

Chapter 10

BDI Logics

John-Jules Ch. Meyer, Jan Broersen and Andreas Herzig

Contents

10.1 Introduction	453
10.2 Bratman’s theory of Belief-Desire-Intention . .	454
10.3 KARO Logic	471
10.4 BDI-modalities in STIT logic	479
10.5 Conclusion	488
10.6 Notes	489
References	493

Abstract This paper presents an overview of so-called BDI logics, logics where the notion of Beliefs, Desires and Intentions play a central role. Starting out from the basic ideas about BDI by Bratman, we consider various formalizations in logic, such as the approach of Cohen and Levesque, slightly remodelled in dynamic logic, Rao & Georgeff’s influential BDI logic based on the branching-time temporal logic CTL*, the KARO framework and BDI logic based on STIT (seeing to it that) logics.

10.1 Introduction

In this chapter we present an overview of so-called BDI (for Beliefs, Desires, and Intentions) logics, that is, logics that describe the mental attitudes of intelligent agents in terms of folk-psychological notions of beliefs, desires and intentions. This theory is based on the work of the philosopher Michael Bratman — as with all chapters in his book, references to the literature are provided in the final section. The chapter is organized as follows: we start with some of the basic ideas in Bratman’s philosophy, which is about practical reasoning (the reasoning about

performing actions) on the basis of the agent's beliefs, desires and, very importantly, intentions, which are special desires to which the agent is committed. Then a number of formalizations of BDI theory in logic is reviewed, the so-called BDI logics. Starting out with Cohen & Levesque's approach, slightly reworked in a dynamic logic setting by Herzig and colleagues. Then we look at Rao & Georgeff's BDI logic, based on the branching-time temporal logic CTL*. Next we discuss the KARO framework which is based on dynamic logic since its conception. We sketch as a small excursion how the KARO framework, which is devised to capture the behaviour of (rational) intelligent agents, also can be used for describing emotional behaviour of agents. We then present a relatively new approach to BDI, based on so-called STIT (seeing to it that) logic. We finally round off with a conclusion section and a section containing the pointers to all bibliographical references.

10.2 Bratman's theory of Belief-Desire-Intention

"What happens to our conception of mind and rational agency when we take seriously future-directed intentions and plans and their roles as inputs into further practical reasoning? This question drives much of this book."

This is how the preface of Michael E. Bratman's famous book "Intention, Plans, and Practical Reason" starts. In this book the author lays down the foundations of what later would be called the BDI (Belief-Desire-Intention) theory of agency, a folk-psychological theory of how humans make decisions and take action (referred to as practical reasoning after Aristotle), and which would lead to a new computing paradigm, agent-oriented programming or agent technology more in general, when AI researchers started to apply it to the specification and implementation of artificial agents.

The main new ingredient in Bratman's theory is that of *intention*. Beliefs and desires were already known to be of importance in human behaviour. For instance, Daniel Dennett's intentional stance, the strategy of interpreting the behaviour of an entity by treating it as if it were a rational agent that governed its choice of action by a consideration of its beliefs and desires, already mentions the role of beliefs and desires in 1987. But Bratman claims that to fully understand the practical reasoning of humans also the notion of intention is needed. An intention is not just a mere desire but something the agent is committed to, that is, not given up too soon by the agent. For instance, if I have an intention to give a lecture in Amsterdam tomorrow, it is not a mere wish to do so, but I'm really taking measures (making plans, e.g., cancelling other plans or making sure my laptop will be in my bag) to do it and unless something happens that seriously interferes with my intention to give that lecture tomorrow, I really will do so. Thus Bratman takes intention to be a first-class citizen, and not something that can be reduced to beliefs and desires. In other words, a reduction of intention to a theory of only beliefs and desires is rejected by Bratman.

Another important notion in Bratman's theory is that of a *plan*. As he explains, rational agents need plans for two reasons. First of all, agents need to "allow deliberation and rational reflection to influence action beyond the present",

since agents have only limited resources (including time) to deliberate at the time of action. Secondly agents need to coordinate their actions, both on an intrapersonal and a interpersonal level, and plans help agents with that, too. As to the relation between plans and intentions, Bratman says that 'our intentions concerning our future actions are typically elements in larger plans'. Bratman focuses on *future-oriented intentions*. Such intentions differ from present-directed intentions, alias intentions-in-action, which accompany an agent's actions (more precisely, an agent's intentional actions), and pertain to what to do beginning now.

To explain the differences between beliefs, desires and intentions Bratman introduces the notion of a *pro-attitude*. A pro-attitude is an agent's mental attitude directed toward an action under a certain description. It plays a motivational role. So desires and intentions are both pro-attitudes while beliefs typically are not. But although desires and intentions are both pro-attitudes they differ. Intentions are conduct-controlling pro-attitudes, while ordinary desires are merely potential influences of action. The '*volitional*' dimension of the commitment involved in future-directed intentions comes from the conduct-controlling nature of intentions: as a conduct-controlling pro-attitude an intention involves a special commitment to action that ordinary desires do not.

Besides identifying intentions as conduct-controlling pro-attitudes, Bratman argues that intentions also have other properties: they have inertia and they serve as inputs into further practical reasoning. By the former is meant that intentions resist reconsideration. This has to do again with the resource-boundedness of realistic cognitive agents. The agent normally simply lacks the time to compute, at any given time, the optimal plan of action given his beliefs and desires, so it has to form future-directed intentions and store them in his mental agenda and use them to avoid computing plans all the time. Once an intention has been formed (and a commitment to action has been made) the intention will normally remain intact until the time of action: it has a characteristic stability / inertia. By this Bratman means that intentions made influence further reasoning about (decisions about) action, where also refinements of intentions (intentions to do more concrete actions) may play a role. For example, if I have the intention to speak in Amsterdam tomorrow, I can form a more refined intention to take the car driving to Amsterdam in order to speak. As a consequence after the second intention it won't be rational anymore to consider time tables for trains going to Amsterdam, while it was so after the first intention. All this has led to seeing intentions as distinctive states of mind, distinct from beliefs and desires, and to a belief-desire-intention model rather than a desire-belief model of practical reasoning.

As we have seen, Bratman describes how prior intentions and plans provide a filter of admissibility on options. This is what later by Cohen & Levesque has been called a 'screen of admissibility'. The basis of this role of intentions in further practical reasoning is the need for consistency in one's web of intentions and beliefs, as Bratman calls it: other things being equal, it should be possible for me to do all that I intend in a world in which my beliefs are true. But as Bratman explains this is not as simple as it looks. In particular "not every option that is incompatible with what the agent already intends and believes is inadmissible." In short, it depends on whether beliefs can be forced to be changed by the new intention so

that the inconsistency disappears or not. In the former case the new intention is admissible, in the latter it is not.

In Bratman's view there is an intrinsic relation between intentions and plans. Plans are intentions. They share the properties of intentions: they resist reconsideration and have inertia, they are conduct controllers and not merely conduct influencers, and they provide crucial inputs for further practical reasoning and planning. But they have increased complexity as compared to simple intentions: they are typically partial in the sense of incomplete (typically I have a partial plan to do something and fill in the details later) and have a hierarchical structure (plans concerning ends embed plans concerning means and preliminary steps, and more general intentions embed more specific ones).

To sum it up, according to Bratman, future-oriented intentions have the following characteristics:

- An intention is a *high-level plan*.
- An intention guides deliberation and triggers further planning: it typically leads to the *refinement* of a high-level plan into a more and more precise plan.
- An intention comes with the agent's *commitment* to achieve it.
- An agent abandons an intention only under the following conditions:
 - the intention has been achieved;
 - he believes it is impossible to achieve it;
 - he abandons another intention for which it is instrumental.

Let us illustrate Bratman's future-oriented intentions by an example. Suppose we are in autumn and I desire to go to Paris next spring. Under certain conditions—such as the importance of that desire and my beliefs about its feasibility—, that desire will make me form an intention to travel to Paris next spring. This is a very high-level plan: I do not settle the exact dates, I do not decide by which means of transportation I am going to go to Paris, and I do not know where to stay yet. I am however committed to that plan: during the following months I will stick to my intention to go to Paris, unless I learn that it is impossible to go to Paris in spring (say because my wife wants to spend our spring holidays in Spain, or because I changed my mind due to an invitation to give a talk at an important conference). During the next months I am going to refine my high-level plan: I will decide to go on a particular weekend, I will decide to go by train and not by plane, and I will book a hotel for the weekend under concern. This more elaborated plan is going to be refined further as time goes by: I decide to take the 7am train and not the 9am train, and I decide to go to the train station by metro and not by taxi, etc. Finally, once I have spent that weekend in Paris I no longer pursue that goal and drop it.

Bratman's theory might be called semi-formal: while he isolates the fundamental concepts and relates them, he does not provide a formal semantics. This was both undertaken by Phil Cohen and Hector Levesque and, more or less at the same time, by Anand Rao and Michael Georgeff. In the next section we will go

into the details of how they casted Bratman's theory into a logic of intention and we present some subsequent modifications and extensions of their original logic. This will be followed by a section on Rao and Georgeff's approach.

10.2.1 Cohen and Levesque's approach to intentions

We have seen that the concepts of belief, desire, time and action play an important role in Bratman's theory of intention. A logical analysis of that theory should involve combining a logic of belief, a logic of desire, a logic of time and a logic of action.

Belief, time, and action play a fundamental role in Cohen and Levesque's logic. However, the concept of desire is somewhat neglected: Cohen and Levesque rather base their logic on the concept of *realistic preference*. The latter can be viewed as a desire that has already been filtered by the agent's beliefs about its realisability. This is highlighted by the property that belief implies realistic preference: when I am convinced that φ is true then I also have to prefer that φ is true. (I might however prefer that φ be false at some point in the future.)

Cohen and Levesque's analysis amounts to a *reduction* of the concept of intention to those of belief, realistic preference, time and action: they define intention in terms of the latter four concepts. The reader may note that this is actually a surprising move, given that Bratman had strongly argued that intentions are independent of desires and cannot be reduced to them.

In the next two sections we present the four building blocks of Cohen and Levesque's logic, grouping together action and time on the one hand, and belief and realistic preference on the other.

Action and time

The basic building block of Cohen and Levesque's logic is a linear version of propositional dynamic logic PDL. The semantics of linear PDL allows to also interpret the temporal operators of linear-time temporal logic LTL.

Standard PDL Standard PDL is not about actions but about events. It has a set \mathcal{A} of atomic event names. Cohen and Levesque add agents to the picture and provide an *agentive version* of PDL. Let us write i, j , etc. for agents from some set of individuals \mathcal{I} . Then atomic actions are elements of $\mathcal{I} \times \mathcal{A}$. We write them $i:\alpha$ where $\alpha \in \mathcal{A}$ is an atomic event and $i \in \mathcal{I}$. Formulas of language of PDL are built from atomic formulas and atomic actions by means of modal operators Poss_π , where π is an action. The formula $\text{Poss}_\pi\varphi$ reads "there is a possible execution of π after which φ is true".¹ This reading highlights that the standard version of PDL allows for several possible executions of π in order to account for indeterminism.

While Poss_π quantifies existentially over the executions of π , the dual modal operator After_π quantifies universally. It is definable from Poss_π as follows:

$$\text{After}_\pi\varphi \stackrel{\text{def}}{=} \neg\text{Poss}_\pi\neg\varphi$$

¹ The standard notation is $\langle\pi\rangle\varphi$; we here deviate in order to be able to distinguish actual action from potential action.

Let us consider the case where φ is truth \top or falsity \perp : $\mathbf{Poss}_\pi\top$ has to be read “ π is executable”, while $\mathbf{After}_\pi\perp$ has to be read “ φ is inexecutable”.

The semantics of PDL is based on *transition systems* where an atomic action $i:\alpha$ can be interpreted as a set of edges. Such a transition system is a pair $\langle W, R \rangle$ where W is a non-empty set of possible worlds and R maps every action π to an accessibility relation $R_\pi \subseteq W \times W$ relating possible worlds. An edge from world w to world u that is labeled π means that it is possible to execute π in w and that u is a possible outcome world when π is executed. The set of all these π edges makes up the accessibility relation R_π interpreting the action π .

A PDL *model* is a transition system together with a valuation V mapping atomic formulas p from the set of propositional variables \mathcal{P} to their extension $V(p) \subseteq W$, i.e., to the set of worlds $V(p)$ where p is true.

Models allow to give truth values to formulas. In particular, $\mathbf{Poss}_\pi\varphi$ is true at a world w if there is a couple (w, w') in R_π such that φ is true at world w' :

$$M, w \models \mathbf{Poss}_\pi\varphi \quad \text{iff} \quad \text{there is a } u \in W \text{ such that } wR_\pi u \text{ and } M, u \models \varphi$$

The formula $\mathbf{Poss}_\pi\varphi$ therefore expresses a weak notion of ability: the action π might occur and φ could be true afterwards.

Linear PDL Probably Cohen and Levesque were the first to adapt PDL in order to model actual agency. The modalities are interpreted in *linear* PDL models. In such models, for every possible world w there is at most one successor world u that is temporally related to w . The accessibility relation linking w to u may be labelled by several atomic actions. Formally, a transition system $\langle W, R \rangle$ is linear if for every world $w \in W$ such that $\langle w, u_1 \rangle \in R_{\pi_1}$ and $\langle w, u_2 \rangle \in R_{\pi_2}$ we have $u_1 = u_2$. An edge from world w to world u that is labeled π means that π is executed in w and that u will be the result. (The reader might note the difference with the above standard PDL.) This allows for the simultaneous performance of two different actions; they must however lead to the same outcome world. The models of linear PDL are the class of linear transition systems.

In order to distinguish the modal operators of actual action from the modal operators of possible action we write the former as $\mathbf{Happ}_\pi\varphi$, read “ π is going to be performed, and φ is true afterwards”. Just as \mathbf{After}_π is the dual of \mathbf{Poss}_π , we define a modal operator \mathbf{IfHapp}_π that is the dual of \mathbf{Happ}_π by stipulating:

$$\mathbf{IfHapp}_\pi\varphi \stackrel{\text{def}}{=} \neg\mathbf{Happ}_\pi\neg\varphi$$

$\mathbf{Happ}_\pi\varphi$ and $\mathbf{IfHapp}_\pi\varphi$ say different things: the first formula says that π is executable and that φ is true after it, while the second says that *if* π is executable then φ is true after it. The former should therefore imply the latter.

The truth condition for \mathbf{Happ}_π is:

$$M, w \models \mathbf{Happ}_\pi\varphi \quad \text{iff} \quad \text{there is a } u \in W \text{ such that } wR_\pi u \text{ and } M, u \models \varphi$$

So it has exactly the same form as that for \mathbf{Poss}_π . We changed the name of the modal operator in order to better suit the linearity of the models.

The following axiom schema characterises linear PDL models:

$$(\mathbf{Happ}_{i:\alpha}\top \wedge \mathbf{Happ}_{j:\alpha'}\varphi) \rightarrow \mathbf{Happ}_{i:\alpha}\varphi \quad (10.1)$$

Beyond atomic events, PDL also has complex events such as sequential and nondeterministic composition, test, and iteration. We will however not refer to them in our present introduction.

Cohen and Levesque's logic has the temporal operators "eventually" (noted **F**), "henceforth" (noted **G**), and "until" (noted **U**). These operators are interpreted in linear PDL models in the obvious way. Let us give the truth condition for the 'eventually' operator, for example:

$$M, w \models \mathbf{F}\varphi \quad \text{iff} \quad \text{there is an integer } n \text{ and there are } v_1, \dots, v_n \in W \text{ such} \\ \text{that } v_1 = w, \langle v_k, v_{k+1} \rangle \in R_{\pi_k} \text{ for some } \pi_k, \text{ and } M, v_n \models \varphi$$

Cohen and Levesque also need existential quantification over actions. We here present their account in terms of an operator fusing existential quantification \exists over events α with the dynamic operator $\mathbf{Happ}_{i:\alpha}$. Its truth condition is as follows:

$$M, w \models \exists \alpha \mathbf{Happ}_{i:\alpha} \varphi \quad \text{iff} \quad \text{there are } \alpha \in \mathcal{A}, u \in W \text{ such that} \\ \langle w, u \rangle \in R_{i:\alpha} \text{ and } M, u \models \varphi$$

Belief and preference

Cohen and Levesque's account of belief is standard, while their account of preference is in terms of the somewhat unusual notion of strong realistic preference.

Belief Cohen and Levesque have modal operators of belief \mathbf{Bel}_i , one per agent i . The modal logic of each of these operators is the standard logic of belief KD45 (see also Chapter 1 of this handbook). Such operators can be interpreted if we add accessibility relations B_i to the transition systems of linear PDL, one per agent i . The set of worlds $B_i(w) = \{u : \langle w, u \rangle \in B_i\}$ is the set of those worlds that are possible for agent i at world w : the set of worlds that are compatible with his beliefs at w .

In order to be an accessibility relation for KD45, each of these relations has to satisfy the following constraints:

- for every $w \in W$ there is at least one $u \in W$ such that $\langle w, u \rangle \in B_i$ (seriality);
- if $\langle w, u \rangle \in B_i$ and $\langle u, v \rangle \in B_i$ then $\langle w, v \rangle \in B_i$ (transitivity);
- if $\langle w, u \rangle \in B_i$ and $\langle w, v \rangle \in B_i$ then $\langle u, v \rangle \in B_i$ (Euclideaness).

These constraints make that the following implications become valid:

- $\mathbf{Bel}_i \varphi \rightarrow \neg \mathbf{Bel}_i \neg \varphi$ (consistency of belief, axiom D)
- $\mathbf{Bel}_i \varphi \rightarrow \mathbf{Bel}_i \mathbf{Bel}_i \varphi$ (positive introspection, axiom 4)
- $\neg \mathbf{Bel}_i \varphi \rightarrow \mathbf{Bel}_i \neg \mathbf{Bel}_i \varphi$ (negative introspection, axiom 5)

Preference For Cohen and Levesque, intentions are particular *strong realistic preferences*. The latter are true in a subset of the worlds that are doxastically possible for an agent. There is a modal operator \mathbf{Pref}_i , one per agent i , and the formula $\mathbf{Pref}_i\varphi$ reads “ i chooses φ to be true”.² Such a notion of preference is strongly realistic in the sense that belief logically implies preference. Semantically, strong realistic preference can be modelled by accessibility relations P_i , one per agent $i \in \mathcal{I}$, such that $P_i \subseteq B_i$. The latter constraint implements realism: a world that is compatible with agent i ’s preferences cannot be incompatible with i ’ beliefs. In other words, at world w agents have to select their preferred worlds among the worlds that are epistemically possible for them at w .

The logic of action, time, belief, and preference

Let us sum up Cohen and Levesque’s semantics. A *frame* is a quadruple $M = \langle W, R, B, P \rangle$ where

- W is a non-empty set of possible worlds;
- $R : (\mathcal{I} \times \mathcal{A}) \longrightarrow (W \times W)$ maps actions π to accessibility relations R_π ;
- $B : \mathcal{I} \longrightarrow (W \times W)$ maps agents i to accessibility relations B_i ;
- $P : \mathcal{I} \longrightarrow (W \times W)$ maps agents i to accessibility relations P_i ;

These frames have to satisfy the following constraints:

- $\langle W, R \rangle$ is a linear transition system;
- every B_i is serial, transitive and Euclidean;
- $P_i \subseteq B_i$, for every $i \in \mathcal{I}$.

Let us call \mathcal{CL} that class of frames. As usual, a model is a frame together with a valuation $V : \mathcal{P} \longrightarrow 2^W$ mapping atomic formulas p to their extension $V(p) \subseteq W$. Validity and satisfiability in \mathcal{CL} frames are defined as usual.

We are now ready to formulate Cohen and Levesque’s reduction of intention.

Defining intention

Cohen and Levesque define a modal operator of intention by means of a cascade of definitions. We here reproduce them in a slightly simplified form. We then discuss them and finally comment on the modifications.

1. i has the *achievement goal* that φ if i prefers that φ is eventually true and believes that φ is currently false. Formally:

$$\mathbf{AGoal}_i\varphi \stackrel{\text{def}}{=} \mathbf{Pref}_i\mathbf{F}\varphi \wedge \mathbf{Bel}_i\neg\varphi$$

² Cohen and Levesque’s original notation is \mathbf{Goal}_i instead of \mathbf{Pref}_i (while they actually refer to it in their title as ‘choice’). We moved to our notation in order to avoid confusion with the concept of choice in STIT theory.

2. i has the *persistent goal* that φ if i has the achievement goal that φ and will keep that goal until it is either fulfilled or believed to be out of reach. Formally:

$$\text{PGoal}_i\varphi \stackrel{\text{def}}{=} \text{AGoal}_i\varphi \wedge (\text{AGoal}_i\varphi) \text{U} (\text{Bel}_i\varphi \vee \text{Bel}_i\text{G}\neg\varphi)$$

3. i has the *intention* that φ if i has the persistent goal that φ and believes he can achieve φ by an action of his. This requires to quantify over i 's actions by means of the fused operator quantifying over events. Formally:

$$\text{Intend}_i\varphi \stackrel{\text{def}}{=} \text{PGoal}_i\varphi \wedge \text{Bel}_i\text{F}\exists\alpha\text{Happ}_{i:\alpha}\varphi$$

Some valid and invalid principles for intention Cohen and Levesque's construction guarantees several desirable properties and avoids some that are unwanted. Here are two of them.

First, i 's intention that φ logically implies i 's belief that $\neg\varphi$. Formally this writes:

$$\text{Intend}_i\varphi \rightarrow \text{Bel}_i\neg\varphi$$

Second, the formula schema $\text{Bel}_i(\varphi \rightarrow \psi) \rightarrow (\text{Intend}_i\varphi \rightarrow \text{Intend}_i\psi)$ is invalid: i 's intention that φ together with i 's belief that φ implies ψ does not logically imply i 's intention that ψ . This is crucial both for Bratman and for Cohen and Levesque. Here is a famous example illustrating why the principle should not be valid: if I intend to go to the dentist and believe that going to the dentist will cause pain then I do not necessarily intend to have pain. This is called the Side-Effect-Free Principle by Su *et al.* They have proposed a logic with a semantics in terms of linear neighbourhood structures (instead of accessibility relations) in order to interpret a modal operator of preference. Such structures allow to validate the principle of consistency of beliefs and intentions ($\text{Bel}_i\varphi \rightarrow \neg\text{Pref}_i\neg\varphi$) while guaranteeing the Side-Effect-Free Principle. The price for that is that preference no longer satisfy the monotony axiom $\text{Pref}_i(\varphi \wedge \psi) \rightarrow (\text{Pref}_i\varphi \wedge \text{Pref}_i\psi)$.

Comments on the simplifications We have simplified the definition of a persistent goal. Cohen and Levesque's original definition allows agents to abandon a persistent goal for some other, superior reason. Their definition is

$$\text{PGoal}_i\varphi \stackrel{\text{def}}{=} \text{AGoal}_i\varphi \wedge (\text{AGoal}_i\varphi) \text{U} (\text{Bel}_i\varphi \vee \text{Bel}_i\text{G}\neg\varphi \vee \psi)$$

where ψ is an unspecified condition accounting for that other reason.

This definition stipulates that a persistent goal φ is also abandoned if some other condition ψ becomes true. This leaves room for the abandonment of persistent goals that are instrumental for some goal that obtains without the agent's intervention. A classical example is a student i coming back late in the night to the dorms who forgot to take the key of the entrance door: his overall goal is to be able to get into the building and he starts to plan to climb over the wall (φ), but then some other student who also comes home late happens to pass just in front of him and opens the door (ψ), thus enabling i to drop φ . Hence, i abandons

his persistent goal φ although it has neither been achieved nor turned out to be impossible.

Dealing with such a general condition ψ however makes it difficult to go beyond specific cases.

Variants and extensions

We now overview an extension of the basic logic, together with several alternatives to Cohen and Levesque's incremental definition of intention.

Introspection of intention Cohen and Levesque do not assume principles of positive and negative introspection of preference. However, they seem natural in the following form

- $\text{Pref}_i\varphi \rightarrow \text{Bel}_i\text{Pref}_i\varphi$
- $\neg\text{Pref}_i\varphi \rightarrow \text{Bel}_i\neg\text{Pref}_i\varphi$

They correspond to the following constraints on the accessibility relations for preference and belief:

- if $\langle w, u \rangle \in B_i$ and $\langle u, v \rangle \in P_i$ then $\langle w, v \rangle \in P_i$;
- if $\langle w, u \rangle \in B_i$ and $\langle w, v \rangle \in P_i$ then $\langle u, v \rangle \in P_i$.

This allows to prove principles of positive and negative introspection of goals, achievement goals, persistent goals, and intention. For instance, they validate $\text{Intend}_i\varphi \rightarrow \text{Bel}_i\text{Intend}_i\varphi$ and $\neg\text{Intend}_i\varphi \rightarrow \text{Bel}_i\neg\text{Intend}_i\varphi$. For example, both when I intend to go to Paris and when I don't intend to go to Paris then I am aware of this.

Weakly realistic preference Sadek has argues for a slightly different notion of realistic preference. The latter does not demand that all preference-accessible worlds be in the set of belief-accessible worlds, but only requires that they have a non-empty intersection. (This is sometimes called weak realism, as opposed to Cohen and Levesque's strong realism.) In frames that are weakened in this way the somewhat counterintuitive principle $\text{Bel}_i\varphi \rightarrow \text{Pref}_i\varphi$ is no longer valid. Instead, the weaker $\text{Bel}_i\varphi \rightarrow \neg\text{Pref}_i\neg\varphi$ is valid, which can be reformulated as

$$\neg(\text{Bel}_i\varphi \wedge \text{Pref}_i\neg\varphi)$$

It says that one cannot simultaneously believe that φ and prefer that $\neg\varphi$.

An epistemic version of achievement goals Herzig and Longin have advocated a different definition of an achievement goal. It is weaker than Cohen and Levesque's in that they only require that i does not believe that φ is currently true (instead of i 's belief that φ is currently false). It is stronger in that they replace i 's goal that φ will be true by i 's goal that φ will be *believed*. This gives the following definition:

$$\text{AGoal}_i^w\varphi \stackrel{\text{def}}{=} \text{Pref}_i\text{FBel}_i\varphi \wedge \neg\text{Bel}_i\varphi$$

They start by arguing for the replacement of $\text{Pref}_i F\varphi$ by $\text{Pref}_i F\text{Bel}_i\varphi$. As they point out, it is the *raison d'être* of an intention to be abandoned at some stage, and an agent can only do so if he believes that he has achieved that goal. So the agent's goal cannot just be that φ be true, but should be that he *believes* that φ is true. Using the same reasons they then argue for the replacement of $\text{Bel}_i\neg\varphi$ by $\neg\text{Bel}_i\varphi$: as long as φ is not believed to be true the agent should stick to his achievement goal φ , so $\neg\text{Bel}_i\varphi$ is better in line with this than $\text{Bel}_i\neg\varphi$. They illustrate the first replacement by a variant of the Byzantine generals example. Let r mean that a message of general i has been received by general j . Suppose i initially believes that j has not received the message yet, i.e., $\text{Bel}_i\neg r$. Suppose moreover that i believes that he will actually *never* know whether j received the message or not, i.e., $\text{Bel}_i G(\neg\text{Bel}_i r \wedge \neg\text{Bel}_i\neg r)$. (This differs from the original Byzantine generals example, where it is possible that the messengers get through and where it is just possible for i that he will never know.) If we express i 's achievement goal that r as $\text{Pref}_i F r$ then Cohen and Levesque make us conclude that $\text{AGoal}_i r$, i.e., i has the achievement goal that φ although he believes that he will never be able to abandon that goal. In contrast, if we express i 's achievement goal that r as $\text{Pref}_i F\text{Bel}_i r$ then we have $\neg\text{AGoal}_i^w r$: i cannot have the achievement goal that r .³

Weaker link between action and goal Sadek and Bretier point out that the definition of intention is too strong in particular in cooperative situations where agent i 's action need not directly achieve his goal φ : it is enough that i triggers a subsequent action of another agent j which will achieve i 's goal. Their modification can be formulated as follows:

$$\text{Intend}_i\varphi \stackrel{\text{def}}{=} \text{PGoal}_i\varphi \wedge \text{Pref}_i F(\exists\alpha \text{Happ}_{i,\alpha} F\varphi)$$

Stronger commitment Sadek and Bretier also discuss a stronger definition of intention where the agent is committed to do all he can to achieve his goal. They express this by a universal quantification over events.⁴ We formulate their definition as follows:

$$\text{Intend}_i\varphi \stackrel{\text{def}}{=} \text{PGoal}_i\varphi \wedge \text{Pref}_i \forall\alpha(\text{Bel}_i \text{Happ}_{i,\alpha} F\varphi \rightarrow \text{Pref}_i F\text{Happ}_{i,\alpha} \top)$$

That definition was criticised in the literature, in particular by Herzig and Long, as being too strong. Indeed, it postulates that agents want to achieve their intentions by all possible means, including illegal actions and actions with a huge cost for them. For example, it might commit me to steal a car if this is the only means to go to Paris on that spring weekend (say because there is a train strike).

Attempts Lorini and Herzig complement Cohen and Levesque's approach by integrating the concept of an *attempt* to perform an action. The motivation is that

³ Our hypothesis $\text{Bel}_i\neg\varphi$ implies the second condition $\neg\text{Bel}_i\varphi$ because the logic of belief contains the D axiom, and $\text{Bel}_i G(\neg\text{Bel}_i r \wedge \neg\text{Bel}_i\neg r)$ implies $\neg\text{Pref}_i F\text{Bel}_i r$, which is the negation of the first condition.

⁴ This is therefore not a fused operator. In order to save space we do not give the the details of the semantics of that quantifier and rely on the reader's intuitions about it.

intentions typically make an agent *try* to perform an action, while the successful performance of that action is not guaranteed. The central principle there is “can and attempts implies does” : if i intends to (attempt to) perform α and α is feasible then α will indeed take place. This requires a logic with both modal operators of possible action Poss_π and modal operators of actual action Happ_π .

Conclusion

Cohen and Levesque succeeded in providing a fine-grained analysis of intention by relating that concept to action, belief and realistic preference. A central point in Bratman’s theory their logic does not account for is the refinement of intentions. According to Bratman, an agent starts by forming high-level intentions such as going to Paris in a month, and as time goes by he makes that intention more precise: he first starts to intend to go to Paris by train and not by plane; at a later stage he decomposes the intention to go to Paris by train into the intention to take a taxi to the train station (instead of a bus), then take the TGV to Paris, and then take the metro. It is probably an interesting direction of future research to integrate intention refinement mechanisms e.g. by resorting to dynamic epistemic logics (see Chapter 6).

10.2.2 Rao & Georgeff’s BDI logic

As mentioned earlier, besides Cohen and Levesque, also Rao and Georgeff, more or less at the same time, published a formalisation of the ground-breaking work of Bratman on the philosophy of intelligent (human) agents. As we have seen, Bratman made a case for the notion of *intention* besides belief and desire, to describe the behaviour of rational agents. Intentions force the agent to commit to certain desires and to really ‘go for them’. So focus of attention is an important aspect here, which also enables the agent to monitor how s/he is doing and take measures if things go wrong. Rao & Georgeff stress that in the case of resource-bounded agents it is imperative to focus on desires / goals and make choices. This was also observed by Cohen & Levesque, who try to formalize the notion of intention in a linear-time temporal logic (or, as we have seen in the previous section, a linear version of dynamic logic) in terms of the notion of a (persistent) goal.

Here we treat Rao & Georgeff’s approach who base it on branching-time temporal logic framework CTL^* to give a formal-logical account of BDI theory. The reader may also like to look at Chapter 5 of this book, the chapter that relates knowledge and time. Like Cohen & Levesque’s approach, BDI logic has influenced many researchers (including Rao & Georgeff themselves) to think about architectures of agent-based systems in order to realize these systems. Rao & Georgeff’s BDI logic is more liberal than that of Cohen & Levesque in the sense that they *a priori* regard each of the three attitudes of belief, desire and intention as primitive: they introduce separate modal operators for belief, desire and intention, and then study possible relations between them.

(The language of) BDI logic is constructed as follows. Two types of formulas are distinguished: state formulas and path formulas. We assume some given first-

order signature. Furthermore, we assume a set E of event types with typical element e . The operators BEL , $GOAL$, $INTEND$ have as obvious intended reading the belief, goal and intention of an agent, respectively, while U, \diamond, O are the usual temporal operators, viz. until, eventually and next, respectively.

Definition 10.1 (State and path formulas)

1. The set of *state formulas* is the smallest closed under:
 - any first-order formula w.r.t. the given signature is a state formula
 - if φ_1 and φ_2 are state formulas then also $\neg\varphi_1, \varphi_1 \vee \varphi_2, \exists x\varphi_1(x)$ are state formulas
 - if e is an event type, then $succeeded(e), failed(e)$ are state formulas
 - if φ is a state formula, then $BEL(\varphi), GOAL(\varphi), INTEND(\varphi)$ are state formulas
 - if ψ is a *path formula*, then $optional(\psi)$ is a state formula
2. The set of *path formulas* is the smallest set closed under:
 - any state formula is a path formula
 - if ψ_1, ψ_2 are path formulas, then $\neg\psi_1, \psi_1 \vee \psi_2, \psi_1 U\psi_2, \diamond\psi_1, O\psi_1$ are path formulas

State formulas are interpreted over a state, that is a (state of the) world at a particular point in time, while path formulas are interpreted over a path of a time tree (representing the evolution of a world). In the sequel we will see how this will be done formally. Here we just give the informal readings of the operators.

The operators *succeeded* and *failed* are used to express that events have (just) succeeded and failed, respectively. Next there are the modal operators for belief, goal and intend. (In the original version of BDI theory, desires are represented by goals, or rather a GOAL operator. In a later paper the GOAL operator was replaced by DES for desire.) The optional operator states that there is a future (represented by a path) where the argument of the operator holds. Finally, there are the familiar (linear-time) temporal operators, such as the ‘until’, ‘eventually’ and ‘nexttime’, which are to be interpreted along a linear time path.

Definition 10.2

The following abbreviations are defined:

1. $\Box\psi = \neg\diamond\neg\psi$ (always)
2. $inevitable(\psi) = \neg optional(\neg\psi)$
3. $done(e) = succeeded(e) \vee failed(e)$
4. $succeeds(e) = inevitableO(succeeded(e))$
5. $fails(e) = inevitableO(failed(e))$
6. $does(e) = inevitableO(done(e))$

The ‘always’ operator is the familiar one from (linear-time) temporal logic. The ‘inevitability’ operator expresses that its argument holds along all possible futures (paths from the current time). The ‘done’ operator states that an event occurs (action is done) no matter whether it is succeeding or not. The final three operators state that an event succeeds, fails, or is done iff it is inevitable (i.e. in any possible future) it is the case that at the next instance the event has succeeded, failed, or has been done, respectively (note that this means that an event, succeeding or failing, is supposed to take one unit of time).

Definition 10.3 (Semantics)

The semantics is given w.r.t. models of the form $\mathcal{M} = \langle W, E, T, \prec, \mathcal{U}, B, G, I, \Phi \rangle$, where

- W is a set of possible worlds
- E is a set of primitive event types
- T is a set of time points
- \prec is a binary relation on time points, which is serial, transitive and backwards linear
- \mathcal{U} is the universe of discourse
- Φ is a mapping of first-order entities to \mathcal{U} , for any world and time point
- $B, G, I \subseteq W \times T \times W$ are accessibility relations for *BEL*, *GOAL*, *INTEND*, respectively –

The semantics of BDI logic, Rao & Georgeff-style, is rather complicated. Of course, we have possible worlds again, but as we will see below, these are not just unstructured elements, but they are each time trees, describing possible flows of time. So, we also need time points and an ordering on them. As BDI logic is based on branching time, the ordering need not be linear in the sense that all time points are related in this ordering. However, it is stipulated that the time ordering is serial (every time point has a successor in the time ordering), the ordering is transitive and backwards-linear, which means that every time point has only one direct predecessor. The accessibility relations for the ‘BDI’-modalities are standard apart from the fact that they are also time-related, that is to say that worlds are (belief/goal/intend-)accessible with respect to a time point. Another way of viewing this is that – for all three modalities – for every time point there is a distinct accessibility relation between worlds.

In order to obtain reasonable properties for beliefs, desires and intentions, a number of constraints on the accessibility relations are stipulated. First of all, a *world / time point compatibility* requirement is assumed for all of the B, G, I accessibility relations: for $R = B, G, I$:

$$\text{If } w' \in R(w, t) \text{ then } t \in w \text{ and } t \in w'$$

where $R(w, t) = \{w' \mid R(w, t, w')\}$ for $R = B, G, I$. This requirement is needed for the semantic clauses for the BEL, GOAL and INTEND modalities that we will give below to work. And next there are the usual requirements of the B

accessibility relation to satisfy seriality, transitivity and Euclideaness in order to obtain the familiar KD45 properties of belief: beliefs are consistent, and satisfy positive and negative introspection. As to the G and I accessibility relations we require seriality in order to obtain the well-known KD property of consistent goals and intentions.

Next we elaborate on the structure of the possible worlds.

Definition 10.4 (Possible worlds)

Possible worlds in W are assumed to be *time trees*: an element $w \in W$ has the form $w = \langle T_w, A_w, S_w, F_w \rangle$ where

- $T_w \subseteq T$ is the set of time points in world w
- A_w is the restriction of the relation \prec to T_w
- $S_w : T_w \times T_w \rightarrow E$ maps adjacent time points to (successful) events
- $F_w : T_w \times T_w \rightarrow E$ maps adjacent time points to (failing) events
- the domains of the functions S_w and F_w are disjoint ⊢

As announced before, a possible world itself is a time tree, a temporal structure representing possible flows of time. The definition above is just a technical one stating that the time relation within a possible world derives naturally from the *a priori* given relation on time points. Furthermore it is indicated by means of the functions S_w and F_w how events are associated with adjacent time points.

Now we come to the formal interpretation of formulas on the above models. Naturally we distinguish state formulas and path formulas, since the former should be interpreted on states whereas the latter are interpreted on paths. In the sequel we use the notion of a *fullpath*: a fullpath in a world w is an *infinite* sequence of time points such that, for all i , $(t_i, t_{i+1}) \in A_w$. We denote a fullpath in w by $(w_{t_0}, w_{t_1}, \dots)$, and define *fullpaths*(w) as the set of all fullpaths occurring in world w (i.e., all fullpaths that start somewhere in the time tree w).

Definition 10.5 (Interpretation of formulas)

The interpretation of formulas w.r.t. a model $\mathcal{M} = \langle W, E, T, \prec, \mathcal{U}, B, G, I, \Phi \rangle$ is now given by:

1. (state formulas)
 - $\mathcal{M}, v, w_t \models q(y_1, \dots, y_n)$ iff $(v(y_1), \dots, v(y_n)) \in \Phi(q, w, t)$
 - $\mathcal{M}, v, w_t \models \neg\varphi$ iff $\mathcal{M}, v, w_t \not\models \varphi$
 - $\mathcal{M}, v, w_t \models \varphi_1 \vee \varphi_2$ iff $\mathcal{M}, v, w_t \models \varphi_1$ or $\mathcal{M}, v, w_t \models \varphi_2$
 - $\mathcal{M}, v, w_t \models \exists x\varphi$ iff $\mathcal{M}, v\{d/x\}, w_t \models \varphi$ for some $d \in \mathcal{U}$
 - $\mathcal{M}, v, w_{t_0} \models \text{optional}(\psi)$ iff exists fullpath $(w_{t_0}, w_{t_1}, \dots)$ such that $\mathcal{M}, v, (w_{t_0}, w_{t_1}, \dots) \models \psi$
 - $\mathcal{M}, v, w_t \models \text{BEL}(\varphi)$ iff for all $w' \in B(w, t) : \mathcal{M}, v, w'_t \models \varphi$
 - $\mathcal{M}, v, w_t \models \text{GOAL}(\varphi)$ iff for all $w' \in G(w, t) : \mathcal{M}, v, w'_t \models \varphi$
 - $\mathcal{M}, v, w_t \models \text{INTEND}(\varphi)$ iff for all $w' \in I(w, t) : \mathcal{M}, v, w'_t \models \varphi$

- $\mathcal{M}, v, w_t \models \text{succeeded}(e)$ iff exists t_0 such that $S_w(t_0, t) = e$
- $\mathcal{M}, v, w_t \models \text{failed}(e)$ iff exists t_0 such that $F_w(t_0, t) = e$ ¬

where $v\{d/x\}$ denotes the function v modified such that $v\{d/x\}(x) = d$. (Note that clauses for BEL, GOAL and INTEND are well-defined due to the world / time point compatibility requirement that we have assumed to hold.)

2. (path formulas)

- $\mathcal{M}, v, (w_{t_0}, w_{t_1}, \dots) \models \varphi$ iff $\mathcal{M}, v, w_{t_0} \models \varphi$, for φ state formula
- $\mathcal{M}, v, (w_{t_0}, w_{t_1}, \dots) \models O\varphi$ iff $\mathcal{M}, v, (w_{t_1}, w_{t_2}, \dots) \models \varphi$
- $\mathcal{M}, v, (w_{t_0}, w_{t_1}, \dots) \models \diamond\varphi$ iff $\mathcal{M}, v, (w_{t_k}, \dots) \models \varphi$ for some $k \geq 0$
- $\mathcal{M}, v, (w_{t_0}, w_{t_1}, \dots) \models \psi_1 U \psi_2$ iff
either there exists $k \geq 0$ such that $\mathcal{M}, v, (w_{t_k}, \dots) \models \psi_2$ and for all
 $0 \leq j < k : \mathcal{M}, v, (w_{t_j}, \dots) \models \psi_1$, or
for all $j \geq 0 : \mathcal{M}, v, (w_{t_j}, \dots) \models \psi_1$

Most of the above clauses should be clear, including those concerning the modal operators for belief, goal and intention. The clause for the ‘optional’ operator expresses exactly that optionally ψ is true if ψ is true in one of the possible futures represented by fullpaths starting at the present time point. The interpretation of the temporal operators is as usual.

Rao & Georgeff now discuss a number of properties that may be desirable to have as axioms. In the following we use α to denote so-called *O-formulas*, which are formulas that contain no positive occurrences of the ‘inevitable’ operator (or negative occurrences of ‘optional’) outside the scope of the modal operators *BEL*, *GOAL* and *INTEND*.

1. $GOAL(\alpha) \rightarrow BEL(\alpha)$
2. $INTEND(\alpha) \rightarrow GOAL(\alpha)$
3. $INTEND(\text{does}(e)) \rightarrow \text{does}(e)$
4. $INTEND(\varphi) \rightarrow BEL(INTEND(\varphi))$
5. $GOAL(\varphi) \rightarrow BEL(GOAL(\varphi))$
6. $INTEND(\varphi) \rightarrow GOAL(INTEND(\varphi))$
7. $\text{done}(e) \rightarrow BEL(\text{done}(e))$
8. $INTEND(\varphi) \rightarrow \text{inevitable} \diamond (\neg INTEND(\varphi))$

In order to render these formulas validities further constraints should be put on the models, since in the general setting above these are not yet valid.

For reasons of space we only consider the first two. In order to define constraints on the models such that these two become valid, we introduce the relation \triangleleft on worlds, as follows:

$w'' \triangleleft w' \Leftrightarrow \text{fullpaths}(w'') \subseteq \text{fullpaths}(w')$. So $w'' \triangleleft w'$ means that there the world (time tree) w'' represents less choices than w' .

Now we define the *B-G condition* as the property that the following holds:

$$\forall w' \in B(w, t) \exists w'' \in G(w, t) : w'' \triangleleft w'$$

Informally, this condition says that for any belief accessible world there is a goal accessible world that contains less choices. It is now easy to show the following proposition.

Proposition 10.1

Let \mathcal{BG} be the class of models of the above form that satisfy the B-G condition. Then: $\mathcal{BG} \models GOAL(\alpha) \rightarrow BEL(\alpha)$ for O-formulas α . \dashv

Similarly one can define the *G-I condition* as

$$\forall w' \in G(w, t) \exists w'' \in I(w, t) : w'' \triangleleft w'$$

and obtain:

Proposition 10.2

Let \mathcal{GI} be the class of models of the above form that satisfy the G-I condition. Then: $\mathcal{GI} \models INTEND(\alpha) \rightarrow GOAL(\alpha)$ for O-formulas α . \dashv

Let us now consider the properties deemed desirable by Rao & Georgeff again. The first formula describes Rao & Georgeff's notion of 'strong realism' and constitutes a kind of belief-goal compatibility: it says that the agent believes he can optionally achieve his goals. There is some controversy on this. Interestingly, but confusingly, Cohen & Levesque adhere to a form of realism that renders more or less the converse formula $BELp \rightarrow GOALp$. But we should be careful and realize that Cohen & Levesque have a different logic in which one cannot express options as in the branching-time framework of Rao & Georgeff. Furthermore, it seems that in the two frameworks there is a different understanding of goals (and beliefs) due to the very difference in ontologies of time employed: Cohen & Levesque's notion of time could be called 'epistemically nondeterministic' or 'epistemically branching', while 'real' time is linear: the agents envisage several future courses of time, each of them being a linear history, while in Rao & Georgeff's approach also 'real' time is branching, representing options that are available to the agent.

The second formula is a similar one to the first. This one is called goal-intention compatibility, and is defended by Rao & Georgeff by stating that if an optionality is intended it should also be wished (a goal in their terms). So, Rao & Georgeff have a kind of selection filter in mind: intentions (or rather intended options) are filtered / selected goals (or rather goal (wished) options), and goal options are selected believed options. If one views it this way, it looks rather close to Cohen & Levesque's 'Intention is choice (chosen / selected wishes) with commitment', or loosely, wishes that are committed to. Here the commitment acts as a filter.

The third one says that the agent really does the primitive actions that s/he intends to do. This means that if one adopts this as an axiom the agent is not allowed to do something else (first). (In our opinion this is rather strict on the agent, since it may well be that postponing its intention for a while is also an option.) On the other hand, as Rao & Georgeff say, the agent may also do things

that are not intended since the converse does not hold. And also nothing is said about the intention to do complex actions.

The fourth, fifth and seventh express that the agent is conscious of its intentions, goals and what primitive action he has done in the sense that he believes what he intends, has as a goal and what primitive action he has just done.

The sixth one says something like that intentions are really wished for: if something is an intention then it is a goal that it is an intention.

The eighth formula states that intentions will inevitably (in every possible future) be dropped eventually, so there is no infinite deferral of its intentions. This leaves open, whether the intention will be fulfilled eventually, or will be given up for other reasons. Below we will discuss several possibilities of giving up intentions according to different types of commitment an agent may have.

BDI-logical expressions can be used to characterize different types of agents. Rao & Georgeff mention the following possibilities:

1. (blindly committed agent) $INTEND(inevitable \diamond \varphi) \rightarrow inevitable(INTEND(inevitable \diamond \varphi) \cup BEL(\varphi))$
2. (single-minded committed agent) $INTEND(inevitable \diamond \varphi) \rightarrow inevitable(INTEND(inevitable \diamond \varphi) \cup (BEL(\varphi) \vee \neg BEL(optional \diamond \varphi)))$
3. (open minded committed agent) $INTEND(inevitable \diamond \varphi) \rightarrow inevitable(INTEND(inevitable \diamond \varphi) \cup (BEL(\varphi) \vee \neg GOAL(optional \diamond \varphi)))$

A blindly committed agent maintains his intentions to inevitably obtaining eventually something until he actually believes that that something has been fulfilled. A single-minded committed agent is somewhat more flexible: he maintains his intention until he believes he has achieved it *or he does not believe that it can be reached (i.e. that it is still an option in some future) anymore*. Finally, the open minded committed agent is even more flexible: he can also drop his intention if it is not a goal (desire) anymore.

Rao & Georgeff are then able to obtain results under which conditions the various types of committed agents will reach their intentions. For example, for a blindly committed agent it holds that under the assumption of the axioms we have discussed earlier plus an axiom that expresses no infinite deferral of intentions:

$$INTEND(\varphi) \rightarrow inevitable \diamond \neg INTEND(\varphi)$$

that

$$INTEND(inevitable(\diamond \varphi)) \rightarrow inevitable(\diamond BEL(\varphi))$$

expressing that if the agent intends to eventually obtain φ it will inevitably eventually believe that it has succeeded in achieving φ .

The branching-time setup of the approach as opposed to a linear-time one is much more expressive and is shown to solve problems such as the *Little Nell problem*. This is about a girl, Little Nell, that is in mortal peril, and a rescue agent that reasons like this: I intend to rescue Little Nell, and therefore I believe (because I'm confident that my actions will succeed) that she will be safe, but then I can drop my intention to rescue her just because she will be safe...! In a linear-time approach – if one is not very careful – this scenario results in a contradictory

(or unintuitive) representation (basically because there is only one future in which apparently Little Nell will be safe), while in a branching-time approach such as Rao and Georgeff's this presents no problem at all. In fact in CTL_{BDI} the scenario comes down to something like (here φ stands for "Little Nell is safe")

$$\begin{aligned} & \text{INTEND}(\text{inevitable} \diamond \varphi) \rightarrow \\ & \text{inevitable}(\text{INTEND}(\text{inevitable} \diamond \varphi) \text{UBEL}(\text{optional} \diamond \varphi)) \end{aligned}$$

informally saying that since the agent believes that there is a way (by performing its plan) to eventually reaching the goal φ , it may drop its intention to perform the plan to achieve eventually φ , which is definitely not valid in CTL_{BDI} ! Intuitively, this is the case, because there may be other branches along which Little Nell will not be safe, so that there is no reason to give up the intention to rescue her.

In the next section we will look at yet another approach, based on (non-linear) dynamic logic, which may perhaps be viewed as an amalgam of those of Cohen & Levesque (using dynamic logic) but allowing for non-linear, i.e. branching, structures.

10.3 KARO Logic

In this section we review the KARO formalism, in which *action*, together with knowledge / belief, is the primary concept, on which other agent notions are built. Historically, the KARO approach was the first approach truly based on dynamic logic, although as we have seen, in retrospect, we may view Cohen & Levesque's approach as being based on a linear variant of PDL (Propositional Dynamic Logic). There are differences, though. We will see that in KARO the fact that it is based on a logic of action is even more employed than in Cohen & Levesque: besides BDI-like notions such as knowledge, belief, desires, and goals that are operators that take formulas as arguments, and are quite similar in nature as the notions that are in Cohen & Levesque's approach, in KARO there are also operators taking actions as arguments such as ability and commitment, and operators that take both actions and formulas as arguments, such as a Can operator and a (possible) intention operator. All these operators are used to describe the mental state of the agent. But even more importantly, in the KARO framework (dedicated) actions are used to *change* the mental state of the agent. So there are revise, commit and uncommit actions to revise beliefs and update the agenda (the commitments) of the agent. In this sense KARO is related to dynamic epistemic logic, the topic of Chapter 6 in this handbook.

KARO logic for rational agents

The KARO formalism is an amalgam of dynamic logic and epistemic / doxastic logic, augmented with several additional (modal) operators in order to deal with the motivational aspects of agents. So, besides operators for knowledge (**K**), belief (**B**) and action ($[\alpha]$, "after performance of α it holds that"), there are additional operators for ability (**A**) and desires (**D**).

Assume a set \mathcal{A} of atomic actions and a set \mathcal{P} of atomic propositions.

Definition 10.6 (Language)

The language $\mathcal{L}_{\text{KARO}}$ of KARO-formulas is given by the BNF grammar:

$$\begin{aligned} \varphi & ::= p(\in \mathcal{P}) \mid \neg\varphi \mid \varphi_1 \wedge \varphi_2 \mid \dots \\ & \quad \mathbf{K}\varphi \mid \mathbf{B}\varphi \mid \mathbf{D}\varphi \mid [\alpha]\varphi \mid \mathbf{A}\alpha \\ \alpha & ::= a(\in \mathcal{A}) \mid \varphi? \mid \alpha_1; \alpha_2 \mid \alpha_1 + \alpha_2 \mid \alpha^* \quad \dashv \end{aligned}$$

Here the formulas generated by the second (α) part are referred to as actions (or rather action expressions). We use the abbreviations $\mathbf{tt} \equiv p \vee \neg p$ (for some fixed $p \in \mathcal{P}$) and $\mathbf{ff} \equiv \neg\mathbf{tt}$. Conditional and while-action are introduced by the usual abbreviations: **if** φ **then** α_1 **else** α_2 **fi** $\equiv (\varphi?; \alpha_1) + (\neg\varphi?; \alpha_2)$ and **while** φ **do** α **od** $\equiv (\varphi?; \alpha)^*; \neg\varphi?$.

Thus formulas are built by means of the familiar propositional connectives and the modal operators for knowledge, belief, desire, action and ability. Actions are the familiar ones from imperative programming: atomic ones, tests, sequential composition, (nondeterministic) choice and repetition.

Definition 10.7 (KARO models)

1. The semantics of the knowledge, belief and desires operators is given by means of Kripke structures of the following form: $\mathcal{M} = \langle W, \vartheta, R_K, R_B, R_D \rangle$, where
 - W is a non-empty set of states (or worlds)
 - ϑ is a truth assignment function per state
 - R_K, R_B, R_D are accessibility relations for interpreting the modal operators $\mathbf{K}, \mathbf{B}, \mathbf{D}$. The relation R_K is assumed to be an equivalence relation, while the relation R_B is assumed to be euclidean, transitive and serial. Furthermore we assume that $R_B \subseteq R_K$. No special constraints are assumed for the relations R_D .
2. The semantics of actions is given by means of structures of type $\langle \Sigma, \{R_a \mid a \in \mathcal{A}\}, \mathcal{C}, Ag \rangle$, where
 - Σ is the set of possible model/state pairs (i.e. models of the above form, together with a state appearing in that model)
 - R_a ($a \in \mathcal{A}$) are relations on Σ encoding the behaviour of atomic actions
 - \mathcal{C} is a function that gives the set of actions that the agent is able to do per model/state pair
 - Ag is a function that yields the set of actions that the agent is committed to (the agent's 'agenda') per model/state pair. \dashv

We have elements in the structures for interpreting the operators for knowledge, belief, and desire. Actions are modelled as model/state pair transformers to emphasize their influence on the mental state (that is, the complex of knowledge, belief and desires) of the agent rather than just the state of the world. Both (cap)abilities and commitments are given by functions that yield the relevant information per model / state pair.

Definition 10.8 (Interpretation of formulas in KARO)

In order to determine whether a formula $\varphi \in \mathcal{L}$ is true in a model/state pair (M, w) (if so, we write $(M, w) \models \varphi$), we stipulate:

- $\mathcal{M}, w \models p$ iff $\vartheta(w)(p) = \text{true}$, for $p \in \mathcal{P}$
- The logical connectives are interpreted as usual.
- $\mathcal{M}, w \models \mathbf{K}\varphi$ iff $\mathcal{M}, w' \models \varphi$ for all w' with $R_K(w, w')$
- $\mathcal{M}, w \models \mathbf{B}\varphi$ iff $\mathcal{M}, w' \models \varphi$ for all w' with $R_B(w, w')$
- $\mathcal{M}, w \models \mathbf{D}\varphi$ iff $\mathcal{M}, w' \models \varphi$ for all w' with $R_D(w, w')$
- $\mathcal{M}, w \models [\alpha]\varphi$ iff $\mathcal{M}', w' \models \varphi$ for all M', w' with $R_\alpha((\mathcal{M}, w), (\mathcal{M}', w'))$
- $\mathcal{M}, w \models \mathbf{A}\alpha$ iff $\alpha \in \mathcal{C}(\mathcal{M}, w)$ ⁵
- $\mathcal{M}, w \models \mathbf{Com}(\alpha)$ iff $\alpha \in \text{Ag}(\mathcal{M}, w)$ ⁶ ⊣

Here R_α is defined as usual in dynamic logic by induction from the basic case R_a , but now on model/state pairs rather than just states. So, e.g. $R_{\alpha_1 + \alpha_2} = R_{\alpha_1} \cup R_{\alpha_2}$, $R_{\alpha^*} = R_\alpha^*$, the reflective transitive closure of R_α , and $R_{\alpha_1; \alpha_2}$ is the relational product of R_{α_1} and R_{α_2} . Likewise the function \mathcal{C} is lifted to complex actions. We call an action α *deterministic* if $\#\{w' \mid R_\alpha(w, w')\} \leq 1$ for any $w \in W$, and *strongly deterministic* if $\#\{w' \mid R_\alpha(w, w')\} \leq 1$. (Here $\#$ stands for cardinality.)

We have clauses for knowledge, belief and desire. The action modality gets a similar interpretation: something (necessarily) holds after the performance / execution of action α if it holds in all the situations that are accessible from the current one by doing the action α . The only thing which is slightly nonstandard is that, as stated above, a situation is characterised here as a model / state pair. The interpretations of the ability and commitment operators are rather trivial in this setting (but see the footnotes): an action is enabled (or rather: the agent is able to do the action) if it is indicated so by the function \mathcal{C} , and, likewise, an agent is committed to an action α if it is recorded so in the agent's agenda.

Furthermore, we will make use of the following syntactic abbreviations serving as auxiliary operators:

Definition 10.9

- (dual) $\langle \alpha \rangle \varphi = \neg[\alpha]\neg\varphi$, expressing that the agent has the opportunity to perform α resulting in a state where φ holds.
- (opportunity) $\mathbf{O}\alpha = \langle \alpha \rangle \text{tt}$, i.e., an agent has the opportunity to do an action iff there is a successor state w.r.t. the R_α -relation;

⁵ In fact, the ability operator can alternatively be defined by means of a second accessibility relation for actions, in a way analogous to the opportunity operator below.

⁶ The agenda is assumed to be closed under certain conditions such as taking 'prefixes' of actions (representing initial computations). Details are omitted here, but see Section 10.6 for references.

- (practical possibility) $\mathbf{P}(\alpha, \varphi) = \mathbf{A}\alpha \wedge \mathbf{O}\alpha \wedge \langle \alpha \rangle \varphi$, i.e., an agent has the practical possibility to do an action with result φ iff it is both able and has the opportunity to do that action and the result of actually doing that action leads to a state where φ holds;
- (can) $\mathbf{Can}(\alpha, \varphi) = \mathbf{KP}(\alpha, \varphi)$, i.e., an agent can do an action with a certain result iff it knows it has the practical possibility to do so;
- (realisability) $\diamond \varphi = \exists a_1, \dots, a_n \mathbf{P}(a_1; \dots; a_n, \varphi)$ ⁷, i.e., a state property φ is realisable iff there is a finite sequence of atomic actions of which the agent has the practical possibility to perform it with the result φ ;
- (goal) $\mathbf{G}\varphi = \neg \varphi \wedge \mathbf{D}\varphi \wedge \diamond \varphi$, i.e., a goal is a formula that is not (yet) satisfied, but desired and realisable.⁸
- (possible intend) $\mathbf{I}(\alpha, \varphi) = \mathbf{Can}(\alpha, \varphi) \wedge \mathbf{KG}\varphi$, i.e., an agent (possibly) intends an action with a certain result iff the agent can do the action with that result and it moreover knows that this result is one of its goals. \dashv

Remark 10.1

- The dual of the (box-type) action modality expresses that there is at least a resulting state where a formula φ holds. It is important to note that in the context of *deterministic* actions, i.e. actions that have at most one successor state, this means that the *only* state satisfies φ , and is thus in this particular case a stronger assertion than its dual formula $[\alpha]\varphi$, which merely states that if there are any successor states they will (all) satisfy φ . Note also that if atomic actions are assumed to be deterministic all actions including the complex ones will be deterministic.
- Opportunity to do an action is modelled by having at least one successor state according to the accessibility relation associated with the action.
- Practical possibility to to an action with a certain result is modelled as having both ability and opportunity to do the action with the appropriate result. Note that $\mathbf{O}\alpha$ in the formula $\mathbf{A}\alpha \wedge \mathbf{O}\alpha \wedge \langle \alpha \rangle \varphi$ is actually redundant since it already follows from $\langle \alpha \rangle \varphi$. However, to stress the opportunity aspect it is added.
- The Can predicate applied to an action and formula expresses that the agent is ‘conscious’ of its practical possibility to do the action resulting in a state where the formula holds.

⁷We abuse our language here slightly, since strictly speaking we do not have quantification in our object language. See our references to KARO in Section 10.6 for a proper definition.

⁸In fact, here we simplify matters slightly. One might stipulate that a goal should be explicitly selected somehow from the desires it has, which could be modelled by means of an additional modal operator. Here we leave this out for simplicity’s sake.

- A formula φ is realisable if there is a ‘plan’ consisting of (a sequence of) atomic actions of which the agent has the practical possibility to do them with φ as a result.
- A formula φ is a goal in the KARO framework if it is not true yet, but desired and realisable in the above meaning, that is, there is a plan of which the agent has the practical possibility to realise it with φ as a result.
- An agent is said to (possibly) intend an action α with result φ if it ‘Can’ do this (knows that it has the practical possibility to do so), and, moreover, knows that φ is a goal. ¬

In order to manipulate both knowledge / belief and motivational matters special actions **revise**, **commit** and **uncommit** are added to the language. (We assume that we cannot nest these operators. So, e.g., **commit(uncommit α)** is not a well-formed action expression.) The semantics of these are again given as model/state transformers (We only do this here in a very abstract manner, viewing the accessibility relations associated with these actions as functions. For further details we refer the reader to inspect some of the KARO references mentioned in Section 10.6.

Definition 10.10 (Accessibility of revise, commit and uncommit actions)

1. $R_{\text{revise}\varphi}(\mathcal{M}, w) = \text{update_belief}(\varphi, (\mathcal{M}, w))$.
2. $R_{\text{commit}\alpha}(\mathcal{M}, w) = \text{update_agenda}^+(\alpha, (\mathcal{M}, w))$, if $\mathcal{M}, w \models \mathbf{I}(\alpha, \varphi)$ for some φ , otherwise $R_{\text{commit}\alpha}(\mathcal{M}, w) = \emptyset$ (indicating failure of the commit action).
3. $R_{\text{uncommit}\alpha}(\mathcal{M}, w) = \text{update_agenda}^-(\alpha, (\mathcal{M}, w))$, if $\mathcal{M}, w \models \mathbf{Com}(\alpha)$, otherwise $R_{\text{uncommit}\alpha}(\mathcal{M}, w) = \emptyset$ (indicating failure of the uncommit action);
4. $\text{uncommit}\alpha \in \mathcal{C}(\mathcal{M}, w)$ iff $\mathcal{M}, w \models \neg\mathbf{I}(\alpha, \varphi)$ for all formulas φ , that is, an agent is able to uncommit to an action if it is not intended to do it (any longer) for any purpose. ¬

Here update_belief , update_agenda^+ and update_agenda^- are functions that update the agent’s belief and agenda (by adding or removing an action), respectively. Details are omitted here, but essentially these actions are model/state transformers again, representing a change of the mental state of the agent (regarding beliefs and commitments, respectively). The $\text{update_belief}(\varphi, (\mathcal{M}, w))$ function changes the model \mathcal{M} in such a way that the agent’s belief is updated with the formula φ , while $\text{update_agenda}^+(\alpha, (\mathcal{M}, w))$ changes the model \mathcal{M} such that α is added to the agenda, and likewise for the update_agenda^- function, but now with respect to removing an action from the agenda. The **revise** operator can be used to cater for revisions due to observations and communication with other agents, which we will not go into further here.

The interpretation of formulas containing revise and (un)commit actions is now done using the accessibility relations above. One can now define validity as usual with respect to the KARO-models. One then obtains the following validities (of course, in order to be able to verify these one should use the proper model and not the abstraction / simplification we have presented here.) Typical properties of this framework, called the KARO logic, include:

Proposition 10.3

1. $\models \Box(\varphi \rightarrow \psi) \rightarrow (\Box\varphi \rightarrow \Box\psi)$, for $\Box \in \{\mathbf{K}, \mathbf{B}, \mathbf{D}, [\alpha]\}$
2. $\models \langle \alpha \rangle \varphi \rightarrow [\alpha] \varphi$, for deterministic α
3. $\models \Box\varphi \rightarrow \Box\Box\varphi$, for $\Box \in \{\mathbf{K}, \mathbf{B}\}$
4. $\models \neg\Box\varphi \rightarrow \Box\neg\Box\varphi$, for $\Box \in \{\mathbf{K}, \mathbf{B}\}$
5. $\models \mathbf{K}\varphi \rightarrow \varphi$
6. $\models \neg\mathbf{B}ff$
7. $\models \mathbf{O}(\alpha; \beta) \leftrightarrow \langle \alpha \rangle \mathbf{O}\beta$
8. $\models \mathbf{Can}(\alpha; \beta, \varphi) \leftrightarrow \mathbf{Can}(\alpha, \mathbf{P}(\beta, \varphi))$
9. $\models \mathbf{I}(\alpha, \varphi) \rightarrow \mathbf{K}\langle \alpha \rangle \varphi$
10. $\models \mathbf{I}(\alpha, \varphi) \rightarrow \langle \text{commit} \alpha \rangle \mathbf{Com}(\alpha)$
11. $\models \mathbf{I}(\alpha, \varphi) \rightarrow \neg \mathbf{Auncommit}(\alpha)$
12. $\models \mathbf{Com}(\alpha) \rightarrow \langle \text{uncommit}(\alpha) \rangle \neg \mathbf{Com}(\alpha)$
13. $\models \mathbf{Com}(\alpha) \wedge \neg \mathbf{Can}(\alpha, \top) \rightarrow \mathbf{Can}(\text{uncommit}(\alpha), \neg \mathbf{Com}(\alpha))$
14. $\models \mathbf{Com}(\alpha) \rightarrow \mathbf{KCom}(\alpha)$
15. $\models \mathbf{Com}(\alpha_1; \alpha_2) \rightarrow \mathbf{Com}(\alpha_1) \wedge \mathbf{K}[\alpha_1] \mathbf{Com}(\alpha_2)$
16. $\models \mathbf{Com}(\text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi}) \wedge \mathbf{K}\varphi \rightarrow \mathbf{Com}(\varphi?; \alpha_1)$
17. $\models \mathbf{Com}(\text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi}) \wedge \mathbf{K}\neg\varphi \rightarrow \mathbf{Com}(\neg\varphi?; \alpha_2)$
18. $\models \mathbf{Com}(\text{while } \varphi \text{ do } \alpha \text{ od}) \wedge \mathbf{K}\varphi \rightarrow \mathbf{Com}((\varphi?; \alpha); \text{while } \varphi \text{ do } \alpha \text{ od}) \quad \dashv$

The first of these properties says that all the modalities mentioned are ‘normal’ in the sense that they are closed under implication. The second states that the dual operator $\langle \alpha \rangle$ is stronger than the operator $[\alpha]$ in case the action α is deterministic: if there is at most one successor state after performing α and we know that there is at least one successor state satisfying φ then *all* successor states satisfy φ . The third and fourth properties are the so-called introspection properties for knowledge and belief. The fifth property says that knowledge is true, while the sixth states that belief (may not be true but) is not inconsistent. The seventh property states that having the opportunity to do a sequential composition of two actions amounts to having the opportunity of doing the first action first and then having the opportunity to do the second. The eighth states that an agent that *can* do a sequential composition of two actions with result φ iff the agent can do the first actions resulting in a state where it has the practical possibility to do the second with φ as result. The ninth states that if one possibly intends to do α with result φ then one knows that there is a possibility of performing α resulting in a state where φ holds. The tenth asserts that if an agent possibly intends to do α with some result φ , it has the opportunity to commit to α with result that it is committed to α (i.e. α is put into its agenda). The eleventh says

that if an agent intends to do α with a certain purpose, then it is unable to uncommit to it (so, if it is committed to α it has to persevere with it). This is the way persistence of commitment is represented in KARO. Note that this is much more ‘concrete’ (also in the sense of computability) than the persistence notions in the other approaches we have seen, where temporal operators pertaining to a possibly infinite future were employed to capture them...! We think it is no coincidence that Hindriks *et al.* established an almost perfect match in the sense of a correspondence between the agent programming language GOAL and Cohen & Levesque’s logic of intention, the main difference being the inability of GOAL to express the persistence properties of intentions in this logic...!) In KARO we have the advantage of having dedicated actions in the action language dealing with the change of commitment that can be used to express persistence without referring to the (infinite) future, rendering the notion of persistence much ‘more computable’. The twelfth property says that if an agent is committed to an action and it has the opportunity to uncommit to it, as result then indeed the commitment is removed. The thirteenth says that whenever an agent is committed to an action that is no longer known to be practically possible, it knows that it can undo this impossible commitment. The fourteenth property states that commitments are known to the agent. The last four properties have to do with commitments to complex actions. For instance, the fifteenth says that if an agent is committed to a sequential composition of two actions then it is committed to the first one, and it knows that after doing the first action it will be committed to the second action.

KARO logic for emotional agents

In this subsection we look at a recent application of BDI logic that deals with agent behaviour that is strictly beyond the scope of the original aim of BDI logic, viz. describing the behaviour of rational agents. We will sketch how the KARO framework can be used for describing emotional agents. Although it is perhaps a bit paradoxical to describe emotions and emotional behaviour with logic, one should bear in mind that we are dealing with behaviour here, and this can be described in logic, especially a logic that deals with actions such as the KARO framework. Furthermore, as we shall see, emotional behaviour will turn out to be complimentary rather than opposed to rational behaviour of agents, something that is also acknowledged by recent work in cognitive science through the work of for instance Damasio. Our presentation here is inspired by two psychological theories: that of Oatley & Jenkins and that of OCC. Since the latter is much more involved (treating 22 emotions, while the former only treats 4 basic emotions), we here mainly follow the ideas of Oatley & Jenkins, and say a few words on modelling OCC later.

According to Oatley & Jenkins, the 4 basic emotions, happiness, sadness, anger and fear, have the following characteristics:

- Happiness results when in the process of trying to achieve a goal, things go ‘right’, as expected, i.e., subgoals are achieved thus far.
- Sadness results when in the process of trying to achieve a goal, things go ‘wrong’, i.e., not as expected, i.e., subgoals are not being achieved.

- Anger is the result of frustration about not being able to execute the current plan, and makes the agent try harder to execute the plan.
- Fear results when a ‘maintenance goal’ is threatened, so that the agent will make sure that this maintenance goal is restored before going on with other activities.

It is directly obvious from these descriptions that these emotions are BDI-related notions! So it is not so strange to use a BDI-logic like KARO to describe them, which is what Meyer did.

Let us take sadness as an example. For simplicity, assume that plans consist of sequences of atomic actions. In KARO we can then express the trigger condition for sadness as follows:

$$\mathbf{I}(\pi, \varphi) \wedge \mathbf{Com}(\pi) \wedge \mathbf{B}([\alpha]\psi) \rightarrow \\ [\alpha](\mathbf{B}\neg\psi \wedge \mathbf{Com}(\pi \setminus \alpha)) \rightarrow \mathit{sad}(\pi \setminus \alpha, \varphi)$$

where α is a prefix of plan π . Intuitively, this says that if the agent has the (possible) intention to perform plan π with goal φ , it is committed to π (so it has a true intention to do π), and it believes that after doing the initial fragment α of the plan π it holds that ψ , then after doing α if it believes that ψ does not hold while it is still committed to the rest of the plan, it is sad (with respect to the rest of the plan and goal φ). In a similar way the trigger conditions of the other emotions can be formalised.

Also, together with Steunebrink and Dastani, Meyer looked at modeling OCC. Particularly, they show how to formalise the (trigger conditions of) emotions in OCC in three steps: first by presenting a more general logical structure of the emotions, which are later refined in terms of doxastic logic and finally in the full-blown BDI logic KARO again. The way emotions get a semantics based on BDI models is quite intricate and beyond the scope of this chapter, but one of the properties that can be proven valid in this approach is the following, using KARO’s (possible) intend operator (here parametrized by an agent):

$$\mathbf{I}_i(\alpha, \varphi) \rightarrow [i : \alpha](\mathit{Pride}_i^T(i : \alpha) \wedge \mathit{Joy}_i^T(\varphi) \wedge \mathit{Gratification}_i^T(i : \alpha, \varphi))$$

(Here $i : \alpha$ in the dynamic logic box refers to the action of i performing α , and the superscript T placed at the emotion operators pertains to the idea that we are considering triggering / elicitation forms of the emotions concerned.) Informally this reads that if the agent i has the possible intention to do α with goal φ , then if he has performed α he is proud (triggered pride) of his action, has (triggered) joy about the achievement of the goal φ and has (triggered) gratification with respect to action $i : \alpha$ and goal φ .

Finally, let us also mention here the strongly related work by Adam *et al.* Also this work is devoted to a formalisation of OCC emotions in BDI terms. There are differences with the work of Steunebrink *et al.*, though. For instance, Adam simply defines joy as a conjunction of belief and desire: $\mathit{Joy} \varphi =_{def} \mathit{Bel} \varphi \wedge \mathit{Des} \varphi$. This seems to express a ‘state of joy’ (experience) rather than a trigger for joy. This raises a confusion of emotion elicitation (triggering) and experience, which is kept separately in the approach of Steunebrink. This confusion also appears at other

places in the work of Adam, e.g. where she defines gratification as the conjunction of pride (which pertains to triggering) and joy (which is about experience as we saw earlier). In later work by Adam, this issue is improved upon and it is explained that the above definition of Joy is solely about the triggering of joy, not the experience. However, the confusion of triggering versus experience is still not resolved completely since it is still present in the introspective properties of emotional awareness $Emotion\varphi \leftrightarrow BelEmotion\varphi$ and $\neg Emotion\varphi \leftrightarrow Bel\neg Emotion\varphi$, which hold in Adam's framework, for any Emotion. This is counterintuitive if Emotion should capture the triggering of the associated emotion, since an agent may not be aware of this triggering.

10.4 BDI-modalities in STIT logic

The principles of BDI logics reflect rationality postulates for agent modalities. In particular, the BDI principles model how B, D and I modalities interact with each other over time (well known are the so called 'commitment strategies' of Cohen & Levesque, stating under which belief and desire conditions intentions have to be dropped, see Section 10.2.1). BDI logics are not meant for knowledge representation but for agent specification: ideally concretely built agents will some day be verified against the logic principles of BDI-logics (how exactly this could ever be done is a question we set aside here).

An essential component of any BDI logic is then its dynamic part. Traditionally, either the dynamic part is formed by a dynamic logic fragment (Cohen & Levesque, KARO) or a temporal fragment (Rao & Georgeff). Recently a third alternative has been considered: STIT (seeing to it that) logic. STIT logics can be said to be in between dynamic logic and temporal logic. Where dynamic logic sees actions as the steps of a program, and temporal logic leaves actions entirely out of the picture, STIT logic sees action as a relation between agents and the effects they can see to. STIT logic achieves this by generalizing temporal structures to choice structures. The most distinguishing feature of STIT logic is that truth of formulas often expresses information about the dynamics of the world. For instance, the STIT logic formula $[ag\ stit]X(at_station)$ says that agent ag currently sees to it that next it is at the station. But, it does *not* say that the agent *can* see to it that next it is at the station (this is however a logical consequence). Abilities are truths about static conditions, and not about dynamic conditions.

In the present section we will discuss how in recent years several authors have aimed to combine STIT logic and BDI notions. There are two parts. In the first part, Section 10.4.1, we focus on classical instantaneous STIT logics and the BDI extensions that have been suggested for them. In the second part, Section 10.4.2, we consider dynamic variants of the BDI modalities and discuss the notion of 'knowingly doing' within a version of STIT where effects of actions take effect in next states: XSTIT.

There is a strong connection, both conceptually and technically, between the family of logics STIT and *Alternating-time Temporal Logic*, which is at the center of attention in Chapter 11.

10.4.1 BDI modalities in instantaneous stit

Traditionally STIT logics encompass operators for agency that assume that an agentive choice performance is something that takes no time. So, an instantaneous stit operator $[ag \text{ stit}]\varphi$ typically obeys the success axiom $[ag \text{ stit}]\varphi \rightarrow \varphi$ to capture the intuition concerning instantaneity saying that if *ag now* sees to it that φ holds, then φ must be true *now*. Before putting forward an alternative to this view, where the central agency operator has a built-in step to a next moment in time, we give the formal definition of standard (Chellas) instantaneous STIT logic and discuss its logical properties. We will here use a slightly different syntax and semantics than used by Chellas himself and also different from that of Horty, but, the logic is the same.

CSTIT

Definition 10.11

Given a countable set of propositions P and $p \in P$, and given a finite set Ags of agent names, and $ag \in Ags$, the formal language $\mathcal{L}_{\text{CSTIT}}$ is:

$$\varphi := p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \Box\varphi \mid [ag \text{ Cstit}]\varphi \quad \dashv$$

Besides the usual propositional connectives, the syntax of CSTIT comprises two modal operators. The operator $\Box\varphi$ expresses ‘historical necessity’, and plays the same role as the well-known path quantifiers in logics such as CTL and CTL*. Another way of talking about this operator is to say that it expresses that φ is ‘settled’. We abbreviate $\neg\Box\neg\varphi$ by $\Diamond\varphi$. The operator $[ag \text{ Cstit}]\varphi$ stands for ‘agent *ag* sees to it that φ ’ (the ‘C’ referring to Chellas). $\langle ag \text{ Cstit} \rangle\varphi$ abbreviates $\neg[ag \text{ Cstit}]\neg\varphi$.

The semantics given in Definition 10.12 below is an alternative to the semantics given by Belnap and colleagues in terms of BT+AC structures (Branching Time + Agentive Choice structures). The differences are not essential. Where the branching of time in BT + AC structures is represented by tree-like orderings of moments, the structures below use ‘bundles’ of linearly ordered sets of moments. We use the latter to be uniform in the semantic structures across different STIT formalisms in this section.

Definition 10.12

A CSTIT-frame is a tuple $\langle S, H, R_{ag} \rangle$ such that⁹:

1. S is a non-empty set of static states. Elements of S are denoted s, s' , etc.
2. H is a non-empty set of possible system histories $\dots s_{-2}, s_{-1}, s_0, s_1, s_2, \dots$ with $s_x \in S$ for $x \in \mathbb{Z}$. Elements of H are denoted h, h' , etc.
3. Dynamic states are tuples $\langle s, h \rangle$, where $s \in S, h \in H$ and s appears on h . Now the relations R_{ag} are ‘effectivity’ equivalence classes over dynamic states such that $\langle s, h \rangle R_{ag} \langle s', h' \rangle$ only if $s = s'$. For any state s and agent ag , the relation R_{ag} defines a partition of the dynamic states built with s .

⁹In the meta-language we use the same symbols both as constant names and as variable names, and we assume universal quantification of unbound meta-variables.

The partition models the possible choices $C_{ag}^s, C'_{ag}^s, C''_{ag}^s, \dots$ of ag in s . A choice profile $\langle C_{ag_1}^s, C_{ag_2}^s \dots C_{ag_n}^s \rangle$ at s is a particular combination of choices $C_{ag_i}^s$ at s , one for each agent ag_i in the system. For any s the intersection of choices in any choice profile is non-empty: $\bigcap_{ag_i \in Ags} C_{ag_i}^s \neq \emptyset$

In Definition 10.12 above, we refer to the states s as ‘static states’. This is to distinguish them from ‘dynamic states’, which are combinations $\langle s, h \rangle$ of a static state and a history. Dynamic states function as the elementary units of evaluation of the logic. This means that the basic notion of ‘truth’ in the semantics of this logic is about dynamic conditions concerning choice performances.

We now define models by adding a valuation of propositional atoms to the frames of Definition 10.12.

Definition 10.13

A frame $\mathcal{F} = \langle S, H, R_{ag} \rangle$ is extended to a model $\mathcal{M} = \langle S, H, R_{ag}, V \rangle$ by adding a valuation V of atomic propositions:

- V is a valuation function $V : P \rightarrow 2^{S \times H}$ assigning to each atomic proposition the set of state history pairs relative to which they are true.

We evaluate truth with respect to dynamic states.

Definition 10.14

Relative to a model $\mathcal{M} = \langle S, H, R_{ag}, V \rangle$, truth $\langle s, h \rangle \models \varphi$ of a formula φ in a dynamic state $\langle s, h \rangle$, with $s \in h$, is defined as:

$$\begin{aligned}
 \langle s, h \rangle \models p & \quad \Leftrightarrow \quad \langle s, h \rangle \in V(p) \\
 \langle s, h \rangle \models \neg \varphi & \quad \Leftrightarrow \quad \text{not } \langle s, h \rangle \models \varphi \\
 \langle s, h \rangle \models \varphi \wedge \psi & \quad \Leftrightarrow \quad \langle s, h \rangle \models \varphi \text{ and } \langle s, h \rangle \models \psi \\
 \langle s, h \rangle \models \Box \varphi & \quad \Leftrightarrow \quad \forall h' : \text{if } s \in h' \text{ then } \langle s, h' \rangle \models \varphi \\
 \langle s, h \rangle \models [ag \text{ Cstit}] \varphi & \quad \Leftrightarrow \quad \forall h' : \text{if } \langle s, h \rangle R_{ag} \langle s, h' \rangle \text{ then } \langle s, h' \rangle \models \varphi \quad \dashv
 \end{aligned}$$

Satisfiability, validity on a frame and general validity are defined as usual.

Now we proceed with the axiomatization.

Theorem 10.2

The following axiom schemes, in combination with a standard axiomatization for propositional logic, and the standard rules (like necessitation) for the normal modal operators, define a complete Hilbert system for CSTIT:

$$\begin{aligned}
 & \text{The S5 axioms for } \Box \\
 & \text{For each } ag \text{ the S5 axioms for } [ag \text{ Cstit}] \\
 (SettC) & \quad \Box \varphi \rightarrow [ag \text{ Cstit}] \varphi \\
 (Indep) & \quad \Diamond [ag_1 \text{ Cstit}] \varphi \wedge \dots \wedge \Diamond [ag_n \text{ Cstit}] \psi \rightarrow \\
 & \quad \Diamond ([ag_1 \text{ Cstit}] \varphi \wedge \dots \wedge [ag_n \text{ Cstit}] \psi) \\
 & \quad \text{for } Ags = \{ag_1, \dots, ag_n\} \quad \dashv
 \end{aligned}$$

Balbani *et al.* propose an alternative axiomatization and a semantics whose units of evaluation are not two dimensional pairs $\langle s, h \rangle$ but one dimensional worlds w . Here we have chosen to give a two-dimensional semantics to emphasize the relation with the XSTIT semantics in section 10.4.2.

BDI-stit

Semmling and Wansing add BDI modalities to a basic Chellas stit logic as the one just defined. Their BDI-stit formalism extends the syntax as follows (we take the liberty of using our own notation for the BDI operators and to define an alternative but equivalent semantics).

Definition 10.15

Given a countable set of propositions P and $p \in P$, and given a finite set Ags of agent names, and $ag \in Ags$, the formal language $\mathcal{L}_{\text{bdi-stit}}$ is:

$$\varphi := p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \Box\varphi \mid [ag \text{ Cstit}]\varphi \mid \langle [ag \text{ bel}]\rangle\varphi \mid \langle [ag \text{ des}]\rangle\varphi \mid \langle [ag \text{ int}]\rangle\varphi \quad \neg$$

To emphasize their weak modal character, we denote the introduced belief, desire, intention and possibility operators with a combination of sharp and square brackets. This alludes to the combination of first order existential and universal quantifications that is present in any first order simulation of a weak modal operator. The reading of the operators speaks for itself; they express belief, desire and intention concerning a proposition φ .

Definition 10.16

A bdi-stit-frame is a tuple $\langle S, H, R_{ag}, N_b, N_d, N_i \rangle$ such that:

1. $\langle S, H, R_{ag} \rangle$ is a CSTIT-frame
2. N_b, N_d and N_i are neighborhood functions of the form $N : S \times H \times Ags \mapsto 2^{2^{S \times H}}$ mapping any combination of a dynamic state $\langle s, h \rangle$ and an agent ag to a set of neighborhoods of $\langle s, h \rangle$. Semmling & Wansing then impose constraints on neighborhood frames that are equivalent to:
 - a. All three functions N_b, N_d and N_i obey $\emptyset \notin N(s, h, ag)$
 - b. All three functions N_b, N_d and N_i obey that if $N \in N(s, h, ag)$ and $N \subset N'$ then $N' \in N(s, h, ag)$
 - c. $N \in N_i(s, h, ag)$ and $N' \in N_i(s, h, ag)$ implies $N \cap N' \neq \emptyset$

The intuition underlying neighborhood functions is the following. $N_b(s, h, ag)$ gives for agent ag in situation $\langle s, h \rangle$ the clusters of possible worlds (situations / dynamic states) the joint possibility of which it believes in. Since clusters and propositions correspond to each other one-to-one (modulo logic equivalence of the propositions), it will also be convenient to look at the clusters or neighborhoods as propositions and to say that if $N \in N_b(s, h, ag)$ the agent ag believes the proposition (modulo logical equivalence) corresponding to N , that $N \in N_d(s, h, ag)$ holds if ag desires the proposition and that $N \in N_i(s, h, ag)$ holds if ag intends the proposition.

Now **a.** says that there is no belief, desire or intention for impossible states of affairs. **b.** says that belief, desire and intention are closed under weakening of the propositions believed, desired or intended. **c.** says that intentions are consistent in the sense that it is not possible to hold at the same time an intention for a proposition and for its negation.

Definition 10.17 (Truth conditions BDI operators)

Relative to a model $\langle S, H, R_{ag}, N_b, N_d, N_i, V \rangle$, truth of belief, desire and intention operators is defined as ($\llbracket \varphi \rrbracket$ is the truth set of φ , that is, the subset of all dynamic elements in $S \times H$ satisfying φ):

$$\begin{aligned} \langle s, h \rangle \models \langle [ag \text{ bel}] \rangle \varphi &\Leftrightarrow \llbracket \varphi \rrbracket \in N_b(s, h, ag) \\ \langle s, h \rangle \models \langle [ag \text{ des}] \rangle \varphi &\Leftrightarrow \llbracket \varphi \rrbracket \in N_d(s, h, ag) \\ \langle s, h \rangle \models \langle [ag \text{ int}] \rangle \varphi &\Leftrightarrow \llbracket \varphi \rrbracket \in N_i(s, h, ag) \end{aligned} \quad \dashv$$

An axiomatization of a probabilistic epistemic logic is obtained by formulating axioms corresponding to the conditions on neighborhood functions.

Theorem 10.3 (Hilbert system BDI operators)

Relative to the semantics following from definitions 10.16 and 10.17 we define the following Hilbert system. We assume the standard derivation rules for the weak modalities, like closure under logical equivalence.

$$\begin{array}{ll} (BelPos) & \neg \langle [ag \text{ bel}] \rangle \perp \\ (DesPos) & \neg \langle [ag \text{ des}] \rangle \perp \\ (IntPos) & \neg \langle [ag \text{ int}] \rangle \perp \\ (BelWk) & \langle [ag \text{ bel}] \rangle \varphi \rightarrow \langle [ag \text{ bel}] \rangle (\varphi \vee \psi) \\ (DesWk) & \langle [ag \text{ des}] \rangle \varphi \rightarrow \langle [ag \text{ des}] \rangle (\varphi \vee \psi) \\ (IntWk) & \langle [ag \text{ int}] \rangle \varphi \rightarrow \langle [ag \text{ int}] \rangle (\varphi \vee \psi) \\ (IntD) & \langle [ag \text{ int}] \rangle \varphi \rightarrow \neg \langle [ag \text{ int}] \rangle \neg \varphi \end{array} \quad \dashv$$

Relative to their own version of the semantics, Semmling & Wansing prove completeness of their logic. Here the completeness of the axiomatization relative to the frames of Definition 10.16 follows from general results in neighborhood semantics and monotonic modal logic. One can check that the conditions on the frames correspond one-to-one with the axioms in the axiomatization.

As can be seen from the axioms and conditions we have shown above, Semmling and Wansing chose to make their BDI-stit logic rather weak, trying to commit only to a minimum of logical properties. But even with this minimalistic approach there is room for debate. For instance, the condition on intentions only looks at pairwise consistency of intentions, but conflicts are still possible in case there is a combination of three intentions: $\{\langle [ag \text{ int}] \rangle \varphi, \langle [ag \text{ int}] \rangle (\varphi \rightarrow \psi), \langle [ag \text{ int}] \rangle \neg \psi\}$ is satisfiable. If we do not want that, it is straightforward to adapt the condition on neighborhood functions (demand that any finite number of neighborhoods has a state in common), but it is unclear how to axiomatise it.

Even though the STIT framework's notion of truth refers to the dynamics of a system of agents (which is why we talk about 'dynamic states'), Semmling and Wansing do not focus on dynamic interpretations of the BDI attitudes. A formula like $\langle [ag \text{ bel}] \rangle [ag \text{ Cstit}] \varphi$ must express something like "ag believes that it sees to it that φ ", but the inherent dynamic aspect of this notion is not analyzed. In particular, no interactions between STIT and BDI modalities are studied. This is likely to be due to the fact that explicit dynamic temporal operators are absent in the logic and because agency is instantaneous. In Section 10.4.2 we will report on the study of the inherent dynamic aspect of such combinations of operators.

BDI, STIT and regret

Lorini and Schwarzenruber use STIT logic as the basis for investigations into what they call counterfactual emotions¹⁰. The typical counterfactual emotion is ‘regret’. Regret can be described as a discrepancy between what actually occurs and what could have happened. Based on this they argue that for a definition of regret in the STIT framework, it needs to be extended with modalities for knowledge and desire. They consider three different STIT formalisms to base their definitions on, but here it will suffice to discuss their ideas using the Chellas STIT logic given earlier.

Knowledge is added to their STIT framework in a straightforward way: a normal S5 knowledge operator $[ag_i \text{ kno}]$ is added for every agent ag_i in the system. The interpretation is in terms of equivalence classes over the basic units of evaluation (for the stit language we consider here: dynamic states). So, for knowledge we get the truth condition $\langle s, h \rangle \models [ag_i \text{ kno}]\varphi$ iff for all s', h' such that $\langle s, h \rangle \sim \langle s', h' \rangle$ we have $\langle s', h' \rangle \models \varphi$. The second operator we take from their system is an operator for desire, which is defined using propositional constants.

Definition 10.18

Let $good_{ag_i}$ denote a propositional constant, one for each agent ag_i in the system, whose truth expresses that a state is good for that agent. Now the modal operators $[ag_i \text{ good}]\varphi$ and $[ag_i \text{ des}]\varphi$ are defined by:

$$\begin{aligned} [ag_i \text{ good}]\varphi &\equiv_{def} \Box(good_{ag_i} \rightarrow \varphi) \\ [ag_i \text{ des}]\varphi &\equiv_{def} [ag_i \text{ kno}][ag_i \text{ good}]\varphi \quad \dashv \end{aligned}$$

The counterfactual aspect is introduced by the definition of the notion of "could have prevented" (CHP). For definitions of such counterfactual properties, the STIT framework is more suited than dynamic logic or situation calculus frameworks, since in STIT we reason about actual performances of actions which also makes it possible to reason about choices that are not (f)actual. Lorini and Schwarzenruber define their notion of CHP in a group stit framework. Here we have only defined individual agency. Therefore we will only consider the two agent case, for which there is no essential difference between group operators and operators for individual agents. The two agent case will enable us to discuss the ideas properly.

In the two agent setting, Lorini and Schwarzenruber’s intuition can be described as follows: agent ag_1 could have prevented φ if and only if (1) φ is currently true, and (2) agent ag_2 does not see to it that φ . They reformulate the second condition as ‘provided that agent ag_2 sticks to its choice, for agent ag_1 there is the possibility to act otherwise in a way that does not guarantee the outcome φ ’¹¹.

¹⁰However, it is not the emotions that are counterfactual; the theory is about factual emotions based on beliefs about counterfactual conditions

¹¹The reformulation is equivalent, but suggests a meaning of group action that is not standard: if ag_1 only has an alternative under the provision that the choice of ag_2 is kept fixed, then a standard interpretation would be that ag_1 and ag_2 have an alternative in *cooperation*. But, then the possibility to prevent also relies on cooperation. What seems needed for a notion of ‘could have prevented’ for individual agents, is that such agents have alternatives *individually*.

Definition 10.19

$$\langle [ag_1 \text{ CHP}] \rangle \varphi \equiv_{def} \varphi \wedge \neg [ag_2 \text{ Cstit}] \varphi$$

CHP is not a normal modality (which is why we use the combination of sharp and square brackets, as explained before). Lorini and Schwarzentruher argue that their notion of CHP obeys agglomeration ' $\langle [ag \text{ CHP}] \rangle \varphi \wedge \langle [ag \text{ CHP}] \rangle \psi \rightarrow \langle [ag \text{ CHP}] \rangle (\varphi \wedge \psi)$ ', but not weakening ' $\langle [ag \text{ CHP}] \rangle \varphi \rightarrow \langle [ag \text{ CHP}] \rangle (\varphi \vee \psi)$ '. But, agglomeration is not necessarily an intuitive property for a notion of 'could have prevented': if the guard could have prevented prisoner 1 to escape through exit 1 and if the guard could have prevented prisoner 2 to escape through exit 2, it does not follow that the guard could have prevented both prisoner 1 to escape through exit 1 and prisoner 2 to escape through exit 2. We believe alternative definitions of the notion of 'could have prevented' are possible. In particular it seems promising to look for notions where the counterfactual conditions concerning alternatives are made explicit using the \diamond operator.

With the right concepts in place, finally we are able to define the notion of regret for a proposition φ as the desire for $\neg\varphi$ in conjunction with knowledge about the fact that φ could have been prevented.

Definition 10.20

$$\langle [ag_i \text{ rgt}] \rangle \varphi \equiv_{def} [ag_i \text{ des}] \neg\varphi \wedge [ag_i \text{ kno}] \langle [ag_i \text{ CHP}] \rangle \varphi$$

Definitions like the one above for regret hinge on the possibility to express that certain actions or choices are actually performed. This type of expressivity is not provided by many other action formalisms.

10.4.2 BDI modalities in XSTIT: dynamic attitudes

In the systems discussed in Section 10.4.1, BDI notions were combined with STIT operators. However, in both approaches there was no special attention for a dynamic interpretation of the combination of BDI and STIT operators. Yet this interpretation strongly suggests itself. If an operator like $[ag \text{ Cstit}] \varphi$ expresses that agent ag now exercises his choice to ensure that φ , and a knowledge operator like $[ag \text{ kno}] \varphi$ also has a dynamic reading (note that the truth condition of knowledge is not with respect to static states s but with respect to dynamic states $\langle s, h \rangle$, cf. the clause for $[ag \text{ kno}] \varphi$ just above Definition 10.18), then a natural interpretation of a combination like $[ag \text{ kno}] [ag \text{ Cstit}] \varphi$ is that agent ag "knowingly sees to it that φ ". We suspect that this dynamic reading was not suggested by the authors of the discussed systems because these systems do not contain temporal modalities and because the used version of STIT has instantaneous effects. In this section we report on work studying the notions of knowingly and intentionally doing that is based on a version of STIT where agency inherently involves a move to some next state: XSTIT.

XSTIT

We give here the basic definitions for XSTIT. XSTIT enriches the CSTIT language of Section 10.4.1 with a temporal next operator and replaces the operator $[ag \text{ Cstit}] \varphi$

by the operator $[ag \text{ xstit}] \varphi$ where the effect φ is not instantaneous but occurs in a next state. Further explanations and motivations can be found elsewhere (see Section 10.6).

Definition 10.21

Given a countable set of propositions P and $p \in P$, and given a finite set Ags of agent names, and $ag \in Ags$, the formal language $\mathcal{L}_{\text{XSTIT}}$ is:

$$\varphi := p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \Box\varphi \mid [ag \text{ xstit}]\varphi \mid X\varphi \quad \dashv$$

The operator $[ag \text{ xstit}]\varphi$ stands for ‘agent ag sees to it that φ in the next state’. The operator $X\varphi$ is a standard next time operator. \Box is again the operator for historical necessity (settledness).

Definition 10.22

An XSTIT-frame is a tuple $\langle S, H, E \rangle$ such that:

1. S is a non-empty set of static states. Elements of S are denoted s, s' , etc.
2. H is a non-empty set of possible system histories $\dots s_{-2}, s_{-1}, s_0, s_1, s_2, \dots$ with $s_x \in S$ for $x \in \mathbb{Z}$. Elements of H are denoted h, h' , etc. We denote that s' succeeds s on the history h by $s' = \text{succ}(s, h)$ and by $s = \text{prec}(s', h)$. We have the following bundling constraint on the set H :
 - a. if $s \in h$ and $s' \in h'$ and $s = s'$ then $\text{prec}(s, h) = \text{prec}(s', h')$
3. $E : S \times H \times Ags \mapsto 2^S$ is an h -effectivity function yielding for an agent ag the set of next static states allowed by the choice performed by the agent relative to a history. We have the following constraints on h -effectivity functions:
 - a. if $s \notin h$ then $E(s, h, ag) = \emptyset$
 - b. if $s' \in E(s, h, ag)$ then $\exists h' : s' = \text{succ}(s, h')$
 - c. if $s' = \text{succ}(s, h')$ and $s' \in h$ then $s' \in E(s, h, ag)$
 - d. $E(s, h, ag_1) \cap E(s, h', ag_2) \neq \emptyset$ for $ag_1 \neq ag_2$

Condition **2.a** is the ‘backwards-linear’ requirement that we have seen in Section 10.2.2. Condition **3.b** ensures that next state effectivity as seen from a current state s does not contain states s' that are not reachable from the current state through some history. Condition **3.c** expresses the STIT condition of ‘no choice between undivided histories’. Condition **3.d** above states that simultaneous choices of different agents never have an empty intersection. This is the central condition of ‘independence of agency’. It reflects that a choice performance of one agent can never have as a consequence that some other agent is limited in the choices it can exercise simultaneously.

Again, we evaluate truth with respect to dynamic states.

Definition 10.23

Relative to a model $\mathcal{M} = \langle S, H, E, V \rangle$, truth $\langle s, h \rangle \models \varphi$ of a formula φ in a dynamic state $\langle s, h \rangle$, with $s \in h$, is defined as:

$$\begin{array}{ll}
\langle s, h \rangle \models p & \Leftrightarrow s \in V(p) \\
\langle s, h \rangle \models \neg\varphi & \Leftrightarrow \text{not } \langle s, h \rangle \models \varphi \\
\langle s, h \rangle \models \varphi \wedge \psi & \Leftrightarrow \langle s, h \rangle \models \varphi \text{ and } \langle s, h \rangle \models \psi \\
\langle s, h \rangle \models \Box\varphi & \Leftrightarrow \forall h' : \text{if } s \in h' \text{ then } \langle s, h' \rangle \models \varphi \\
\langle s, h \rangle \models X\varphi & \Leftrightarrow \text{if } s' = \text{succ}(s, h) \text{ then } \langle s', h \rangle \models \varphi \\
\langle s, h \rangle \models [ag \text{ xstit}]\varphi & \Leftrightarrow \forall s', h' : \text{if } s' \in E(s, h, ag) \text{ and} \\
& \quad s' \in h' \text{ then } \langle s', h' \rangle \models \varphi \quad \dashv
\end{array}$$

Satisfiability, validity on a frame and general validity are defined as usual.

Now we proceed with the axiomatization.

Theorem 10.4

The following axiom schemes, in combination with a standard axiomatization for propositional logic, and the standard Hilbert rules (like necessitation) for the normal modal operators, define a complete Hilbert system for XSTIT:

$$\begin{array}{ll}
\text{S5 for } \Box & \\
\text{For each } ag, \text{ KD for } [ag \text{ xstit}] & \\
(Lin) \quad \neg X\neg\varphi \leftrightarrow X\varphi & \\
(Sett) \quad \Box X\varphi \rightarrow [ag \text{ xstit}]\varphi & \\
(XSett) \quad [ag \text{ xstit}]\varphi \rightarrow X\Box\varphi & \\
(Indep) \quad \Diamond [ag_1 \text{ xstit}]\varphi \wedge \dots \wedge \Diamond [ag_n \text{ xstit}]\psi \rightarrow & \\
\quad \Diamond ([ag_1 \text{ xstit}]\varphi \wedge \dots \wedge [ag_n \text{ xstit}]\psi) & \\
\quad \text{for } Ags = \{ag_1, \dots, ag_n\} & \dashv
\end{array}$$

Knowingly doing

To study the notion of knowingly doing, like before, in Section 10.4.1, a normal S5 knowledge operator $[ag_i \text{ kno}]\varphi$ is added for every agent ag_i in the system. The equivalence classes, or ‘information sets’ of these operators contain state-history pairs, which means knowledge concerns information about the dynamics of the system of agents. Now knowingly doing is suitably modeled by the combination of operators $[ag_i \text{ kno}][ag_i \text{ xstit}]\varphi$. For the logic of this notion we can consider several interactions between the contributing modalities. Here we only briefly mention some of the possibilities.

It is a fundamental property of agency that an agent cannot know what other agents choose simultaneously. This is expressed by the following axiom.

Definition 10.24

The property of ignorance about concurrent choices of others is defined by the axiom:

$$(IgnCC) \quad [ag_1 \text{ kno}][ag_2 \text{ xstit}]\varphi \rightarrow [ag_1 \text{ kno}]\Box[ag_2 \text{ xstit}]\varphi \text{ for } ag_1 \neq ag_2 \quad \dashv$$

(IgnCC) expresses that if an agent knows that something results from the choice of another agent, it can only be that the agent knows it is settled that that something results from a choice of the other agent.

Definition 10.25

The property of knowledge about the next state is defined by the axiom:

$$(XK) \quad [ag \text{ kno}]X\varphi \rightarrow [ag \text{ kno}][ag \text{ xstit}]\varphi \quad \dashv$$

The (XK) property expresses that the only things an agent can know about the next state are the things it knows to be seeing to it itself.

Definition 10.26

The property of effect recollection is defined by the axiom:

$$(Rec-Eff) \quad [ag \text{ kno}][ag \text{ xstit}]\varphi \rightarrow [ag \text{ xstit}][ag \text{ kno}]\varphi \quad \dashv$$

(Rec-Eff) expresses that if agents knowingly see to something, then they know that something is the case in the resulting state.

The above three properties for knowingly doing just exemplify some of the possibilities. More properties have been studied. Also the theory on these dynamic attitudes has been extended to beliefs and to intentions. The case of intentional action is particularly interesting because there is an extensive philosophical literature on this notion. One of the philosophical scenario's discussed by Broersen for instance is the side effect problem. Here we can only point to the fact that XSTIT, after addition of the right BDI modalities, seems to be a suitable base logic for the study of such notions.

10.5 Conclusion

In this chapter we have reviewed the use of epistemic logic, extended with other modalities for motivational attitudes such as desires and intentions, for describing (the behaviour of) intelligent agents. What is immediately clear is that although all logical approaches are based on and inspired by Bratman's seminal work on BDI theory for practical reasoning, the formalizations themselves are quite different in nature. They also enjoy different (and sometimes even on first sight contradictory) properties. In our view this means that although Bratman did his uttermost to present a clear philosophy, as is often the case when formalizing these kind of philosophical theories, there is still a lot of ambiguity or, put more positively, freedom to formalise these matters. We have even seen that quite different base logics may be used, such as (branching-time) temporal logic, dynamic logic and also stit logic. This makes the formal logics in themselves hard to compare. We think it depends on the purpose of the formalisation (is it used for better understanding, or does it serve as a basis for computational and executable frameworks) which of the BDI logics will be most appropriate. The latter is especially important for designers and programmers of agent systems.

10.6 Notes

Pointers to the theory of Bratman and colleagues include the work on practical reasoning by Bratman (1987) and by Bratman, Israel, and Pollack (1988), as well as the paper on intentions by Bratman (1990). The concept of pro-attitude is also introduced by Bratman (1987). The term ‘Practical Reasoning’ dates back to Aristotle: For further reading and a contemporary reference list the reader is also advised to read the short introduction ‘Practical Reason’ in the Stanford Encyclopedia of Philosophy (2013), or to consult work on ‘Practical Syllogism’. The term ‘screen of admissibility’ to describe Bratman’s theory of how intentions filter admissible options was coined by Cohen and Levesque (1990). For more on the subtle issue of admissible intentions and inconsistent beliefs and an elaborate example, see pp. 41–42 of the paper by Bratman (1987). Daniel Dennett set out his work on the intentional stance (based on the notions of desire and belief) in (Dennett, 1987). A reduction of intention to beliefs and desires alone is rejected (pp. 6–9) by Bratman (1987).

A pioneering paper on intelligent agents is by Wooldridge and Jennings (1995), and also by Wooldridge (1999). The latter for instance addresses the issue of how logics like those studied in this chapter may eventually be *realised* in agent-based systems. The area of agents and multi-agent systems is now a mature and still active area: a good start for further orientation is provided by the website of IFAAMAS (IFAAMAS), which hosts proceedings of the Autonomous Agents and Multi-Agent Systems conference since 2007. As mentioned in the text, Bratman’s semi-formal theory on intentions was further formalised by Cohen and Levesque (1990), and, more or less at the same time, by Rao and Georgeff (1991). In fact their work won, respectively, the 2006 and 2007 IFAAMAS Awards for Influential Papers in Autonomous Agents and Multiagent Systems.

Cohen and Levesque (1990) were among the first to adapt PDL to model agency. Many responses, extensions and variations of this logic have been proposed. For instance, the Side-Effect-Free Principle was coined by Su, Sattar, Lin, and Reynolds (2007), while the study of introspective properties of preferences was undertaken by Herzig and Longin (2004), Lorini and Demolombe (2008), and by Herzig, Lorini, Hübner, and Vercouter (2010). The first of those three references also put forward a definition of epistemic achievement goal. Weakening the link between action and goals, and at the same time allowing for a stronger commitment (a notion criticised by Herzig and Longin (2004)) was undertaken by Sadek (2000) and by Bretier (1995). Sadek (1992) moreover argues for a notion of realistic preference, also called weak realism. The notion of *attempting* to do an action was formalised by Lorini and Herzig (2008).

Our treatment of Rao and Georgeff’s BDI logic is mainly based on (Rao and Georgeff, 1991) (which includes the GOAL-operator) and (Rao and Georgeff, 1998) (where this operator is replaced by an operator DES for desire). Their approach, in turn, is for its formal model heavily inspired by the framework on branching time logic put forward by Emerson (1990). The abbreviation CTL stands for Computation-Tree Logic, in which one can quantify over branches, and, given a branch, over its points. CTL* is an extension of CTL where some of CTL’s restrictions on the occurrences of branching and linear time quantifiers, are lifted. The

semantic constraints that make some of the desirable axioms of Section 10.2.2 valid, are discussed by Rao and Georgeff (1991, 1998) and by Wooldridge (2000). In particular, the world/time point compatibility requirement is taken from Wooldridge (2000). We have used the little Nell problem as an advocate to use branching time logic, rather than linear time logic. This problem was discussed by McDermott (1982), and also addressed by the papers of Cohen & Levesque and Rao & Georgeff. It is interesting to note that in a series of talks and papers however, Moshe Vardi has argue that from a computer science perspective, it is far from clear that branching time provides a superior model over that of linear time (see for instance the paper by Vardi (2001)).

KARO stands for Knowledge, Ability, Result and Opportunity. The KARO framework was developed in a number of papers by van Linder, van der Hoek, and Meyer (1995, 1997), van der Hoek, van Linder, and Meyer (1998) and Meyer, van der Hoek, and van Linder (1999) and the thesis by van Linder (1996). All the basic operators of KARO can be interpreted as modal operators, cf. also footnote 5, which was proven by van der Hoek, Meyer, and van Schagen (2000). Adding a notion of agenda to KARO, in order to deal with commitments, was proposed and studied by Meyer et al. (1999) (see also footnote 6). Readers interested in knowing the details of our simplification as mentioned in footnote 8 are referred to the paper by Meyer et al. (1999). One of KARO's foundations is dynamic logic, for which the book by Harel (1984) is a standard reference. An early reference to the agent programming language GOAL (Goal-Oriented Agent Language) is provided by Hindriks, Boer, Hoek, and Meyer (2001), while de Boer, Hindriks, van der Hoek, and Meyer (2007) provide, rather than the language for declarative goals itself, both a programming framework and a programming logic for such computational agents. The connection with intention logic of Cohen & Levesque was given by Hindriks, van der Hoek, and Meyer (2012). A starting point to get familiar with the software platform for GOAL is (2013). As an example of how logic frameworks can be beneficial for designers and programmers of agent systems, the KARO framework played an important role in devising the agent programming language 3APL by Hindriks, de Boer, van der Hoek, and Meyer (1999). Moreover, by formalising the relations between the KARO framework and the 3APL language, Hindriks and Meyer (2007) and Meyer (2007) were able to propose a verification logic for 3APL and its derivatives. As a proper treatment of this aspect of BDI logics goes beyond the scope of this chapter we refer to the literature on this issue, e.g., the work by van der Hoek and Wooldridge (2003), Dastani, van Riemsdijk, and Meyer (2007), Alechina, Dastani, Logan, and Meyer (2007), Dastani, Hindriks, and Meyer (2010), and by Hindriks et al. (2012).

Adding emotions to the mix of agents' attitudes is by now a well-respected avenue in agent-based modeling. Emotions are now regarded to complement rationality, rather than being in tension with it. This was acknowledged in for instance cognitive science through the work of Damasio (1994), and is comparable with the rise of 'behavioural economics' complementing 'classical game theory' (for the latter, see also Chapter 9 of this book). See for instance 'Emotions and Economic Theory' by Elster (1998). The work that we presented on emotions was inspired by two psychological theories: that of Oatley and Jenkins (1996) and that of Ortony, Clore, and Collins (1988) (also referred to as OCC). More specifi-

cally, KARO logic for emotional agents as based on Oatley & Jenkins' framework is worked out by Meyer (2006), while the formalisation of OCC theories is done by Steunebrink, Dastani, and Meyer (2007, 2012). References to the work of Adam *et al.* on emotions include the papers by Adam (2007) and by Adam, Herzig, and Longin (2009). Their introspective properties for emotion presented at the end of Section 10.3 are taken from Theorem 13 of the latter reference.

STIT (Seeing To It That) logics are an attempt of treating agency as a modality in an approach originating in a series of papers, beginning with the paper by Belnap and Perloff (1988). Other important contributions in this field are by Chellas (1995) and by Horty (2001). A semantics for STIT based on worlds (rather than state-history pairs) was given by Balbiani, Herzig, and Troquard (2007). The work of Semmling and Wansing (2008) adds BDI modalities to STIT logic. Lorini and Schwarzenruber (2011) use STIT logic to analyse emotions. For further details and motivation for next-state STIT logic XSTIT we refer to work by Broersen (2009, 2011a). The latter, together with Broersen (2011b) (which, among other things, discusses the side effect problem), is the place to further read on knowingly doing. Monotonic modal logic is presented in Hansen (2003). Theorem 10.2 is taken from Xu (1998), and Theorem 10.4 from Broersen (2011a).

There have been a number of further developments, which can be seen as consequences of BDI logics. First of all we mention that on the basis of a BDI logic, Shoham (1993) initiated the field of agent-oriented programming (AOP), which is basically a way of programming the (BDI-like) mental states of agents as a new paradigm of programming. It could be viewed as a successor (or even refinement) of the well-known Object-Oriented (OO) Programming paradigm. Momentarily there is a host of such dedicated agent cognitive (BDI) programming languages (see the work by Bordini, Dastani, Dix, and Seghrouchni (2005, 2009)), including 2APL//3APL by Hindriks *et al.* (1999) and Dastani (2008) which is inspired by KARO.

Another extension that is made in the literature is a combination with reasoning with uncertainty. Casali, I. Godo, and Sierra (2011) propose a so-called Graded BDI logic. This means that beliefs, desires and intentions are not crisp anymore, but fuzzy. This is done by using a multi-valued base logic. In particular, a fuzzy modal operator D^+ is introduced with as reading “ φ is positively desired”. Its truth degree corresponds to the agent's level of satisfaction would φ become true. It satisfies the principle $D^+(\varphi \vee \psi) \leftrightarrow (D^+\varphi \wedge D^+\psi)$. There is also an analogous operator D^- , read “ φ is negatively desired”. Intention is then defined from this. Their logic combines Rational Pavelka Logic (an extension of Lukasiewicz's many-valued logic with infinitely many truth values by rational truth degree constants) as described by Hájek (1998) with a multi-context logic of Giunchiglia and Serafini (1994) in order to account for the agents' mental attitudes. The approach allows for a much more fine-grained deliberation process. In particular, the process of deriving intentions can now be formally expressed as follows: the intention degree to reach a desire φ by means of a plan α is taken as a trade-off between the benefit of reaching this desire and the cost of the plan, weighted by the belief degree r .

Lorini and Demolombe (2008) and Lorini (2011) propose a similar graded approach to BDI, but now based on a sense of plausibility (non-exceptionality) as proposed by Spohn (1998). In this logic both the belief and the goal attitude are

graded. Beyond modal operators of knowledge Know_i , the logic of Lorini (2011) has special atoms exc_h indicating that the degree of exceptionality of the current state is h , where h ranges over a finite set of natural numbers Num . Then graded belief is defined as

$$\text{Bel}_i^{\geq h} \varphi \stackrel{\text{def}}{=} \bigvee_{k \geq h} \left(\neg \text{Know}_i(\text{exc}_k \rightarrow \varphi) \wedge \bigwedge_{l < k} \text{Know}_i(\text{exc}_l \rightarrow \varphi) \right)$$

Graded goals $\text{Goal}_i^{\geq h} \varphi$ are defined in a similar way from special atoms des_h indicating that the degree of desirability of the current state is h . This theory is next applied to the modelling of expectation-based emotions such as hope, fear, disappointment and relief, and their intensity.

Finally, in analogy to the area of belief revision, also a lot of work has appeared on the revision of the other mental states in BDI, notably intentions. In fact, intention reconsideration is already part of the first implementation of Bratman's theory as set out by Bratman et al. (1988) and is present in the BDI architecture of Rao and Georgeff (1992) via the so-called BDI control loop. It also is present in the logical theory KARO of van Linder et al. (1995) and Meyer et al. (1999) that we have seen earlier. Using a dynamic logic it is very natural to also incorporate actions in KARO that are mental state-revising such as belief revision as proposed by van Linder et al. (1995) and changes in commitments as studied by Meyer et al. (1999). But recently there have appeared more fundamental theories of intention revision, such as by van der Hoek, Jamroga, and Wooldridge (2007), where a logic of intention dynamics is developed based on dynamic logic and a dynamic update operator of the form $[\Omega]\varphi$, meaning after the agent has updated on the basis of observations Ω , it must be the case that φ . The approach of beliefs, desires and intentions is rather syntactical (set of sentences without modalities for these mental attitudes interpreted by model structures). Contrary to this, van Ditmarsch, de Lima, and Lorini (2010) present a model-theoretic approach of intentions and intention dynamics. It uses modal operators for time, belief and choice as basis, and intention as a derived operator, and for intention change it employs a dynamic operator called local assignment. This is an operation on the model that changes the truth value of atomic formulae at specific time points. Shoham (2009) discusses the interaction of belief revision and intention revision. This is done more in a traditional AGM sense as put forward by Alchourrón, Gärdenfors, and Makinson (1985), than in a dynamic epistemic modal sense.

Acknowledgement

Thanks to Emiliano Lorini for his useful comments on an earlier version of this chapter.

References

- Adam, C. (2007). *Emotions: from psychological theories to logical formalization and implementation in a BDI agent*. Ph. D. thesis, Institut National Polytechnique de Toulouse, Toulouse.
- Adam, C., A. Herzig, and D. Longin (2009). A logical formalization of the OCC theory of emotion. *Synthese* 168(2), 201–248.
- Alchourrón, C. E., P. Gärdenfors, and D. Makinson (1985). On the logic of theory change: Partial meet contractions and revision functions. *Journal of Symbolic Logic* 50(2), 510–530.
- Alechina, N., M. Dastani, B. Logan, and J.-J. Ch. Meyer (2007). A logic of agent programs. In R. C. Holte and A. E. Howe (Eds.), *Proceedings of AAAI-07*, Vancouver, Canada, pp. 795–800. AAAI Press.
- Balbiani, P., A. Herzig, and N. Troquard (2007). Alternative axiomatics and complexity of deliberative stit theories. *Journal of Philosophical Logic* 37, 387–406.
- Belnap, N. and M. Perloff (1988). Seeing to it that: a canonical form for gentiles. *Theoria* 54, 175–199.
- Bordini, R. H., M. Dastani, J. Dix, and A. E. F. Seghrouchni (Eds.) (2005). *Multi-Agent Programming: Languages, Platforms and Applications*, Volume 15 of *Multiagent Systems, Artificial Societies, and Simulated Organizations*. New York: Springer.
- Bordini, R. H., M. Dastani, J. Dix, and A. E. F. Seghrouchni (Eds.) (2009). *Multi-Agent Programming (Languages, Tools and Applications)*. Dordrecht Heidelberg: Springer.
- Bratman, M. E. (1987). *Intentions, Plans, and Practical Reason*. Massachusetts: Harvard University Press.
- Bratman, M. E. (1990). What is intention? In P. R. Cohen, J. Morgan, and M. E. Pollack (Eds.), *Intentions in Communication*, Chapter 2, pp. 15–31. Cambridge, Massachusetts: MIT Press.
- Bratman, M. E., D. J. Israel, and M. E. Pollack (1988). Plans and resource-bounded practical reasoning. *Computational Intelligence* 4, 349–355.
- Bretier, P. (1995). *La communication orale coopérative: contribution à la modélisation logique et à la mise en oeuvre d'un agent rationnel dialoguant*. Ph. D. thesis, Université Paris Nord, Paris, France.
- Broersen, J. M. (2009). A complete stit logic for knowledge and action, and some of its application. In M. Baldoni, T. C. Son, M. B. van Riemsdijk, and M. Winikoff (Eds.), *Declarative Agent Languages and Technologies VI (DALT 2008)*, Volume 5397 of *LNCS*, Berlin, pp. 47–59. Springer.

- Broersen, J. M. (2011a). Deontic epistemic *stit* logic distinguishing modes of ‘Mens Rea’. *Journal of Applied Logic* 9(2), 127–152.
- Broersen, J. M. (2011b). Making a start with the *stit* logic analysis of intentional action. *Journal of Philosophical Logic* 40, 399–420.
- Casali, A., L. I. Godo, and C. Sierra (2011). A graded BDI agent model to represent and reason about preferences. *Artificial Intelligence* 175(7-8), 1468–1478.
- Chellas, B. F. (1995). On bringing it about. *Journal of Philosophical Logic* 24, 563–571.
- Cohen, P. R. and H. J. Levesque (1990). Intention is choice with commitment. *Artificial Intelligence* 42(3), 213–261.
- Damasio, A. (1994). *Descartes’ Error: Emotion, Reason, and the Human Brain*. New York: Grosset/Putnam Press.
- Dastani, M. (2008). 2APL: a practical agent programming language. *Autonomous Agents and Multi-Agent Systems* 16(3), 214–248.
- Dastani, M., K. V. Hindriks, and J.-J. Ch. Meyer (Eds.) (2010). *Specification and Verification of Multi-Agent Systems*. New York/Dordrecht/Heidelberg/London: Springer.
- Dastani, M., B. van Riemsdijk, and J.-J. Ch. Meyer (2007). A grounded specification language for agent programs. In M. Huhns, O. Shehory, E. H. Durfee, and M. Yokoo (Eds.), *Proceedings of the 6th International Joint Conference On Autonomous Agents and Multi-Agent Systems (AAMAS2007)*, pp. 578–585.
- de Boer, F., K. Hindriks, W. van der Hoek, and J.-J. Meyer (2007). A verification framework for agent programming with declarative goals. *Journal of Applied Logic* 5(2), 277 – 302.
- Dennett, D. C. (1987). *The Intentional Stance*. Cambridge, Massachusetts: MIT Press.
- van Ditmarsch, H., T. de Lima, and E. Lorini (2010). Intention change via local assignments. In *Proceedings of LADS 2010*, Volume 6822 of *LNCS*, Berlin/Heidelberg, pp. 136–151. Springer-Verlag.
- Elster, J. (1998). Emotions and economic theory. *Journal of Economic Literature* 36(1), 47–74.
- Emerson, E. A. (1990). Temporal and modal logic. In J. van Leeuwen (Ed.), *Handbook of Theoretical Computer Science*, Volume B: Formal Models and Semantics, Chapter 14, pp. 996–1072. Amsterdam: Elsevier Science.
- Giunchiglia, F. and L. Serafini (1994). Multilanguage hierarchical logics (or: how we can do without modal logics). *Artificial Intelligence* 65, 29–70.

- Goal (2013). The software platform for GOAL. <http://ii.tudelft.nl/trac/goal>, retrieved July 2013.
- Hájek, P. (1998). *Metamathematics of Fuzzy Logic*, Volume 4 of *Trends in Logic*. Dordrecht: Kluwer Academic Publishers.
- Hansen, H. H. (2003). Monotonic modal logics. Master's thesis, ILLC, Amsterdam.
- Harel, D. (1984). Dynamic logic. In D. Gabbay and F. Guenther (Eds.), *Handbook of Philosophical Logic*, Volume II, pp. 497–604. Dordrecht/Boston: Reidel.
- Herzig, A. and D. Longin (2004). C&L intention revisited. In D. Dubois, C. Welty, and M.-A. Williams (Eds.), *Proceedings of the 9th International Conference on Principles of Knowledge Representation and Reasoning (KR2004)*, pp. 527–535. AAAI Press.
- Herzig, A., E. Lorini, J. F. Hübner, and L. Vercouter (2010). A logic of trust and reputation. *Logic Journal of the IGPL* 18(1), 214–244.
- Hindriks, K. V., F. S. d. Boer, W. v. d. Hoek, and J.-J. Ch. Meyer (2001). Agent programming with declarative goals. In *Proceedings of the 7th International Workshop on Intelligent Agents VII. Agent Theories Architectures and Languages*, ATAL '00, London, UK, UK, pp. 228–243. Springer-Verlag.
- Hindriks, K. V., F. S. de Boer, W. van der Hoek, and J.-J. Ch. Meyer (1999). Agent programming in 3APL. *International Journal of Autonomous Agents and Multi-Agent Systems* 2(4), 357–401.
- Hindriks, K. V. and J.-J. Ch. Meyer (2007). Agent logics as program logics: grounding KARO. In C. Freksa, M. Kohlhase, and K. Schill (Eds.), *29th Annual German Conference on AI, KI 2006*, Volume 4314 of *LNAI*, pp. 404–418. Springer.
- Hindriks, K. V., W. van der Hoek, and J.-J. Ch. Meyer (2012). GOAL agents instantiate intention logic. In A. Artikis, R. Craven, N. K. Çiçekli, B. Sadighi, and K. Stathis (Eds.), *Logic Programs, Norms and Action (Sergot Festschrift)*, Volume 7360 of *LNAI*, Heidelberg, pp. 196–219. Springer.
- van der Hoek, W., W. Jamroga, and M. Wooldridge (2007). Towards a theory of intention revision. *Synthese* 155(2), 265–290.
- van der Hoek, W., J.-J. Ch. Meyer, and J. W. van Schagen (2000). Formalizing potential of agents: the KARO framework revisited,. In M. Faller, S. Kaufmann, and M. Pauly (Eds.), *Formalizing the Dynamics of Information*, Volume 91 of *CSLI Lecture Notes*, pp. 51–67. Stanford: CSLI Publications.
- van der Hoek, W., B. van Linder, and J.-J. Ch. Meyer (1998). An integrated modal approach to rational agents. In M. Wooldridge and A. Rao (Eds.), *Foundations of Rational Agency*, Volume 14 of *Applied Logic Series*, pp. 133–168. Dordrecht: Kluwer.

- van der Hoek, W. and M. Wooldridge (2003). Towards a logic of rational agency. *Logic Journal of the IGPL* 11(2), 133–157.
- Horty, J. F. (2001). *Agency and Deontic Logic*. Oxford University Press.
- IFAAMAS. <http://www.ifaamas.org>, retrieved July 2013.
- van Linder, B. (1996). *Modal Logics for Rational Agents*. Ph. D. thesis, Utrecht University.
- van Linder, B., W. van der Hoek, and J.-J. Ch. Meyer (1995). Actions that make you change your mind: belief revision in an agent-oriented setting. In A. Laux and H. Wansing (Eds.), *Knowledge and Belief in Philosophy and Artificial Intelligence*, pp. 103–146. Berlin: Akademie Verlag.
- van Linder, B., W. van der Hoek, and J.-J. Ch. Meyer (1997). Seeing is believing (and so are hearing and jumping). *Journal of Logic, Language and Information* 6, 33–61.
- Lorini, E. (2011). A dynamic logic of knowledge, graded beliefs and graded goals and its application to emotion modelling. In *Proceedings of LORI 2011*, Volume 6953 of *LNCS*, Berlin / Heidelberg, pp. 165–178. Springer-Verlag.
- Lorini, E. and R. Demolombe (2008). Trust and norms in the context of computer security: toward a logical formalization. In R. V. der Meyden and L. V. der Torre (Eds.), *Proceedings of the International Workshop on Deontic Logic in Computer Science (DEON 2008)*, Volume 5076 of *LNCS*, Berlin/Heidelberg, pp. 50–64. Springer-Verlag.
- Lorini, E. and A. Herzig (2008). A logic of intention and attempt. *Synthese KRA* 163(1), 45–77.
- Lorini, E. and F. Schwarzentruher (2011). A logic for reasoning about counterfactual emotions. *Artificial Intelligence* 175, 814–847.
- McDermott, D. V. (1982). A temporal logic for reasoning about processes and plans. *Cognitive Science* 6, 101–155.
- Meyer, J.-J. Ch. (2006). Reasoning about emotional agents. *International Journal of Intelligent Systems* 21(6), 601–619.
- Meyer, J.-J. Ch. (2007). Our quest for the holy grail of agent verification. In N. Olivetti (Ed.), *Proceedings of TABLEAUX 2007*, Volume 4548 of *LNAI*, Berlin/Heidelberg, pp. 2–9. Springer.
- Meyer, J.-J. Ch., W. van der Hoek, and B. van Linder (1999). A logical approach to the dynamics of commitments. *Artificial Intelligence* 113, 1–40.
- Oatley, K. and J. M. Jenkins (1996). *Understanding Emotions*. Blackwell Publishing.

- Ortony, A., G. L. Clore, and A. Collins (1988). *The Cognitive Structure of Emotions*. Cambridge: Cambridge University Press.
- Rao, A. and M. P. Georgeff (1992). An abstract architecture for rational agents. In B. Nebel, C. Rich, and W. Swartout (Eds.), *Proceedings of the Third International Conference on Principles of Knowledge Representation and Reasoning (KR'92)*, San Mateo, California, pp. 439–449. Morgan Kaufmann.
- Rao, A. S. and M. P. Georgeff (1991). Modeling rational agents within a BDI-architecture. In R. F. J. Allen and E. Sandewall (Eds.), *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning (KR'91)*, San Mateo, California, pp. 473–484. Morgan Kaufmann.
- Rao, A. S. and M. P. Georgeff (1998). Decision procedures for BDI logics. *Journal of Logic and Computation* 8(3), 293–344.
- Sadek, M. D. (1992). A study in the logic of intention. In B. Nebel, C. Rich, and W. Swartout (Eds.), *Proceedings of the Third International Conference on Principles of Knowledge Representation and Reasoning (KR'92)*, San Mateo, California, pp. 462–473. Morgan Kaufmann.
- Sadek, M. D. (2000). Dialogue acts are rational plans. In M. Taylor, F. Nel, and D. Bouwhuis (Eds.), *The structure of multimodal dialogue*, Philadelphia/Amsterdam, pp. 167–188. From ESCA/ETRW, Workshop on The Structure of Multimodal Dialogue (Venaco II), 1991.
- Semmling, C. and H. Wansing (2008). From BDI and *stit* to *bdi-stit* logic. *Logic and Logical Philosophy* 17(1-2), 185–207.
- Shoham, Y. (1993). Agent-oriented programming. *Artificial Intelligence* 60(1), 51–92.
- Shoham, Y. (2009). Logical theories of intention and the database perspective. *Journal of Philosophical Logic* 38(6), 633–647.
- Spohn, W. (1998). Ordinal conditional functions: a dynamic theory of epistemic states. In *Causation in Decision, Belief Change and Statistics*, pp. 105–134. Dordrecht: Kluwer.
- Steunebrink, B., M. Dastani, and J.-J. Ch. Meyer (2007). A logic of emotions for intelligent agents. In R. C. Holte and A. E. Howe (Eds.), *Proceedings of AAAI-07*, Vancouver, Canada, pp. 142–147. AAAI Press.
- Steunebrink, B. R., M. Dastani, and J.-J. Ch. Meyer (2012). A formal model of emotion triggers for BDI agents with achievement goals. *Synthese/KRA* 185(1), 83–129.
- Su, K., A. Sattar, H. Lin, and M. Reynolds (2007). A modal logic for beliefs and pro attitudes. In *Proceedings of the 22nd AAAI Conference on Artificial Intelligence (AAAI 2007)*, pp. 496–501.

- Vardi, M. Y. (2001). Branching vs. linear time: Final showdown. In *Tools and Algorithms for the Construction and Analysis of Systems*, pp. 1–22. Springer.
- Wallace, R. J. (2013). http://en.wikipedia.org/wiki/Practical_reason, retrieved July 2013.
- Wooldridge, M. J. (1999). Intelligent agents. In G. Weis (Ed.), *Multiagent Systems*, pp. 27–77. Cambridge, Massachusetts: The MIT Press.
- Wooldridge, M. J. (2000). *Reasoning about Rational Agents*. Cambridge, Massachusetts: MIT Press.
- Wooldridge, M. J. and N. R. Jennings (Eds.) (1995). *Intelligent Agents*. Berlin: Springer.
- Xu, M. (1998). Axioms for deliberative stit. *Journal of Philosophical Logic* 27(5), 505–552.