# A Logic of Capabilities

(extended abstract)

W. van der Hoek
B. van Linder
J.-J.Ch. Meyer*
Vrije Universiteit Amsterdam
Faculty of Mathematics and Informatics
De Boelelaan 1081a, 1081 HV  Amsterdam
The Netherlands
E-mail: wiebe@cs.vu.nl, bernd@cs.vu.nl, jules@cs.vu.nl

January 26, 1993

## Abstract

In this paper a logic of capabilities is introduced. We start by defining a language in which not only knowledge and actions of agents can be expressed, but also the abilities. The semantics for this language is given by extended Kripke models. Equivalence transformations on actions are defined in such a way that equivalent actions have equivalent results and that the abilities of an agent are closed under these transformations. The planning capacities of the agent are considered by looking at the Can-predicate and the Cannot-predicate. The former indicates that an agent knows that he can achieve a given goal by performing a given action, the latter indicates that the agent knows that that he cannot achieve the given goal by performing the given action. It turns out that the possibility to express the abilities of an agent provides for a flexible and intuitively appealing framework.

## 1. Introduction

The planning of an agent in order to achieve some given goal is a topic of much interest in artificial intelligence. An overview of the various aspects of planning can be found in [GL87]. According to [Pol92], in AI planning is considered to be the process of formulating a program of action to achieve some specified goal. Obviously, it is only possible for an agent to achieve a given goal by performing a given sequence of actions if performing this sequence leads to the goal and if it is possible for the agent

---
* also at the Katholieke Universiteit Nijmegen.

to perform all the actions in the sequence in the correct order.

The logic that we consider in this paper is based on the concepts of the logic defined in [Moo84]. We follow the notation of [MH], which is based on Harel's dynamic logic (see [Har84]). In this paper we will restrict ourselves to a propositional logic. This is basically done for the sake of simplicity. Nevertheless, in principle an extension to first-order logic is certainly possible and might even be necessary for practical purposes.

We introduce a syntactical operator that can be used to indicate whether a given agent is *capable* of performing a given action. The notion of ability to perform actions bears a slight resemblance to the notion of knowledge of formulae in awareness logic as introduced in [FH88]. In awareness logic an agent needs to be aware of a given formula in order to gain explicit belief of this formula whereas in the logic of capabilities an agent needs to be able to perform a given action to make it possible for him to perform that action. Using capability and knowledge of capability in combination with knowledge of the result of an action provides for some rather intuitive results for the planning capacities of an agent, as can be seen in section 6. With regard to planning, only correctness and feasibility of plans is considered, i.e. our approach does not deal with finding heuristics or any other aspects of planning.

The contents of this paper is as follows. In section 2 the language $\mathcal{L}$ is defined. In section 3.1 a class of models $\mathfrak{M}$, and the semantics for the language $\mathcal{L}$ are defined. In section 3.2, the validity of some formulae from $\mathcal{L}$ for $\mathfrak{M}$ is proved. Furthermore some other properties are shown. In section 3.3 the notion of action transformation is introduced and it is shown that the abilities of an agent are

closed under the action transformations given here. In section 4 the axioms and proof rules of the logic of capabilities are given. In section 5 the planning capacities of an agent are considered by giving an alternative definition of the so called Can-predicate, originally introduced in [Moo84]. We furthermore define a Cannot-predicate, which also proves to be useful when considering the planning capacities of a rational agent. In section 6 we sum up some properties of both the Can-predicate and the Cannot-predicate. These properties show that both predicates behave in an intuitively acceptable way. Section 7 ends this paper with a discussion and some guidelines for further research.

## 2. Defining the language $\mathcal{L}$

As stated in section 1, the system defined in this paper is loosely based on the system defined in [Moo84]. The main difference lies in the fact that in Moore's system no attention whatsoever is paid to the physical and mental condition of an agent. It is assumed that whenever an agent knows what some action $\alpha$ is, he is automatically able to perform $\alpha$, regardless of whether the agent satisfies all physical or mental demands that performance of action $\alpha$ asks from him. In our opinion this approach does not seem to be completely in correspondence with the intuition as can be seen in the following example.

2.1. EXAMPLE. Assume that some agent is planning to win a gold medal in the men's 100 metres at the olympics. To achieve this the agent has to run faster than any other competitor. The agent knows that if he runs fastest he will win the gold medal. He also knows what the action 'run fastest' implies. According to the theory of [Moo84], this suffices to conclude that the agent can win the gold medal. Intuitively this seems to be rather strange: it seems to be necessary that the agent is *physically able* to run fastest; it does not suffice that he knows what this implies.
On the other hand, if the agent knows that if he runs fastest he will win the gold medal and he is able to run fastest, it is rather natural to conclude that the agent can make a successful plan to get a gold medal.

The observation made in example 2.1 clearly shows the need for the introduction of a new operator that can be used to express the ability of an agent to perform a given action. For this the syntactical operator $\mathbf{A}$, indexed with an agent $i$, is used in $\mathcal{L}$; the expression $\mathbf{A}_i\alpha$ is used to denote the fact that the agent $i$ is able to perform the action $\alpha$.

2.2. CONVENTION. Throughout this paper we will speak of *actions* to refer to executable codes, programs or algorithms. A typical action is for instance 'open the window'. We will speak of *events* to refer to the execution of some specific action by a specific agent. For instance, the opening of the window by rational agent $i$ is a typical event.

2.3. DEFINITION. The language $\mathcal{L}$ is based on a given set $\Pi$ of propositional symbols, a finite set $At$ of atomic actions and a finite set $\mathcal{A} = \{1, \ldots, n\}$ of agents. The language $\mathcal{L}$ is defined by the following BNF, where $\varphi$ is a typical element of the set of formulae $\mathcal{L}$, $p$ is a typical element of the set of propositional symbols $\Pi$, $\alpha$ is a typical element of the set of actions $Ac$, $a$ is a typical element of the set of atomic actions $At$ and $i$ is a typical element of the set of agents $\mathcal{A}$.

$$
\begin{aligned}
\varphi ::= \ &p & | \\
&\neg\varphi & | \\
&\varphi_1 \wedge \varphi_2 & | \\
&\mathbf{K}_i\varphi & | \\
&[\mathrm{do}_i(\alpha)]\varphi & | \\
&\mathbf{A}_i\alpha &
\end{aligned}
$$

where the class $Ac$ of actions is given by:

$$
\begin{aligned}
\alpha ::= \ &a & | \\
&\texttt{skip} & | \\
&\texttt{fail} & | \\
&\texttt{test } \varphi & | \\
&\alpha_1; \alpha_2 & | \\
&\texttt{if } \varphi \texttt{ then } \alpha_1 \texttt{ else } \alpha_2 \texttt{ fi} & | \\
&\texttt{while } \varphi \texttt{ do } \alpha_1 \texttt{ od} & | \\
&\alpha_1 + \alpha_2 &
\end{aligned}
$$

To provide for optimal flexibility, no demand is made whatsoever upon the use of parentheses in $\mathcal{L}$: whenever one thinks that using parentheses provides for more clarity one is encouraged to use them.

We let $p$ with possible marks denote atomic propositions; $\neg\varphi$ denotes the negation of a formula and $\varphi_1 \wedge \varphi_2$ denotes the conjunction of formulae. The connectives $\vee, \rightarrow$ and $\leftrightarrow$ are defined in the usual way. The formula tt is introduduced as an abbreviation for $p \vee \neg p$ and ff is assumed to denote $\neg$tt. The formula $<\mathrm{do}_i(\alpha)> \varphi$ is defined as the dual of $[\mathrm{do}_i(\alpha)]\varphi$: $<\mathrm{do}_i(\alpha)> \varphi \equiv \neg[\mathrm{do}_i(\alpha)]\neg\varphi$.
The formula $\mathbf{K}_i\varphi$ denotes agent $i$'s knowledge of $\varphi$. Following common practice in the world of computer science –see for instance [HM85], [FH88] or [MH]– the modal system $\mathbf{S5}$ is used for the knowledge part of LCap. The axioms for $\mathbf{S5}$ are given in section 4.

The meaning of the formula $[\mathrm{do}_i(\alpha)]\varphi$ is as follows: if the agent $i$ performs the action $\alpha$ then the formula $\varphi$ will be

true after the performance of $\alpha$ by $i$.

The formula $\mathbf{A}_i\alpha$ denotes the ability of agent $i$ to perform the action $\alpha$. This ability of the agent can be thought of as a combination of physical and mental abilities.

The class of actions $Ac$ can be seen as a *rich-test* extension (see [Har84] or [Har79]) of the class of actions considered in [Hoa69]. The meanings of the respective actions in $Ac$ are subsequently: the atomic action, the empty action, the never succeeding action, a test for the truth of formula $\varphi$, sequential composition, conditional composition, repetitive composition and the nondeterministic choice. For a given finite set of agents $\mathcal{A}$ and a given class of actions $Ac$, the class of events $\mathrm{do}_i(\alpha)$ with $i \in \mathcal{A}$ and $\alpha \in Ac$ is usually denoted by $\mathcal{E}_{\mathcal{A}}^{Ac}$.

The expressive power of language $\mathcal{L}$ is considerable, as can be seen in the following example.

2.4. EXAMPLE (The egocentric agent). Consider the following action: 'If agent $i$ knows that he will feel better after helping his neighbour and he knows that he is able to help him, he will do so, otherwise he will do nothing'. It is possible to denote in the language $\mathcal{L}$ the fact that if agent $i$ is feeling well before performing the egocentric action, he will still feel well or he will feel even better after performing it. Using some obvious abbreviations this formalization is given by:

$$w \to [\mathrm{do}_i(\text{if } \mathbf{K}_i[\mathrm{do}_i(h)]b \land \mathbf{K}_i\mathbf{A}_ih$$
$$\text{then } h \text{ else skip fi})](w \lor b)$$

# 3. Semantics for the language $\mathcal{L}$

In this section a class of models $\mathfrak{M}$ is defined. The models from $\mathfrak{M}$ are used to provide a semantics for the language $\mathcal{L}$. This semantics is an extension of the possible world semantics commonly used to model modal systems. Furthermore the validity of some formulae from $\mathcal{L}$ is proved for the class $\mathfrak{M}$ of models. At the end of this section the notion of action transformation is introduced and equivalences under these transformations are shown.

## 3.1. Basic definitions

In the following definitions it is assumed that some set $\mathbf{bool} = \{1, 0\}$ of truth values is given. It will always be clear from the context whether 0 and 1 denote truth values or natural numbers.

3.1. DEFINITION.
A Kripke model $\mathcal{M}$ is a tuple $< \mathcal{S}, \pi, \mathrm{R}, \mathrm{r}, \mathrm{c} >$ where

(1) $\mathcal{S}$ is a set of possible worlds.

(2) $\pi : \Pi \times \mathcal{S} \to \mathbf{bool}$ is a total function that assigns a truth value to propositional symbols in possible worlds.

(3) $\mathrm{R} : \mathcal{A} \to \wp(\mathcal{S} \times \mathcal{S})$ is the function that yields the accessibility relations for a given agent, i.e. if $(s, s') \in \mathrm{R}(i)$ this means that $s'$ is an epistemic alternative for the agent $i$ in possible world $s$. Since the knowledge part of the logic of capabilities is given by the modal system S5, the function $\mathrm{R}$ is such that for all agents $i$ and for all $s$ and $s'$ from $\mathcal{S}$

- $(s, s) \in \mathrm{R}(i)$
- $(s, s') \in \mathrm{R}(i) \Rightarrow (s', s) \in \mathrm{R}(i)$
- $(s, s') \in \mathrm{R}(i), (s', s'') \in \mathrm{R}(i) \Rightarrow (s, s'') \in \mathrm{R}(i)$

all hold, i.e. $\mathrm{R}(i)$ is an *equivalence relation* (cf. [HM85], [MH]).

(4) $\mathrm{r} : \mathcal{A} \times At \to \mathcal{S} \to \wp(\mathcal{S})$ is such that $\mathrm{r}(i, a)(s)$ yields the result of performing atomic action $a$ by agent $i$ in the possible world $s$. This function is such that for all atomic actions $a$ holds that

$$\forall i \forall a \forall s[|\mathrm{r}(i, a)(s)| \leq 1]$$

where $|V|$ denotes the number of elements of the set $V$.

(5) $\mathrm{c} : \mathcal{A} \times At \to \mathcal{S} \to \mathbf{bool}$ is the capability function such that $\mathrm{c}(i, a)(s)$ indicates whether the agent $i$ is capable of performing the action $a$ in the possible world $s$.

The class of models given by definition 3.1 is denoted by $\mathfrak{M}$.

3.2. DEFINITION (Semantics of the logic of capabilities). Let $\mathcal{M} = < \mathcal{S}, \pi, \mathrm{R}, \mathrm{r}, \mathrm{c} >$ be some Kripke model from $\mathfrak{M}$. The functions $\pi$, $\mathrm{r}$ and $\mathrm{c}$ are extended by simultanous induction as follows.

| | |
|---|---|
| $\pi$ | $: \mathcal{L} \times \mathcal{S} \to \mathbf{bool}$ |
| $\pi(\neg\varphi, s)$ | $= \neg\pi(\varphi, s)$ |
| $\pi(\varphi \land \psi, s)$ | $= \pi(\varphi, s) \land \pi(\psi, s)$ |
| $\pi(\mathbf{K}_i\varphi, s)$ | $= \bigwedge \pi(\varphi, s')$ for all $s'$ with $(s, s') \in \mathrm{R}(i)$. |
| $\pi([\mathrm{do}_i(\alpha)]\varphi, s)$ | $= \bigwedge \pi(\varphi, s')$ for all $s'$ with $s' \in \mathrm{r}(i, \alpha)(s)$. |
| $\pi(\mathbf{A}_i\alpha, s)$ | $= \mathrm{c}(i, \alpha)(s)$ |
| | |
| $\mathrm{r}$ | $: \mathcal{A} \times Ac \to \mathcal{S} \to \wp(\mathcal{S})$ |
| $\mathrm{r}(i, \mathtt{skip})(s)$ | $= \{s\}$ |
| $\mathrm{r}(i, \mathtt{fail})(s)$ | $= \emptyset$ |
| $\mathrm{r}(i, \mathtt{test } \varphi)(s)$ | $= \{s\}$ if $\pi(\varphi, s) = 1$ |

$$r(i, \alpha_1; \alpha_2)(s) \qquad = \emptyset \text{ otherwise}$$
$$= r(i, \alpha_2)(r(i, \alpha_1)(s))$$
$$r(i, \text{if } \varphi \text{ then} \qquad = r(i, \alpha_1)(s)$$
$$\alpha_1 \text{ else } \alpha_2 \text{ fi})(s) \qquad \text{if } \pi(\varphi, s) = 1$$
$$= r(i, \alpha_2)(s) \text{ otherwise}$$
$$r(i, \text{while } \varphi \text{ do } \alpha_1 \text{ od})(s) = \{s' \mid \exists k \exists s_0 \ldots \exists s_k$$
$$[s_0 = s \wedge s_k = s' \wedge$$
$$\forall j < k[s_{j+1} \in r(i, \alpha_1)(s_j)$$
$$\wedge \pi(\varphi, s_j) = 1]$$
$$\wedge \pi(\neg\varphi, s') = 1]\}$$
$$r(i, \alpha_1 + \alpha_2)(s) \qquad = r(i, \alpha_1)(s) \cup r(i, \alpha_2)(s)$$
$$\text{where } r(i, \alpha)(S') \qquad = \bigcup_{s' \in S'} r(i, \alpha)(s')$$
$$\text{for } S' \subseteq S.$$

$$c \qquad : \mathcal{A} \times Ac \to S \to \text{bool}$$
$$c(i, \text{skip})(s) \qquad = 1$$
$$c(i, \text{fail})(s) \qquad = 0$$
$$c(i, \text{test } \varphi)(s) \qquad = \pi(\varphi, s)$$
$$c(i, \alpha_1; \alpha_2)(s) \qquad = c(i, \alpha_1)(s) \wedge$$
$$c(i, \alpha_2)(r(i, \alpha_1)(s))$$
$$c(i, \text{if } \varphi \text{ then} \qquad = c(i, \text{test } \varphi; \alpha_1)(s) \vee$$
$$\alpha_1 \text{ else } \alpha_2 \text{ fi})(s) \qquad c(i, \text{test } \neg\varphi; \alpha_2)(s)$$
$$c(i, \text{while } \varphi \text{ do } \alpha_1 \text{ od})(s) = 1 \text{ if}$$
$$\exists s' \exists k \exists s_0 \ldots \exists s_k[s_0 = s \wedge$$
$$s_k = s' \wedge \forall j < k$$
$$[s_{j+1} \in r(i, \alpha_1)(s_j) \wedge$$
$$c(i, \text{test } \varphi; \alpha_1)(s_j) = 1]$$
$$\wedge c(i, \text{test } \neg\varphi)(s') = 1]$$
$$= 0 \text{ otherwise}$$
$$c(i, \alpha_1 + \alpha_2)(s) \qquad = c(i, \alpha_1)(s) \vee c(i, \alpha_2)(s)$$
$$\text{where } c(i, \alpha)(S') \qquad = \bigwedge_{s' \in S'} c(i, \alpha)(s)$$
$$\text{for } S' \subseteq S.$$

**3.3. REMARK.** The function $\pi$ as given in the definitions 3.1 and 3.2 indicates for a given formula and a given possible world whether or not the formula holds in the possible world. The domain and codomain of the function $r$ are meant to formalize the intuition that $r$ determines for a given agent $i$, action $\alpha$ and possible world $s$, the worlds that are a possible result of the occurrence of the event $do_i(\alpha)$ in $s$. The function $c$ indicates for a given agent, action and possible world whether or not the agent is capable of performing the action in the possible world. One could also think of alternative definitions for the domain and codomain of $\pi$, $r$ and $c$, for instance $\pi : \mathcal{L} \to \wp(S)$, $r : \mathcal{A} \times Ac \to \wp(S \times S)$ and $c : \mathcal{A} \times Ac \to \wp(S)$. For example the function $c$ would then yield, for a given agent and a given action, the set of possible worlds in which the agent is capable of performing the action. Obviously, both definitions denote an intuitively identical concept.

**3.4. REMARK.** Singleton sets $\{s\}$ with $s \in S$ are usually denoted by $s$.

**3.5. REMARK.** The definitions of $r(i, \text{test } \varphi)(s)$ and $c(i, \text{test } \varphi)(s)$ as defined above are based on the idea that a test for $\varphi$ is actually an action that searches for confirmation of $\varphi$; if it is not possible to get confirmation for $\varphi$, $\text{test } \varphi$ fails. This is in fact the approach taken in dynamic logic. One could also think of a different sort of tests, namely tests that are *informative*. If $\text{test } \varphi$ is informative with respect to $\varphi$ this means that after executing $\text{test } \varphi$, the agent either knows that $\varphi$ holds or knows that $\neg\varphi$ holds. The latter approach could be formalized by replacing $\pi(\varphi, s)$ by $\pi(\mathbf{K}_i\varphi \vee \mathbf{K}_i\neg\varphi, s)$ in the definition of $r(i, \text{test } \varphi)(s)$ and $c(i, \text{test } \varphi)(s)$. For the moment we have chosen to define $r$ and $c$ for tests as done above, i.e. analogously to dynamic logic. Nevertheless alternative definitions, possibly based on the assumption that tests are informative, certainly deserve further attention.

**3.6. DEFINITION.** For all models $\mathcal{M} = < S, \pi, R, r, c >$ from $\mathfrak{M}$, for all possible worlds $s \in S$, and for all formulae $\varphi$ we define:

- $\mathcal{M}, s \models \varphi$ iff $\pi(\varphi, s) = 1$.
- The formula $\varphi$ is valid in $\mathcal{M}$, notation $\mathcal{M} \models \varphi$, iff $\mathcal{M}, s \models \varphi$ for all $s \in S$.

**3.7. DEFINITION.** Let $\mathfrak{M}$ be some class of models. For all formulae $\varphi$ we define:

- The formula $\varphi$ is satisfiable in $\mathfrak{M}$ if some Kripke model $\mathcal{M} = < S, \pi, R, r, c > \in \mathfrak{M}$ and $s \in S$ exist such that $\mathcal{M}, s \models \varphi$.
- The formula $\varphi$ is valid in $\mathfrak{M}$, notation $\models \varphi$ or $\models_{\mathfrak{M}} \varphi$, iff $\mathcal{M} \models \varphi$ for all $\mathcal{M} \in \mathfrak{M}$.

The following lemma states some properties that will be useful in proving the theorems of the following sections.

**3.8. LEMMA.** *For all models $\mathcal{M}$, for all possible worlds $s \in S$, for all formulae $\varphi$ and $\psi$ and for all actions $\alpha$ we have:*

- $\mathcal{M}, s \models \neg\varphi \Leftrightarrow \text{not } (\mathcal{M}, s \models \varphi)$
- $\mathcal{M}, s \models \varphi \wedge \psi \Leftrightarrow \mathcal{M}, s \models \varphi \text{ and } \mathcal{M}, s \models \psi$
- $\mathcal{M}, s \models \mathbf{K}_i\varphi \Leftrightarrow \forall s'[(s, s') \in R(i) \Rightarrow \mathcal{M}, s' \models \varphi]$
- $\mathcal{M}, s \models [do_i(\alpha)]\varphi \Leftrightarrow$
  $\forall s'[s' \in r(i, \alpha)(s) \Rightarrow \mathcal{M}, s' \models \varphi]$
- $\mathcal{M}, s \models <do_i(\alpha)>\varphi \Leftrightarrow$
  $\exists s'[s' \in r(i, \alpha)(s) \wedge \mathcal{M}, s' \models \varphi]$

- $\mathcal{M}, s \models \mathbf{A}_i \alpha \Leftrightarrow c(i, \alpha)(s) = 1$

**3.9. REMARK.** The definitions of both $r$ and $c$ for repetitive composition are constructive ones, and are based on a definition given in [HR83]. Our definition should be compared to that of [Moo84], in which only a fixed-point equation is given. As is shown below, a fixed-point equation on its own is not sufficient.
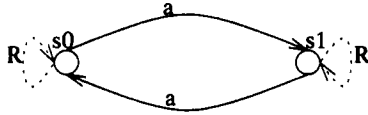
**3.10. DEFINITION.** Functions $r'(i, \text{while } \varphi \text{ do } \alpha_1 \text{ od})(s)$ and $c'(i, \text{while } \varphi \text{ do } \alpha_1 \text{ od})(s)$ are defined by:

- $r'(i, \text{while } \varphi \text{ do } \alpha_1 \text{ od})(s) =$
  $r'(i, \text{if } \varphi \text{ then } \alpha_1; \text{while } \varphi \text{ do } \alpha_1 \text{ od else skip fi})(s)$
- $c'(i, \text{while } \varphi \text{ do } \alpha_1 \text{ od})(s) =$
  $c'(i, \text{if } \varphi \text{ then } \alpha_1; \text{while } \varphi \text{ do } \alpha_1 \text{ od else skip fi})(s)$

Assume that $r'$ and $c'$ are defined analogously to $r$ and $c$ with the exception of the definition for repetitive composition. Then both $r'$ and $c'$ are partial functions.

**3.11. EXAMPLE.** Assume that the set of propositional symbols $\Pi$ is given by $\Pi = \{p\}$, the set of atomic actions is $\{a\}$ and there is only one agent.
Let the Kripke model $\mathcal{M} =< \mathcal{S}, \pi, \mathbf{R}, r, c >$ be such that:

(1) $\mathcal{S} = \{s_0, s_1\}$

(2) $\pi(p, s_0) = \pi(p, s_1) = 1$

(3) R: see the dotted lines in the picture below.

(4) r: see the solid lines in the picture below.

(5) $c(1, a)(s_0) = c(1, a)(s_1) = 1$



We find that the computations of both $r'(1, \text{while } p \text{ do } a \text{ od})(s_0)$ and $c'(1, \text{while } p \text{ do } a \text{ od})(s_0)$ do not terminate. Note that using definition 3.2 no problems concerning non-termination occur:

- $r(1, \text{while } p \text{ do } a \text{ od})(s_0) =$
  $r(1, \text{while } p \text{ do } a \text{ od})(s_1) = \emptyset$
- $c(1, \text{while } p \text{ do } a \text{ od})(s_0) =$
  $c(1, \text{while } p \text{ do } a \text{ od})(s_1) = 0$

Using our definition, the equation of definition 3.10 taken as a proposition actually holds, i.e. we have for all $s$:

- $r(i, \text{while } \varphi \text{ do } \alpha_1 \text{ od})(s) =$
  $r(i, \text{if } \varphi \text{ then } \alpha; \text{while } \varphi \text{ do } \alpha_1 \text{ od else skip fi})(s)$

- $c(i, \text{while } \varphi \text{ do } \alpha_1 \text{ od})(s) =$
  $c(i, \text{if } \varphi \text{ then } \alpha; \text{while } \varphi \text{ do } \alpha_1 \text{ od else skip fi})(s)$

The notion of determinism is defined as usual.

**3.12. DEFINITION (Determinism of events).**
Let $\mathcal{M} =< \mathcal{S}, \pi, \mathbf{R}, r, c >$ be some Kripke model from $\mathfrak{M}$ and let $do_i(\alpha)$ be some event from $\mathcal{E}_{\mathcal{A}}^{Ac}$.

- The event $do_i(\alpha)$ is deterministic in some possible world $s$ from $\mathcal{S}$ if and only if

$$|r(i, \alpha)(s)| \leq 1$$

- The event $do_i(\alpha)$ is deterministic for $\mathcal{M}$ if and only if $do_i(\alpha)$ is deterministic in all $s$ from $\mathcal{S}$.

In definition 3.13 some possible properties dealing with performance of actions by agents and the ability of agents are defined.

**3.13. DEFINITION.** Let $\mathcal{M} =< \mathcal{S}, \pi, \mathbf{R}, r, c >$ be some Kripke model from $\mathfrak{M}$ and let $\alpha$ be some action from $Ac$.

- The action $\alpha$ is $\mathbf{A}_i$-*realizable* in some possible world $s \in \mathcal{S}$ if and only if

$$\mathcal{M}, s \models \mathbf{A}_i \alpha \rightarrow <do_i(\alpha)> tt$$

holds.

- The action $\alpha$ is A-realizable in some possible world $s \in \mathcal{S}$ if and only if $\alpha$ is $\mathbf{A}_i$-realizable in $s$ for all $i \in \mathcal{A}$.

- The event $do_i(\alpha)$ is *accordant* in some possible world $s \in \mathcal{S}$ if and only if

$$\mathcal{M}, s \models \mathbf{K}_i[do_i(\alpha)]\varphi \rightarrow [do_i(\alpha)]\mathbf{K}_i\varphi$$

holds for all formulae $\varphi$.

- The action $\alpha$ is A-realizable for $\mathcal{M}$ if and only if $\alpha$ is A-realizable in every $s$ from $\mathcal{S}$.

- The event $do_i(\alpha)$ is accordant for $\mathcal{M}$ if and only if $do_i(\alpha)$ is accordant in every $s$ from $\mathcal{S}$.

Definition 3.13 states that an event $do_i(\alpha)$ is accordant, if it holds for arbitrary formula $\varphi$ that if agent $i$ knows that $\varphi$ will hold after he performs action $\alpha$, then agent $i$ will know after performing $\alpha$ that $\varphi$ holds. $\mathbf{A}_i$-realizability –and hence A-realizability– is a rather important feature of actions when considered from the point of view of planning. If some action $\alpha$ is $\mathbf{A}_i$-realizable, this implies that whenever agent $i$ is capable of performing $\alpha$, it is possible for him to perform $\alpha$.

## 3.2. Validity of formulae from $\mathcal{L}$

Theorems 3.14 and 3.15 deal with the validity of some formulae from $\mathcal{L}$ for the class of models $\mathfrak{M}$ defined in section 3. Theorem 3.14 shows some properties of events, theorem 3.15 shows some properties of capabilities.

**3.14. THEOREM.** *For all agents $i$, actions $\alpha, \alpha_1$ and $\alpha_2$, and for all formulae $\varphi$ and $\psi$ we have:*

(1) For deterministic events $\mathrm{do}_i(\alpha)$:
$$\models <\mathrm{do}_i(\alpha)> \varphi \to [\mathrm{do}_i(\alpha)]\varphi$$

(2) $\models \varphi \leftrightarrow [\mathrm{do}_i(\texttt{skip})]\varphi$

(3) $\models [\mathrm{do}_i(\texttt{fail})]\mathrm{ff}$

(4) $\models [\mathrm{do}_i(\texttt{test } \varphi)]\varphi$

(5) $\models [\mathrm{do}_i(\alpha_1;\alpha_2)]\varphi \leftrightarrow [\mathrm{do}_i(\alpha_1)]([\mathrm{do}_i(\alpha_2)]\varphi)$

(6) $\models \varphi \wedge [\mathrm{do}_i(\texttt{if } \varphi \texttt{ then } \alpha_1 \texttt{ else } \alpha_2 \texttt{ fi})]\psi$
$\leftrightarrow \varphi \wedge [\mathrm{do}_i(\alpha_1)]\psi$

(7) $\models \neg\varphi \wedge [\mathrm{do}_i(\texttt{if } \varphi \texttt{ then } \alpha_1 \texttt{ else } \alpha_2 \texttt{ fi})]\psi$
$\leftrightarrow \neg\varphi \wedge [\mathrm{do}_i(\alpha_2)]\psi$

(8) $\models [\mathrm{do}_i(\texttt{while } \varphi \texttt{ do } \alpha_1 \texttt{ od})]\neg\varphi$

(9) $\models [\mathrm{do}_i(\alpha_1 + \alpha_2)]\varphi \leftrightarrow [\mathrm{do}_i(\alpha_1)]\varphi \wedge [\mathrm{do}_i(\alpha_2)]\varphi$

**3.15. THEOREM.** *For all agents $i$, for all actions $\alpha_1$ and $\alpha_2$ and for all formulae $\varphi$ we have:*

(1) $\models \mathbf{A}_i\texttt{skip}$

(2) $\models \neg\mathbf{A}_i\texttt{fail}$

(3) $\models \varphi \leftrightarrow \mathbf{A}_i\texttt{test } \varphi$

(4) $\models \mathbf{A}_i\texttt{test } \varphi \leftrightarrow \neg\mathbf{A}_i\texttt{test } \neg\varphi$

(5) $\models \mathbf{A}_i(\alpha_1;\alpha_2) \leftrightarrow \mathbf{A}_i\alpha_1 \wedge [\mathrm{do}_i(\alpha_1)]\mathbf{A}_i\alpha_2$

(6) $\models \varphi \wedge \mathbf{A}_i(\texttt{if } \varphi \texttt{ then } \alpha_1 \texttt{ else } \alpha_2 \texttt{ fi}) \leftrightarrow$
$\varphi \wedge \mathbf{A}_i\alpha_1$

(7) $\models \neg\varphi \wedge \mathbf{A}_i(\texttt{if } \varphi \texttt{ then } \alpha_1 \texttt{ else } \alpha_2 \texttt{ fi}) \leftrightarrow$
$\neg\varphi \wedge \mathbf{A}_i\alpha_2$

(8) $\models \mathbf{A}_i(\alpha_1 + \alpha_2) \leftrightarrow \mathbf{A}_i\alpha_1 \vee \mathbf{A}_i\alpha_2$

Case 2 of theorem 3.15 states intuitively that no agent is capable of failing successfully.

Both for A-realizability and accordancy as defined in section 3.1, correspondences exist. The notion of correspondence is given by the following definition, taken from [Hoe92]. (See [Ben84] for more details on correspondence theory.)

**3.16. DEFINITION.** A frame $\mathcal{F}$ is a Kripke model without valuation $\pi$:

$$\mathcal{F} = < \mathcal{S}, \mathrm{R}, \mathrm{r}, \mathrm{c} >$$

We write $\mathcal{F} \models \varphi$ if and only if $< \mathcal{F}, \pi > \models \varphi$ for all $\pi$. If $\vartheta$ is any –usually first-order definable– property of $\mathcal{F}$ we say that formula $\varphi$ corresponds with $\vartheta$, notation $\varphi \sim_{co} \vartheta$, if for all frames $\mathcal{F}$, $\mathcal{F} \models \varphi \leftrightarrow \mathcal{F}$ satisfies $\vartheta$.

**3.17. THEOREM (Correspondences).**
*The following correspondences hold:*

(1) $\mathbf{A}_i\alpha \to <\mathrm{do}_i(\alpha)> \mathrm{tt} \sim_{co}$
$\forall s[\mathrm{c}(i,\alpha)(s) = 1 \Rightarrow \mathrm{r}(i,\alpha)(s) \neq \emptyset]$

(2) $\mathbf{K}_i[\mathrm{do}_i(\alpha)]\varphi \to [\mathrm{do}_i(\alpha)]\mathbf{K}_i\varphi \sim_{co}$
$\forall s_0 \forall s_1[\exists s_2[s_2 \in \mathrm{r}(i,\alpha)(s_0) \wedge (s_2,s_1) \in \mathrm{R}(i)] \Rightarrow$
$\exists s_3[(s_0,s_3) \in \mathrm{R}(i) \wedge s_1 \in \mathrm{r}(i,\alpha)(s_3)]]$

It turns out that for some given model A-realizability of all actions from $Ac$ can easily be achieved, by imposing a restriction on the functions $\mathrm{r}$ and $\mathrm{c}$ when applied to atomic actions.

**3.18. THEOREM.** *For all models $\mathcal{M} =< \mathcal{S}, \pi, \mathrm{R}, \mathrm{r}, \mathrm{c} >$ from $\mathfrak{M}$ we have:*

$$(\forall s \forall i \forall a \in At[\mathrm{c}(i,a)(s) = 1 \Rightarrow \mathrm{r}(i,a)(s) \neq \emptyset]) \Rightarrow$$
$$(\forall s \forall i \forall \alpha \in Ac[\mathrm{c}(i,\alpha)(s) = 1 \Rightarrow \mathrm{r}(i,\alpha)(s) \neq \emptyset])$$

**3.19. REMARK.** It turns out that A-realizability is a very desirable property. Even stronger, models that do not have the property that all occurring actions are A-realizable provide for rather counterintuitive results. For instance, suppose that in some given possible world $s$ an agent is capable of performing the action $\alpha$ but no possible world $s'$ exists such that $s' \in \mathrm{r}(i,\alpha)(s)$. Then it is not possible for the agent to execute $\alpha$. Furthermore it follows that the agent is capable of performing the action $\alpha; \texttt{fail}$ in $s$.
To avoid any counterintuitive situations we demand from now on that all occurring Kripke models are such that all actions are A-realizable, i.e. all models $< \mathcal{S}, \pi, \mathrm{R}, \mathrm{r}, \mathrm{c} >$ are such that $\forall a \in At[\mathrm{c}(i,a)(s) = 1 \Rightarrow \mathrm{r}(i,a)(s) \neq \emptyset]$ holds for all $i$ and $s$; these models are said to meet the *realizability condition.*

**3.20. REMARK.** Note that the realizability condition does not imply that $< \mathrm{do}_i(\alpha) > \mathrm{tt} \to \mathbf{A}_i\alpha$ holds, i.e. $< \mathrm{do}_i(\alpha) > \mathrm{tt}$ and $\mathbf{A}_i\alpha$ are not equivalent. To let our framework be as general and flexible as possible, we allow for the situation in which both $\neg\mathbf{A}_i\alpha$ and $< \mathrm{do}_i(\alpha) > \mathrm{tt}$ hold in a given possible world. Intuitively, this possible world represents a situation in which execution of $\alpha$ would lead to the truth of the formula $\varphi$, but since the agent is not capable of performing $\alpha$, the event $\mathrm{do}_i(\alpha)$ will not occur.

## 3.3. Transformations on actions

The notion of *action transformation* is a well-known one in the theory of software engineering. In software engineering, action transformations are used to transform programs in other, probably more efficient or better understandable, programs. (For more information on action transformations and transformational programming see [Par90].)

We define two actions to be (action)-equivalent, if they can be transformed into one another by means of a finite number of applications of action transformations. In theorem 3.23 below, it is shown that our semantics is such that the results upon termination of two equivalent actions are logically equivalent, and furthermore that the abilities of an agent are closed under some well-known transformations.

The following definition gives the action transformations that we will take into consideration. Some of these are slightly modified versions of the transformations given in [Par90].

3.21. DEFINITION. The one-step action transformations $\sim_a$ are defined as follows.

(1) $\text{skip}; \alpha \sim_a \alpha$

(2) $\alpha; \text{skip} \sim_a \alpha$

(3) $\text{fail}; \alpha \sim_a \text{fail}$

(4) $\alpha; \text{fail} \sim_a \text{fail}$

(5) $\text{test } \varphi \sim_a \text{if } \varphi \text{ then skip else fail fi}$

(6) $\text{test } \varphi; \text{test } \varphi \sim_a \text{test } \varphi$

(7) $\alpha_1; (\alpha_2; \alpha_3) \sim_a (\alpha_1; \alpha_2); \alpha_3$

(8) $\text{if tt then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \sim_a \alpha_1$

(9) $\text{if ff then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \sim_a \alpha_2$

(10) $\text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_1 \text{ fi} \sim_a \alpha_1$

(11) $\text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \sim_a$
$\text{if } \neg\varphi \text{ then } \alpha_2 \text{ else } \alpha_1 \text{ fi}$

(12) $\text{if } \varphi \text{ then if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi else } \alpha_2 \text{ fi} \sim_a$
$\text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi}$

(13) $\text{while } \varphi \text{ do } \alpha_1 \text{ od} \sim_a$
$\text{if } \varphi \text{ then } \alpha_1; \text{while } \varphi \text{ do } \alpha_1 \text{ od else skip fi}$

(14) $\text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \sim_a$
$\text{if } \varphi \text{ then } \alpha_1 \text{ else fail fi}+$
$\text{if } \neg\varphi \text{ then } \alpha_2 \text{ else fail fi}$

(15) $\text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi} \sim_a$
$\text{test } \varphi; \alpha_1 + \text{test } \neg\varphi; \alpha_2$

(16) $\alpha + \alpha \sim_a \alpha$

(17) $\alpha_1 + \alpha_2 \sim_a \alpha_2 + \alpha_1$

(18) $\alpha_1 + (\alpha_2 + \alpha_3) \sim_a (\alpha_1 + \alpha_2) + \alpha_3$

The actions $\alpha$ and $\alpha'$ from $Ac$ are *a-equivalent* (for action-equivalent), notation $\alpha \equiv_a \alpha'$, if and only if action $\alpha$ can be transformed in $\alpha'$ by application of a finite number of $\sim_a$ steps.

3.22. REMARK. To obtain equivalences under the second transformation concerning the fail-statement (transformation 4), it is used essentially that all occurring models satisfy the realizability condition. If the realizability condition is not met, the capabilities of an agent are not closed under this transformation.

Using definition 3.21 we find the following to hold.

3.23. THEOREM. *For all the actions $\alpha$ and $\alpha'$ from $Ac$, for all the agents $i$ from $\mathcal{A}$, and for all the formulae $\varphi$ we have:*

- $\alpha \equiv_a \alpha' \Rightarrow \models [do_i(\alpha)]\varphi \leftrightarrow [do_i(\alpha')]\varphi$

- $\alpha \equiv_a \alpha' \Rightarrow \models A_i\alpha \leftrightarrow A_i\alpha'$

Theorem 3.23 states that if two actions are equivalent, these actions have the same result upon performing and furthermore that if an agent is able to perform one of these actions he is also able to perform the other one.

## 4. Defining LCap

The logic LCap is defined by the following axioms and proof rules.

THE AXIOMS:

A1. All propositional tautologies and their epistemic and dynamic instances.

A2. $(K_i\varphi \wedge K_i(\varphi \rightarrow \psi)) \rightarrow K_i\psi$

A3. $K_i\varphi \rightarrow \varphi$

A4. $K_i\varphi \rightarrow K_iK_i\varphi$

A5. $\neg K_i\varphi \rightarrow K_i\neg K_i\varphi$

Ac0. $([do_i(\alpha)](\varphi \rightarrow \psi) \wedge [do_i(\alpha)]\varphi) \rightarrow [do_i(\alpha)]\psi$

Ac1. $[do_i(\text{skip})]\varphi \leftrightarrow \varphi$

Ac2. $[do_i(\text{fail})]ff$

Ac3. $[do_i(\text{test } \varphi)]\psi \leftrightarrow (\varphi \rightarrow \psi)$

Ac4. $[do_i(\alpha_1; \alpha_2)]\varphi \leftrightarrow [do_i(\alpha_1)][do_i(\alpha_2)]\varphi$

Ac5. $[do_i(\text{if } \varphi \text{ then } \alpha_1 \text{ else } \alpha_2 \text{ fi})]\psi \leftrightarrow$
$(\varphi \rightarrow [do_i(\alpha_1)]\psi) \wedge (\neg\varphi \rightarrow [do_i(\alpha_2)]\psi)$

Ac6. $[do_i(\text{while } \varphi \text{ do } \alpha_1 \text{ od})]\psi \leftrightarrow ((\neg\varphi \wedge \psi)\vee$
$(\varphi \wedge [do_i(\alpha_1)][do_i(\text{while } \varphi \text{ do } \alpha_1 \text{ od})]\psi))$

Ac7. $[do_i(\alpha_1 + \alpha_2)]\varphi \leftrightarrow [do_i(\alpha_1)]\varphi \wedge [do_i(\alpha_2)]\varphi$

Ab1. $A_i skip$
Ab2. $\neg A_i fail$
Ab3. $A_i test\ \varphi \leftrightarrow \varphi$
Ab4. $A_i(\alpha_1; \alpha_2) \leftrightarrow A_i\alpha_1 \wedge [do_i(\alpha_1)]A_i\alpha_2$
Ab5. $A_i if\ \varphi\ then\ \alpha_1\ else\ \alpha_2\ fi \leftrightarrow$
$\quad (\varphi \to A_i\alpha_1) \wedge (\neg\varphi \to A_i\alpha_2)$
Ab6. $A_i while\ \varphi\ do\ \alpha_1\ od \leftrightarrow (\neg\varphi \vee$
$\quad (\varphi \wedge A_i\alpha_1 \wedge [do_i(\alpha_1)]A_i while\ \varphi\ do\ \alpha_1\ od))$
Ab7. $A_i(\alpha_1 + \alpha_2) \leftrightarrow A_i\alpha_1 \vee A_i\alpha_2$

THE PROOF RULES:

R1. $\dfrac{(\varphi \wedge \psi) \to [do_i(\alpha_1)]\psi}{\psi \to [do_i(while\ \varphi\ do\ \alpha_1\ od)]\psi}$

R2. $\dfrac{\varphi \to (A_i\alpha_1 \wedge <do_i(while\ \varphi\ do\ \alpha_1\ od)> \neg\varphi)}{A_i while\ \varphi\ do\ \alpha_1\ od}$

R3. $\dfrac{\varphi \quad \varphi \to \psi}{\psi}$

R4. $\dfrac{\varphi}{K_i\varphi}$

R5. $\dfrac{\varphi}{[do_i(\alpha)]\varphi}$

The following definition introduces the notion of proofs for LCap.

**4.1. DEFINITION.** A proof of a formula $\varphi$ is a finite sequence of formulae such that

(1) The formula $\varphi$ is the last formula in the sequence, and

(2) each formula in the sequence is either an axiom of LCap or follows by one of the inference rules of LCap from previous formulae in the sequence.

If a proof of the formula $\varphi$ exists, $\varphi$ is called a theorem of LCap and we write $\vdash_{LCap} \varphi$, or $\vdash \varphi$ for short.

**4.2. THEOREM.** *LCap is sound with regard to $\mathfrak{M}$, i.e.*

$\vdash_{LCap} \varphi \Rightarrow \models_{\mathfrak{M}} \varphi$

**4.3. REMARK.** Although we did not yet deal with proving the completeness of LCap for $\mathfrak{M}$, we do not expect this proof to be substantially more difficult than the proofs of completeness for propositional dynamic logic (see for instance [Har84]).

Using the axioms and proof rules of LCap, it is possible to give the proof-theoretical equivalents of the proofs given in the preceding sections. In particular it is possible to prove the proof-theoretical equivalent of theorem 3.23.

**4.4. THEOREM.** *For all $\alpha$ and $\alpha'$ from Ac, for all $i$ from $A$, and for all formulae $\varphi$ we have*

- $\alpha \equiv_a \alpha' \Rightarrow \vdash [do_i(\alpha)]\varphi \leftrightarrow [do_i(\alpha')]\varphi$
- $\alpha \equiv_a \alpha' \Rightarrow \vdash A_i\alpha \leftrightarrow A_i\alpha'$

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

In [Moo84] the Can-predicate is –in a slightly different notation– introduced by:

$$Can_i(\alpha, \varphi) \Leftarrow \exists x[K_i(x = \alpha \wedge [do_i(\alpha)]\varphi)]$$

Intuitively, $Can_i(\alpha, \varphi)$ means that agent $i$ knows how to achieve $\varphi$ by performing action $\alpha$. In terms of [Moo84], the agent knows how to perform $\alpha$ since he knows what action $\alpha$ is, and he knows that performing $\alpha$ will lead to the truth of $\varphi$.

Notice that in the formula above '$\Leftarrow$' may not be replaced by '$\Leftrightarrow$': in carrying out a complex action –for instance a sequence $\alpha; \beta$– an agent does not need to know a priori what action $\alpha; \beta$ is; it suffices that he knows now what action $\alpha$ is and that he knows after performing $\alpha$ what action $\beta$ is.

As stated in section 1, we will define the Can-predicate in a different way, taking into account the abilities of the agent. The Can-predicate is introduced as a syntactic concept whose semantics is given by means of the semantics of the parts it consists of.

The intended intuitive meaning of $Can_i(\alpha, \varphi)$ as we define it is the following:

'Agent $i$ knows that he can achieve $\varphi$ by performing $\alpha$.'

Starting from this intuitive meaning, four possible definitions of the Can-predicate can be given, each of these more or less corresponding to the intended intuitive meaning.

(1) $\mathbf{Can}_i(\alpha, \varphi) \equiv [\mathrm{do}_i(\alpha)]\varphi \wedge \mathbf{A}_i\alpha$

(2) $\mathbf{Can}_i(\alpha, \varphi) \equiv [\mathrm{do}_i(\alpha)]\varphi \wedge \mathbf{K}_i\mathbf{A}_i\alpha$

(3) $\mathbf{Can}_i(\alpha, \varphi) \equiv \mathbf{K}_i[\mathrm{do}_i(\alpha)]\varphi \wedge \mathbf{A}_i\alpha$

(4) $\mathbf{Can}_i(\alpha, \varphi) \equiv \mathbf{K}_i[\mathrm{do}_i(\alpha)]\varphi \wedge \mathbf{K}_i\mathbf{A}_i\alpha$

Since we want the agent to *know* that he can achieve $\varphi$, it seems natural to demand that the agent *knows* that performing $\alpha$ will lead to $\varphi$, i.e. $\mathbf{K}_i[\mathrm{do}_i(\alpha)]\varphi$ seems to be a natural part of the definition of $\mathbf{Can}_i(\alpha, \varphi)$. In particular this means that cases 1 and 2 do not provide for intuitively acceptable definitions of the Can-predicate.

Furthermore it seems reasonable to demand that the agent knows that he is able to perform a given action. If an agent is able to perform some action $\alpha$ but he himself is not aware of this ability, then this action $\alpha$ is obviously not used in any plans of the agent.

This observation leads to the following definition of the Can-predicate.

5.1. DEFINITION. Predicate $\mathbf{Can}_i(\alpha, \varphi)$ is defined by:

$$\mathbf{Can}_i(\alpha, \varphi) \equiv \mathbf{K}_i([\mathrm{do}_i(\alpha)]\varphi \wedge \mathbf{A}_i\alpha)$$

In the spirit of definition 5.1 one could also think of a Cannot-predicate. This predicate is not just the negation of the Can-predicate, but has as its intended meaning that the agent knows that he cannot reach some goal $\varphi$ by performing some action $\alpha$, since he knows that either the action does not lead to the desired goal or he is not capable of performing the action.

5.2. DEFINITION. Predicate $\mathbf{Cannot}_i(\alpha, \varphi)$ is defined by:

$$\mathbf{Cannot}_i(\alpha, \varphi) \equiv \mathbf{K}_i(\neg[\mathrm{do}_i(\alpha)]\varphi \vee \neg\mathbf{A}_i\alpha)$$

Using definitions 5.1 and 5.2 it is possible to divide the class of actions, for a given agent and and a given goal, in three parts:

(1) Actions of which the agent knows that he can perform them to achieve the goal.

(2) Actions of which the agent knows that performing them will not lead to the goal.

(3) Actions of which the agent does not know whether they can be used to achieve the goal or not.

# 6. Properties of the Can-predicate and the Cannot-predicate

In the preceding section the definitions of the Can-predicate and the Cannot-predicate are given. In this section, we sum up some of the properties that these predicates have. In the theorems 6.1 and 6.2, we list the properties of the *means* part of the predicates, this is the $\alpha$ in $\mathbf{Can}_i(\alpha, \varphi)$ and $\mathbf{Cannot}_i(\alpha, \varphi)$. Theorems 6.3 and 6.4 list the properties of the *goal* part of the predicates, the $\varphi$ in $\mathbf{Can}_i(\alpha, \varphi)$ and $\mathbf{Cannot}_i(\alpha, \varphi)$. In theorem 6.5 two relations that exist between the Can-predicate and the Cannot-predicate are given. It furthermore turns out that both the Can-predicate and the Cannot-predicate are closed under equivalence of actions.

6.1. THEOREM ( Properties concerning the means part).
*For all agents $i$, actions $\alpha_1$ and $\alpha_2$ and for all formulae $\varphi$ and $\psi$ the following hold.*

(1) $\models \mathbf{Can}_i(\mathrm{skip}, \varphi) \leftrightarrow \mathbf{K}_i\varphi$

(2) $\models \neg\mathbf{Can}_i(\mathrm{fail}, \mathrm{tt})$

(3) $\models \mathbf{Can}_i(\mathrm{test}\ \varphi, \psi) \leftrightarrow \mathbf{K}_i(\varphi \wedge \psi)$

(4) $\models \mathbf{Can}_i(\alpha_1, \mathbf{Can}_i(\alpha_2, \varphi)) \rightarrow \mathbf{Can}_i(\alpha_1; \alpha_2, \varphi)$

(5) For accordant events $\mathrm{do}_i(\alpha_1)$:
$\models \mathbf{Can}_i(\alpha_1; \alpha_2, \varphi) \rightarrow [\mathrm{do}_i(\alpha_1)]\mathbf{Can}_i(\alpha_2, \varphi)$

(6) $\models \mathbf{Can}_i(\alpha_1 + \alpha_2, \varphi) \rightarrow \mathbf{Can}_i(\alpha_1, \varphi) \vee \mathbf{Can}_i(\alpha_2, \varphi)$

6.2. THEOREM ( Properties concerning the means part).
*For all agents $i$, actions $\alpha_1$ and $\alpha_2$ and for all formulae $\varphi$ and $\psi$ the following hold.*

(1) $\models \mathbf{Cannot}_i(\mathrm{skip}, \varphi) \leftrightarrow \mathbf{K}_i\neg\varphi$

(2) $\models \mathbf{Cannot}_i(\mathrm{fail}, \mathrm{tt})$

(3) $\models \mathbf{Cannot}_i(\mathrm{test}\ \varphi, \psi) \leftrightarrow \mathbf{K}_i(\varphi \rightarrow \neg\psi)$

(4) $\models \mathbf{Cannot}_i(\alpha_1; \alpha_2, \varphi) \rightarrow$
$\mathbf{Cannot}_i(\alpha_1, \mathbf{Can}_i(\alpha_2, \varphi))$

(5) $\models \mathbf{Cannot}_i(\alpha_1, \varphi) \wedge \mathbf{Cannot}_i(\alpha_2, \varphi) \rightarrow$
$\mathbf{Cannot}_i(\alpha_1 + \alpha_2, \varphi)$

6.3. THEOREM (Properties concerning the goal part).
*For all agents $i$, for all actions $\alpha$ and for all formulae $\varphi, \varphi_1$ and $\varphi_2$ the following hold.*

(1) $\models \mathbf{Can}_i(\alpha, \varphi_1 \wedge \varphi_2) \leftrightarrow \mathbf{Can}_i(\alpha, \varphi_1) \wedge \mathbf{Can}_i(\alpha, \varphi_2)$

(2) $\models \mathbf{Can}_i(\alpha, \varphi_1) \vee \mathbf{Can}_i(\alpha, \varphi_2) \rightarrow \mathbf{Can}_i(\alpha, \varphi_1 \vee \varphi_2)$

(3) $\models \mathbf{Can}_i(\alpha, \neg\varphi) \rightarrow \neg\mathbf{Can}_i(\alpha, \varphi)$

(4) $\models \mathbf{Can}_i(\alpha, \varphi) \wedge \mathbf{K}_i[\mathrm{do}_i(\alpha)](\varphi \rightarrow \psi) \rightarrow \mathbf{Can}_i(\alpha, \psi)$

6.4. THEOREM (Properties concerning the goal part).
*For all agents $i$, for all actions $\alpha$ and for all formulae $\varphi, \varphi_1$ and $\varphi_2$ the following hold.*

(1) $\models \mathbf{Cannot}_i(\alpha, \varphi_1) \vee \mathbf{Cannot}_i(\alpha, \varphi_2) \rightarrow$
$\mathbf{Cannot}_i(\alpha, \varphi_1 \wedge \varphi_2)$

(2) $\models \textbf{Cannot}_i(\alpha, \varphi_1 \vee \varphi_2) \rightarrow$
$\quad \textbf{Cannot}_i(\alpha, \varphi_1) \wedge \textbf{Cannot}_i(\alpha, \varphi_2)$

(3) $\models \textbf{Cannot}_i(\alpha, \varphi) \wedge \textbf{Cannot}_i(\alpha, \neg\varphi) \leftrightarrow$
$\quad \textbf{K}_i((<\text{do}_i(\alpha)> \varphi \wedge <\text{do}_i(\alpha)> \neg\varphi) \vee \neg\textbf{A}_i\alpha)$

(4) $\models \textbf{Cannot}_i(\alpha, \varphi) \wedge \textbf{K}_i[\text{do}_i(\alpha)](\psi \rightarrow \varphi) \rightarrow$
$\quad \textbf{Cannot}_i(\alpha, \psi)$

6.5. THEOREM. *For all agents $i$, for all actions $\alpha$ and for all formulae $\varphi$ the following hold.*

(1) $\models \textbf{Can}_i(\alpha, \varphi) \rightarrow \neg\textbf{Cannot}_i(\alpha, \varphi)$

(2) $\models \textbf{Can}_i(\alpha, \varphi) \rightarrow \textbf{Cannot}_i(\alpha, \neg\varphi)$

From case 4 of theorem 6.3 and case 4 of theorem 6.4, the following corollary can be derived.

6.6. COROLLARY. *The following rules are, for all agents $i$, actions $\alpha$ and formulae $\varphi$ and $\psi$, sound for $\mathfrak{M}$:*

$$\frac{\textbf{Can}_i(\alpha, \varphi) \quad \varphi \rightarrow \psi}{\textbf{Can}_i(\alpha, \psi)}$$

$$\frac{\textbf{Cannot}_i(\alpha, \varphi) \quad \psi \rightarrow \varphi}{\textbf{Cannot}_i(\alpha, \psi)}$$

The following corollary is a consequence of theorem 3.23.

6.7. COROLLARY. *For all agents $i$, actions $\alpha$ and $\alpha'$ and for all formulae $\varphi$ we have:*

- $\alpha \equiv_a \alpha' \Rightarrow \models \textbf{Can}_i(\alpha, \varphi) \leftrightarrow \textbf{Can}_i(\alpha', \varphi)$

- $\alpha \equiv_a \alpha' \Rightarrow \models \textbf{Cannot}_i(\alpha, \varphi) \leftrightarrow \textbf{Cannot}_i(\alpha', \varphi)$

## 7. Discussion

In this paper we introduced a new operator that denotes the capabilities of a rational agent. Combining this operator with a logic that combines knowledge and actions provides for an interesting new approach that can serve as a base for further development.

Our future research regarding the topics introduced in this paper will be focused on the following four points:

- Considering various possible definitions of $\pi(\textbf{test } \varphi, s)$, for nontrivial $\varphi$ containing events, and the consequences of modifying this definition on the closure of abilities under the action transformations.

- Proving the completeness of LCap for the class of models $\mathfrak{M}$.

- Defining a first-order extension of $\mathcal{L}$ and extending the semantics and the logic LCap to deal with this.

- Introducing defeasible reasoning and investigate planning in uncertain circumstances.

## References

[Ben84] J. van Benthem. Correspondence theory. In D.M. Gabbay and F. Guenthner, editors, *Handbook of Philosophical Logic*, volume 2, pages 167–247. Reidel, Dordrecht, 1984.

[FH88] R. Fagin and J.Y. Halpern. Belief, awareness and limited reasoning. *Artificial Intelligence*, 34:39–76, 1988.

[GL87] M.P. Georgeff and A.L. Lansky, editors. *Reasoning about Actions and Plans*. The 1986 Workshop, Morgan Kaufmann, 1987.

[Har79] D. Harel. *First-Order Dynamic Logic*, volume 68 of *LNCS*. Springer-Verlag, 1979.

[Har84] D. Harel. Dynamic logic. In D.M. Gabbay and F. Guenthner, editors, *Handbook of Philosophical Logic*, volume 2, pages 497–604. Reidel, Dordrecht, 1984.

[HM85] J.Y. Halpern and Y.O. Moses. A guide to the modal logics of knowledge and belief. In *Proc. 9th IJCAI*, pages 480–490, 1985.

[Hoa69] C.A.R. Hoare. An axiomatic basis for computer programming. *Communications of the ACM*, 12:576–580, 1969.

[Hoe92] W. van der Hoek. *Modalities for Reasoning about Knowledge and Quantities*. PhD thesis, Vrije Universiteit Amsterdam, 1992.

[HR83] J. Halpern and J. Reif. The propositional dynamic logic of deterministic, well-structured programs. *Theoretical Computer Science*, 27:127–165, 1983.

[MH] J.-J. Ch. Meyer and W. van der Hoek. Epistemic logic for AI and computer science. Manuscript.

[Moo84] R.C. Moore. A formal theory of knowledge and action. Technical Report 320, SRI International, 1984.

[Par90] H.A. Partsch. *Specification and Transformation of Programs*. Springer Verlag, 1990.

[Pol92] M.E. Pollack. The uses of plans. *Artificial Intelligence*, 57:43–68, 1992.

159