) \cup (\neg \varphi?; \preceq_{i,\tau}; \neg \varphi?) \cup (\neg \varphi?; \equiv_i; \varphi?) \\ F^{\uparrow_{i,\tau} \varphi}(\preceq_{j,\tau'}) &= \preceq_{j,\tau'} \text{ if } i \neq j \text{ or } \tau \neq \tau' \\ F^{\uparrow_{i,\tau} \varphi}(\preceq_{i,\tau}^\sim) &= (\varphi?; \preceq_{i,\tau}^\sim; \varphi?) \cup (\neg \varphi?; \preceq_{i,\tau}^\sim; \neg \varphi?) \cup (\varphi?; \equiv_i; \neg \varphi?) \\ F^{\uparrow_{i,\tau} \varphi}(\preceq_{j,\tau'}^\sim) &= \preceq_{j,\tau'}^\sim \text{ if } i \neq j \text{ or } \tau \neq \tau' \\ F^{\uparrow_{i,\tau} \varphi}(\pi; \pi') &= F^{\uparrow_{i,\tau} \varphi}(\pi); F^{\uparrow_{i,\tau} \varphi}(\pi') \\ F^{\uparrow_{i,\tau} \varphi}(\pi \cup \pi') &= F^{\uparrow_{i,\tau} \varphi}(\pi) \cup F^{\uparrow_{i,\tau} \varphi}(\pi') \\ F^{\uparrow_{i,\tau} \varphi}(\pi \cap \pi') &= F^{\uparrow_{i,\tau} \varphi}(\pi) \cap F^{\uparrow_{i,\tau} \varphi}(\pi') \\ F^{\uparrow_{i,\tau} \varphi}(-\pi) &= -F^{\uparrow_{i,\tau} \varphi}(\pi) \\ F^{\uparrow_{i,\tau} \varphi}(\psi?) &= [\uparrow_{i,\tau} \varphi] \psi? \end{aligned}$$

## 6.2 Conservative Attitude Revision

Let us move from radical attitude change to conservative attitude change by introducing radical revision operators of type  $[\uparrow_{i,\tau} \varphi]$ , with  $\tau \in \{P, D\}$ . The formula  $[\uparrow_{i,P} \varphi] \psi$  (resp.  $[\uparrow_{i,D} \varphi] \psi$ ) is meant to stand for “ $\psi$  holds, after agent  $i$  has conservatively revised her beliefs (resp. desires) with  $\varphi$ ”. As for radical revision, we assume that conservative revision operations are public, i.e., if an agent conservatively revises her beliefs (resp. desires) with  $\varphi$ , then this is common knowledge among all agents. The semantic interpretation of such operators relative to a MCM  $M = (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V)$  be a MCM and a world  $w$  in  $W$  is as follows:

$$M, w \models [\uparrow_{i,\tau} \varphi] \psi \iff M^{\uparrow_{i,\tau} \varphi}, w \models \psi,$$

where:

$$\begin{aligned} M^{\uparrow_{i,P} \varphi} &= (W, (\preceq_{i,P}^{\uparrow_{i,P} \varphi})_{i \in \text{Agt}}, (\preceq_{i,D})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V), \\ M^{\uparrow_{i,D} \varphi} &= (W, (\preceq_{i,P})_{i \in \text{Agt}}, (\preceq_{i,D}^{\uparrow_{i,D} \varphi})_{i \in \text{Agt}}, (\equiv_i)_{i \in \text{Agt}}, V), \end{aligned}$$

with:

$$\begin{aligned} \preceq_{i,P}^{\uparrow_{i,P} \varphi} &= \left\{ (w, v) \in W \times W : \left( (w \in \text{Best}_{i,P}(\varphi, w) \text{ iff } v \in \text{Best}_{i,P}(\varphi, w)) \text{ and } w \preceq_{i,P} v \right) \text{ or } \right. \\ &\quad \left. (w \notin \text{Best}_{i,P}(\varphi, w), v \in \text{Best}_{i,P}(\varphi, w) \text{ and } w \equiv_i v) \right\}, \\ \preceq_{i,D}^{\uparrow_{i,D} \varphi} &= \left\{ (w, v) \in W \times W : \left( (w \in \text{Worst}_{i,D}(\neg\varphi, w) \text{ iff } v \in \text{Worst}_{i,D}(\neg\varphi, w)) \text{ and } w \preceq_{i,D} v \right) \text{ or } \right. \\ &\quad \left. (w \in \text{Worst}_{i,D}(\neg\varphi, w), v \notin \text{Worst}_{i,D}(\neg\varphi, w) \text{ and } w \equiv_i v) \right\}, \end{aligned}$$

and  $\preceq_{j,\tau}^{\uparrow_{i,\tau} \varphi} = \preceq_{j,\tau}$  for all  $j \in \text{Agt}$  such that  $i \neq j$ .

Unlike radical revision, plausibility update and desirability update are asymmetric under conservative revision. Agent  $i$ 's conservative belief revision with  $\varphi$  replaces the current plausibility ordering  $\preceq_{i,P}$  with the new plausibility ordering  $\preceq_{i,P}^{\uparrow_{i,P} \varphi}$ . It promotes the most plausible  $\varphi$ -worlds to the highest plausibility rank, but apart from that, the old plausibility ordering remains. Agent  $i$ 's conservative desire revision with  $\varphi$  replaces the current desirability ordering  $\preceq_{i,D}$  with the new desirability ordering  $\preceq_{i,D}^{\uparrow_{i,D} \varphi}$ . In particular, it demotes the least desirable  $\neg\varphi$ -worlds to the lowest desirability rank, but apart from that, the old desirability ordering remains.

Conservative attitude revision satisfies a weak form of success principle which guarantees the formation of a belief (resp. a desire), after a belief (resp. desire) revision is performed. Let  $\varphi, \mathcal{L}_{\text{PL}}(\text{Atm})$ . Then,

$$\models_{\text{MCM}} \neg \text{B}_i(\varphi, \perp) \rightarrow [\uparrow_{i,P} \varphi] \text{B}_i \varphi \quad (20)$$

$$\models_{\text{MCM}} \neg \text{D}_i(\varphi, \top) \rightarrow [\uparrow_{i,D} \varphi] \text{D}_i \varphi \quad (21)$$

According to the previous validities, if an agent does not believe a contradiction conditional on  $\varphi$  then, after she has conservatively revised her beliefs with  $\varphi$ , she will believe that  $\varphi$ . If an agent does not desire a tautology conditional on  $\varphi$  then, after she

has conservatively revised her desires with  $\varphi$ , she will desire that  $\varphi$ . But, unlike radical attitude revision, conservative attitude revision does not necessarily guarantee the formation of a strong belief (resp. a strong desire), after a belief (resp. desire) revision is performed. Indeed, we have the following, for  $\varphi \in \mathcal{L}_{PL}(Atm)$ :

$$\not\models_{MCM} [\uparrow_{i,P} \varphi](B_i \varphi \rightarrow SB_i \varphi) \quad (22)$$

$$\not\models_{MCM} [\uparrow_{i,D} \varphi](D_i \varphi \rightarrow SD_i \varphi) \quad (23)$$

This means that the formation of a belief (resp. desire) through conservative belief (resp. desire) revision does not necessarily entail the formation of a strong belief (resp. strong desire) with the same content.

**Example (cont.)** *Let us illustrate the conservative belief revision mechanism with the help of the crossroad game. Suppose agent 1 informs agent 2 that “if they both lose time, then there will no collision” and, as a consequence, 2 performs a conservative belief revision operation with input  $(lo_1 \wedge lo_2) \rightarrow \neg co$ . By the previous validity (20) we can prove that, under the hypothesis  $\varphi_1$  stated in Section 3.4 and the assumption that 2 does not believe a contradiction conditional on 1’s assertion, 2 believes that there will be no collision, after its conservative belief operation:*

$$\models_{MCM} \left( \neg B_2((lo_1 \wedge lo_2) \rightarrow \neg co, \perp) \wedge \varphi_1 \right) \rightarrow [\uparrow_{2,P} (lo_1 \wedge lo_2) \rightarrow \neg co] B_2 \neg co \quad (24)$$

As for radical revision, we have reduction axioms which allow us to eliminate conservative attitude revision operators from a formula.

**Proposition 17** *The following equivalences are valid:*

$$\begin{aligned} [\uparrow_{i,\tau} \varphi] p &\leftrightarrow p \\ [\uparrow_{i,\tau} \varphi] x &\leftrightarrow x \\ [\uparrow_{i,\tau} \varphi] \neg \psi &\leftrightarrow \neg [\uparrow_{i,\tau} \varphi] \psi \\ [\uparrow_{i,\tau} \varphi] (\psi_1 \wedge \psi_2) &\leftrightarrow ([\uparrow_{i,\tau} \varphi] \psi_1 \wedge [\uparrow_{i,\tau} \varphi] \psi_2) \\ [\uparrow_{i,\tau} \varphi] [\pi] \psi &\leftrightarrow [F^{\uparrow_{i,\tau} \varphi}(\pi)] [\uparrow_{i,\tau} \varphi] \psi \end{aligned}$$

where for all  $j \in \text{Agt}$  and for all  $\tau, \tau' \in \{P, D\}$ :

$$\begin{aligned}
F^{\uparrow i, \tau \varphi}(\equiv_j) &= \equiv_j \\
F^{\uparrow i, P \varphi}(\preceq_{i, P}) &= ((\varphi \wedge [\prec_{i, P}] \neg \varphi)?; \preceq_{i, \tau}; (\varphi \wedge [\prec_{i, P}] \neg \varphi)?) \cup \\
&\quad (\neg(\varphi \wedge [\prec_{i, P}] \neg \varphi)?; \preceq_{i, \tau}; \neg(\varphi \wedge [\prec_{i, P}] \neg \varphi)?) \cup \\
&\quad (\neg(\varphi \wedge [\prec_{i, P}] \neg \varphi)?; \equiv_i; (\varphi \wedge [\prec_{i, P}] \neg \varphi)?) \\
F^{\uparrow i, D \varphi}(\preceq_{i, D}) &= ((\neg \varphi \wedge [\succ_{i, D}] \varphi)?; \preceq_{i, \tau}; (\neg \varphi \wedge [\succ_{i, D}] \varphi)?) \cup \\
&\quad (\neg(\neg \varphi \wedge [\succ_{i, D}] \varphi)?; \preceq_{i, \tau}; \neg(\neg \varphi \wedge [\succ_{i, D}] \varphi)?) \cup \\
&\quad ((\neg \varphi \wedge [\succ_{i, D}] \varphi)?; \equiv_i; \neg(\neg \varphi \wedge [\succ_{i, D}] \varphi)?) \\
F^{\uparrow i, \tau \varphi}(\preceq_{j, \tau'}) &= \preceq_{j, \tau'} \text{ if } i \neq j \text{ or } \tau \neq \tau' \\
F^{\uparrow i, P \varphi}(\preceq_{i, P}^{\sim}) &= ((\varphi \wedge [\prec_{i, P}] \neg \varphi)?; \preceq_{i, \tau}; (\varphi \wedge [\prec_{i, P}] \neg \varphi)?) \cup \\
&\quad (\neg(\varphi \wedge [\prec_{i, P}] \neg \varphi)?; \preceq_{i, \tau}; \neg(\varphi \wedge [\prec_{i, P}] \neg \varphi)?) \cup \\
&\quad ((\varphi \wedge [\prec_{i, P}] \neg \varphi)?; \equiv_i; \neg(\varphi \wedge [\prec_{i, P}] \neg \varphi)?) \\
F^{\uparrow i, D \varphi}(\preceq_{i, D}^{\sim}) &= ((\neg \varphi \wedge [\succ_{i, D}] \varphi)?; \preceq_{i, \tau}; (\neg \varphi \wedge [\succ_{i, D}] \varphi)?) \cup \\
&\quad (\neg(\neg \varphi \wedge [\succ_{i, D}] \varphi)?; \preceq_{i, \tau}; \neg(\neg \varphi \wedge [\succ_{i, D}] \varphi)?) \cup \\
&\quad (\neg(\neg \varphi \wedge [\succ_{i, D}] \varphi)?; \equiv_i; \neg(\neg \varphi \wedge [\succ_{i, D}] \varphi)?) \\
F^{\uparrow i, \tau \varphi}(\preceq_{j, \tau'}^{\sim}) &= \preceq_{j, \tau'}^{\sim} \text{ if } i \neq j \text{ or } \tau \neq \tau' \\
F^{\uparrow i, \tau \varphi}(\pi; \pi') &= F^{\uparrow i, \tau \varphi}(\pi); F^{\uparrow i, \tau \varphi}(\pi') \\
F^{\uparrow i, \tau \varphi}(\pi \cup \pi') &= F^{\uparrow i, \tau \varphi}(\pi) \cup F^{\uparrow i, \tau \varphi}(\pi') \\
F^{\uparrow i, \tau \varphi}(\pi \cap \pi') &= F^{\uparrow i, \tau \varphi}(\pi) \cap F^{\uparrow i, \tau \varphi}(\pi') \\
F^{\uparrow i, \tau \varphi}(-\pi) &= -F^{\uparrow i, \tau \varphi}(\pi) \\
F^{\uparrow i, \tau \varphi}(\psi?) &= [\uparrow_{i, \tau} \varphi] \psi?
\end{aligned}$$

### 6.3 Dynamic Logic of Cognitive Attitudes and their Change

The modal language  $\mathcal{L}_{\text{DLCAC}}(\text{Atm}, \text{Nom}, \text{Agt})$ , or simply  $\mathcal{L}_{\text{DLCAC}}$ , for the Dynamic Logic of Cognitive Attitudes and their Change (DLCAC) extends the language  $\mathcal{L}_{\text{DLCA}}$  of the logic DLCA by dynamic operators of type  $[\uparrow_{i, \tau} \varphi]$  and  $[\uparrow_{i, \tau} \varphi]$ . It is defined by the following grammar:

$$\varphi ::= p \mid x \mid \neg \varphi \mid \varphi \wedge \varphi' \mid [\pi] \varphi \mid [\uparrow_{i, \tau} \varphi] \psi \mid [\uparrow_{i, \tau} \varphi] \psi$$

where  $\pi$  ranges over the language of cognitive programs  $\mathcal{P}$ ,  $p$  ranges over  $\text{Atm}$ ,  $x$  ranges over  $\text{Nom}$ ,  $i$  ranges over  $\text{Agt}$  and  $\tau$  ranges over  $\{P, D\}$ .

**Definition 20** We define DLCAC to be the extension of DLCA given by the reduction principles of Proposition 16 and Proposition 17 and the following rule of replacement of equivalents

$$\frac{\psi_1 \leftrightarrow \psi_2}{\varphi \leftrightarrow \varphi[\psi_1/\psi_2]} \quad (\text{REP})$$



where  $\varphi[\psi_1/\psi_2]$  is the formula that results from  $\varphi$  by replacing zero or more occurrences of  $\psi_1$ , in  $\varphi$ , by  $\psi_2$ .

As the rule of replacement of equivalents preserves validity, the equivalences of Propositions 16 and 17 together with this allow to reduce every formula of the language  $\mathcal{L}_{\text{DLCAC}}$  to an equivalent formula of the language  $\mathcal{L}_{\text{DLCA}}$ . 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 18** *Let  $\varphi$  be a formula in the language of  $\mathcal{L}_{\text{DLCAC}}$ . Then*

- *red( $\varphi$ ) has no dynamic operators  $[\uparrow_{i,\tau} \varphi]$  or  $[\uparrow_{i,\tau} \varphi]$ , and*
- *red( $\varphi$ )  $\leftrightarrow \varphi$  is valid relative to the class of MCMs.*

The first item of Proposition 18 is clear. The second item is proved using the equivalences of Propositions 16 and 17 and the rule of replacement of equivalents.

The following theorem is a direct consequence of Theorem 1 and Proposition 18.

**Theorem 3** *The logic DLCA is sound and complete for the class of multi-agent cognitive models.*

## 7 Conclusion and perspectives

We have presented a logical framework for modelling a rich variety of cognitive attitudes of both epistemic type and motivational type. We have presented two extensions of the basic setting, one by the notion of choice and the other by dynamic operators for belief change and desire change. We have applied the former to the analysis of games under incomplete information. We have provided sound and complete axiomatizations for the basic setting and for its two extensions. Directions of future research are manifold and are briefly discussed in the rest of this section.

**Decidability and complexity** The present paper is devoted to study the proof-theoretic aspects of the proposed logics. In future work, we plan to investigate their computational aspects including decidability of their satisfiability problems and, at a later stage, complexity. In order to prove decidability, we expect to be able to use existing filtration techniques from modal logic. Note that once we have proved decidability of the static setting DLCA, we can use the reduction axioms of Propositions 16 and 17 to prove decidability of the dynamic setting DLCAC.

We plan to study complexity of the satisfiability problems for interesting fragments of the language  $\mathcal{L}_{\text{DLCA}}$  by reducing them to satisfiability problems of existing logics. For instance, consider the following single-agent (*sa*) fragment of the language  $\mathcal{L}_{\text{DLCA}}$  where only atomic programs (*ap*) are allowed, noted  $\mathcal{L}_{\text{DLCA}}^{\text{sa,ap}}$ :

$$\varphi ::= p \mid x \mid \neg\varphi \mid \varphi \wedge \varphi' \mid [\preceq_{1,p}]\varphi \mid [\preceq_{1,D}]\varphi \mid [\equiv_1]\varphi$$

where 1 is an arbitrary agent in  $Agt$ . We can observe that the satisfiability problem for this fragment is EXPTIME-hard. Indeed, because of Constraint C1 in Definition 1, the modality  $[\equiv_1]$  plays the role of the universal modality with respect to the modalities  $[\preceq_{1,P}]$  and  $[\preceq_{1,D}]$ . As shown in [27], adding the universal modality to a multimodal logic with independent modalities, such as the S4-modalities  $[\preceq_{1,P}]$  and  $[\preceq_{1,D}]$ , causes EXPTIME-hardness.

Consider now the following intersection-free (*if*) and complement-free (*cf*) fragment of  $\mathcal{L}_{DLCA}$ , noted  $\mathcal{L}_{DLCA}^{if,cf}$ :

$$\begin{aligned}\pi &::= \equiv_i | \preceq_{i,P} | \preceq_{i,D} | \pi; \pi' | \pi \cup \pi' | \neg \pi | \varphi? \\ \varphi &::= p | x | \neg \varphi | \varphi \wedge \varphi' | [\pi] \varphi\end{aligned}$$

Our first conjecture is that we can find a polysize reduction of the satisfiability problem for  $\mathcal{L}_{DLCA}^{if,cf}$  to the satisfiability problem of converse propositional dynamic logic (PDL) with nominals, also called converse combinatory propositional dynamic logic (CcPDL).<sup>10</sup> The latter problem is known to be EXPTIME-complete [21]. Therefore, if our conjecture is true, we will be able to conclude that the satisfiability problems for the fragments  $\mathcal{L}_{DLCA}^{sa,ap}$  and  $\mathcal{L}_{DLCA}^{if,cf}$  are both EXPTIME-complete.

We also intend to study complexity of the nominal-free (*nf*) fragment of  $\mathcal{L}_{DLCA}$ , noted  $\mathcal{L}_{DLCA}^{nf}$ . Nominals play a technical role in the logic DLCA by making it easier the task of axiomatizing intersection and complement of programs (Axioms **Add2** $_{\cap}$  and **Comp1** $_{\sim}$  in Definition 13). Our second conjecture is that the language  $\mathcal{L}_{DLCA}$  is strictly more expressive than its nominal-free fragment  $\mathcal{L}_{DLCA}^{nf}$ . Our third conjecture is that we can find a polysize reduction of the satisfiability problem for  $\mathcal{L}_{DLCA}^{nf}$  to the satisfiability problem of boolean modal logic with a bounded number of modal parameters which is known to be EXPTIME-complete [36]. We leave the proof of the previous three conjectures to future work. We leave to future work (i) the proof of the previous three conjectures, and (ii) the development of tableau-based automated reasoning procedures for the language  $\mathcal{L}_{DLCA}$  and for its fragments  $\mathcal{L}_{DLCA}^{sa,ap}$ ,  $\mathcal{L}_{DLCA}^{if,cf}$  and  $\mathcal{L}_{DLCA}^{nf}$  which can be used for programming artificial agents endowed with cognitive attitudes.

**Well-foundedness** Future work will also be devoted to study a variant of our logic DLCA under the assumption of converse well-foundedness for the relation  $\preceq_{i,P}$  and well-foundedness for the relation  $\preceq_{i,D}$ . As emphasized in Section 3, these properties are required to make agents' beliefs and desires consistent, namely, to guarantee that the formulas  $\neg(B_i \varphi \wedge B_i \neg \varphi)$ ,  $\neg B_i \perp$ ,  $\neg(D_i \varphi \wedge D_i \neg \varphi)$  and  $\neg D_i \top$  become valid. We will define the logic  $DLCA^{wf}$  to be the extension of the logic DLCA of Definition 13 by the following two axioms:

$$\begin{aligned}\langle \equiv_i \rangle \psi &\rightarrow \langle \equiv_i \rangle (\psi \wedge [\prec_{i,P}] \neg \psi) & (\mathbf{CWF}_{\prec_{i,P}}) \\ \langle \equiv_i \rangle \psi &\rightarrow \langle \equiv_i \rangle (\psi \wedge [\succ_{i,D}] \neg \psi) & (\mathbf{WF}_{\succ_{i,D}})\end{aligned}$$

<sup>10</sup> The main idea of the polynomial embedding is to exploit the iteration construct  $*$  of PDL for the translation  $tr$  of the cognitive programs, by stipulating that  $tr(\equiv_i) = (any_i \cup \neg any_i)^*$ ,  $tr(\preceq_{i,P}) = P_i^*$ ,  $tr(\preceq_{i,D}) = D_i^*$ , and homomorphic otherwise, where  $\mathcal{A}_i$  is agent  $i$ 's set of atomic programs (or actions),  $\mathcal{A} = \bigcup_{i \in Agt} \mathcal{A}_i$  is the set of PDL atomic programs,  $any_i = \bigcup_{a_i \in \mathcal{A}_i} a_i$  and, finally,  $P_i^*$  and  $D_i^*$  are special atomic programs in  $\mathcal{A}_i$ .

Such axioms are variants of the so-called Gödel-Löb (**GL**) axiom from provability logic [14]. Our conjecture is that the logic  $DLCA^{wf}$  so defined is sound and complete for the class of multi-agent cognitive models (MCMs) whose relations  $\preceq_{i,D}$  and  $\preceq_{i,P}$  are, respectively, well-founded and conversely well-founded.

**Ceteris paribus preference** We also plan to study a *ceteris paribus* notion of dyadic preference in the sense of Von Wright [55], which has been recently formalized in a modal logic setting by van Benthem et al. [48]. According to Von Wright, for an agent to have a preference of  $\varphi$  over  $\psi$ , she should prefer a situation in which  $\varphi$  is true to a situation in which  $\psi$  is true, *all other things being equal*.<sup>11</sup> Our aim is to show that the DLCA framework is expressive enough to capture both the static and the dynamic aspects of this notion of *ceteris paribus* preference.

## Acknowledgments

This work was supported by the ANR project CoPains (“Cognitive Planning in Persuasive Multimodal Communication”). Support from the ANR-3IA Artificial and Natural Intelligence Toulouse Institute is also gratefully acknowledged.

## References

- [1] C. E. Alchourrón, P. Gardenfors, and D. Makinson. On the logic of theory change: Partial meet contraction and revision functions. *The Journal of Symbolic Logic*, 50:510–530, 1985.
- [2] N. B. Amor, H. Fargier, R. Sabbadin, and M. Trabelsi. Possibilistic games with incomplete information. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI 2019)*, pages 1544–1550, 2019.
- [3] A. J. Anglberger, N. Gratzl, and O. Roy. Obligation, free choice, and the logic of weakest permissions. *The Review of Symbolic Logic*, 8:807–827, 2015.
- [4] G. Aucher. A combined system for update logic and belief revision. In *Intelligent Agents and Multi-Agent Systems, 7th Pacific Rim International Workshop on Multi-Agents (PRIMA 2004)*, volume 3371 of *LNCS*, pages 1–17. Springer, 2005.
- [5] R. Aumann. Interactive epistemology I: Knowledge. *International Journal of Game Theory*, 28(3):263–300, 1999.
- [6] R. Aumann and A. Brandenburger. Epistemic conditions for Nash equilibrium. *Econometrica*, 63:1161–1180, 1995.
- [7] P. Balbiani, A. Herzig, and N. Troquard. Alternative axiomatics and complexity of deliberative stit theories. *Journal of Philosophical Logic*, 37(4):387–406, 2008.

---

<sup>11</sup>See also [53] for a “ceteris paribus” interpretation of the notion of goal.

- [8] A. Baltag, L. Moss, and S. Solecki. The logic of public announcements, common knowledge and private suspicions. In *Proceedings of TARK'98*, pages 43–56. Morgan Kaufmann, 1998.
- [9] A. Baltag and L. S. Moss. Logics for epistemic programs. *Synthese*, 139(2):165–224, 2004.
- [10] A. Baltag and S. Smets. A qualitative theory of dynamic interactive belief revision. In *Proceedings of LOFT 7*, volume 3 of *Texts in Logic and Games*, pages 13–60. Amsterdam University Press, 2008.
- [11] A. Baltag and S. Smets. Talking your way into agreement: Belief merge by persuasive communication. In *Proceedings of the Second Multi-Agent Logics, Languages, and Organisations Federated Workshops (MALLOW)*, volume 494. CEUR, 2009.
- [12] N. Belnap, M. Perloff, and M. Xu. *Facing the future: agents and choices in our indeterminist world*. Oxford University Press, 2001.
- [13] S. Benferhat, D. Dubois, H. Prade, and M.-A. Williams. A practical approach to revising prioritized knowledge bases. *Studia Logica*, 70:105–130, 2002.
- [14] G. Boolos. *The Logic of Provability*. Cambridge University Press, 1993.
- [15] C. Boutilier. Revision sequences and nested conditionals. In *Proceedings of the 13th International Joint Conference on Artificial Intelligence (IJCAI 1993)*, pages 519–525. Morgan Kaufmann, 1993.
- [16] C. Boutilier. Towards a logic for qualitative decision theory. In *Proceedings of International Conference on Principles of Knowledge Representation and Reasoning (KR' 94)*, pages 75–86. AAAI Press, 1994.
- [17] R. I. Brafman and M. Tennenholtz. On the foundations of qualitative decision theory. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence (AAAI'96)*, pages 1291–1296. AAAI Press, 1996.
- [18] R. I. Brafman and M. Tennenholtz. An axiomatic treatment of three qualitative decision criteria. *Journal of the ACM*, 47(3):452–482, 2000.
- [19] P. R. Cohen and H. J. Levesque. Intention is choice with commitment. *Artificial Intelligence*, 42:213–261, 1990.
- [20] A. Darwiche and J. Pearl. On the logic of iterated belief revision. *Artificial Intelligence*, 89(1-2):1–29, 1997.
- [21] G. De Giacomo. *Decidability of Class-Based Knowledge Representation Formalisms*. PhD thesis, Università di Roma “La Sapienza”, 1995.
- [22] J. Doyle and R. Thomason. Background to qualitative decision theory. *The AI Magazine*, 20(2):55–68, 1999.

- [23] D. Dubois, E. Lorini, and H. Prade. The strength of desires: a logical approach. *Minds and Machines*, 27(1):199–231, 2017.
- [24] R. Fagin, J. Halpern, Y. Moses, and M. Vardi. *Reasoning about Knowledge*. MIT Press, Cambridge, Massachusetts, 1995.
- [25] G. Gargov and V. Goranko. Modal logic with names. *Journal of Philosophical Logic*, 22:607–636, 1993.
- [26] D. Harel, D. Kozen, and J. Tiuryn. *Dynamic Logic*. MIT Press, Cambridge, Massachusetts, 2000.
- [27] E. Hemaspaandra. The price of universality. *Notre Dame Journal of Formal Logic*, 37(2):174–203, 1996.
- [28] J. Hintikka. *Knowledge and belief: an introduction to the logic of the two notions*. Cornell University Press, 1962.
- [29] I. L. Humberstone. Direction of fit. *Mind*, 101(401):59–83, 1992.
- [30] T. F. Icard, E. Pacuit, and Y. Shoham. Joint revision of beliefs and intention. In *Proceedings of the Twelfth International Conference on Principles of Knowledge Representation and Reasoning (KR 2010)*, pages 572–574. AAAI Press, 2010.
- [31] D. Lewis. A problem about permission. In *Essays in honour of Jaakko Hintikka*, pages 163–175. 1979.
- [32] F. Liu. *Reasoning about Preference Dynamics*. Springer, 2011.
- [33] E. Lorini. Temporal STIT logic and its application to normative reasoning. *Journal of Applied Non-Classical Logics*, 23(4):372–399, 2013.
- [34] E. Lorini. Logics for games, emotions and institutions. *If-CoLog Journal of Logics and their Applications*, 4(9):3075–3113, 2017.
- [35] E. Lorini. Reasoning about cognitive attitudes in a qualitative setting. In *Proceedings of the 16th European Conference on Logics in Artificial Intelligence (ECAI 2019)*, volume 11468 of *LNCS*, pages 726–743. Springer, 2019.
- [36] C. Lutz and U. Sattler. The complexity of reasoning with boolean modal logics. In *Proceedings of the Third Conference on Advances in Modal logic (AiML 3)*, pages 329–348. World Scientific, 2000.
- [37] K. Marsh and H. Wallace. The influence of attitudes on beliefs: Formation and change. In D. Albarracin, B. T. Johnson, and M. P. Zanna, editors, *The Handbook of Attitudes*, pages 369–395. Lawrence Erlbaum Ass., 2005.
- [38] J. J. Ch. Meyer, W. van der Hoek, and B. van Linder. A logical approach to the dynamics of commitments. *Artificial Intelligence*, 113(1-2):1–40, 1999.
- [39] S. Passy and T. Tinchev. An essay in combinatorial dynamic logic. *Information and Computation*, 93:263–332, 1991.

- [40] C. Paternotte. Rational choice theory. In I. Jarvie and J. Zamora-Bonilla, editors, *SAGE Handbook for the Philosophy of Social Sciences*, pages 307–321. SAGE Publications Inc., 2011.
- [41] M. Platts. *Ways of meaning*. Routledge, and Kegan Paul, 1979.
- [42] H. Rott. Shifting priorities: Simple representations for 27 iterated theory change operators. In D. Makinson, J. Malinowski, and H. Wansing, editors, *Towards Mathematical Philosophy: Papers from the Studia Logica conference Trends in Logic IV*, pages 269–296. Springer, 2009.
- [43] J. Searle. *Expression and meaning*. Cambridge University Press, 1979.
- [44] Y. Shoham. Logical theories of intention and the database perspective. *Journal of Philosophical Logic*, 38(6):633–647, 2009.
- [45] W. Spohn. How to make sense of game theory. In *Philosophy of Economics*, volume 2, pages 239–270. 1982.
- [46] W. Spohn. Ordinal conditional functions: a dynamic theory of epistemic states. In *Causation in decision, belief change and statistics*, pages 105–134. Kluwer, 1988.
- [47] J. van Benthem. Dynamic logic for belief revision. *Journal of Applied Non-Classical Logics*, 17(2):129–155, 2007.
- [48] J. van Benthem, P. Girard, and O. Roy. Everything else being equal: A modal logic for ceteris paribus preferences. *Journal of Philosophical Logic*, 38:83–125, 2009.
- [49] J. van Benthem and F. Liu. Dynamic logic of preference upgrade. *Journal of Applied Non-Classical Logics*, 17(2):157–182, 2007.
- [50] H. van Ditmarsch, W. van der Hoek, and B. Kooi. *Dynamic epistemic logic*, volume 337. Synthese Library, Springer, 2007.
- [51] H. P. van Ditmarsch. Prolegomena to dynamic logic for belief revision. *Synthese*, 147(2):229–275, 2005.
- [52] J. van Eijck. Yet more modal logics of preference change and belief revision. In *New Perspectives on Games and Interaction*, volume 4 of *Texts in Logic and Games*, pages 81–104. Amsterdam University Press, 2008.
- [53] M. P. Wellman and J. Doyle. Preferential semantics for goals. In *Proceedings of the Ninth National conference on Artificial intelligence (AAAI’91)*, pages 698–703, 1991.
- [54] M. Wooldridge. *Reasoning about rational agents*. MIT Press, Cambridge, 2000.
- [55] G. H. Von Wright. *The logic of preference*. Edinburgh University Press, 1963.
- [56] G. H. Von Wright. The logic of preference reconsidered. *Theory and Decision*, 3:140–169, 1972.