

# ECASL : A Model of Rational Agency for Communicating Agents

Shakil M. Khan  
Department of Computer Science  
York University  
Toronto, ON, Canada M3J 1P3  
skhan@cs.yorku.ca

Yves Lespérance  
Department of Computer Science  
York University  
Toronto, ON, Canada M3J 1P3  
lesperan@cs.yorku.ca

## ABSTRACT

The Cognitive Agent Specification Language (CASL) is a framework for specifying and verifying complex communicating multiagent systems. In this paper, we develop an extended version, ECASL, which incorporates a formal model of means-ends reasoning suitable for a multiagent context. In particular, we define a simple model of cooperative ability, give a definition of rational plans, and show how an agent's intentions play a role in determining her next actions. This bridges the gap between intentions to achieve a goal and intentions to act. We also show that in the absence of interference, an agent that is able to achieve a goal, intends to do so, and is acting rationally will eventually achieve it.

## Categories and Subject Descriptors

I.2 [Artificial Intelligence]: Problem Solving, Control Methods, and Search; I.2.11 [Artificial Intelligence]: Distributed Artificial Intelligence—*Multiagent systems*

## General Terms

Theory, Design, Verification

## Keywords

Agent Theory, Rationality, Intentions, Agent Communication

## 1. INTRODUCTION

Most agent theories [1, 13] suffer from a similar problem: they axiomatize the relation between the different mental attitudes of the agents and the physical states of the world, but they do not account for how the agents will achieve their goals, how they plan and commit to plans. Ideally, an agent's intention to achieve a state of affairs in a situation should drive the agent to intend to execute a plan that she thinks is rational in that situation. In other words, an

agent's future directed intentions should lead her to adopt rational plans and eventually achieve her intentions.

Another recent thread in agent theory introduces a procedural component to the framework in an attempt to close the gap between agents' intentions to achieve a state of affairs and their intentional actions, as well as to support the modeling of complex multiagent systems. One example of this is the Cognitive Agent Specification Language (CASL) [20, 21], which is a framework for specifying and verifying complex communicating multiagent systems. However, it is somewhat restricted in the sense that it requires the modeler to specify agent behavior explicitly, and the program that controls the agent's actions need not be consistent with the agent's intentions, or do anything to achieve them.

In this paper, we propose a solution to this problem by developing an extended version of CASL, ECASL. In particular, we define rational plans and ability in a multiagent context, and use these notions to link future and present directed intentions. We introduce a special action, the *commit* action, that makes the agent commit to a plan, and define a meta-controller *BehaveRationallyUntil* that has the agent act rationally to achieve a specific goal by choosing and committing to a rational plan, and carrying it out. Then we show that given that an agent has an intention, she will act to achieve it provided that she is able to do so.

The paper is organized as follows: in the next section, we outline previous work on CASL. In Section 3, we develop a simple formalization of cooperative ability for agents working in a multiagent setting. In Section 4, we define rational plans, relate future and present directed intentions, and discuss what it means for an agent to behave rationally. We also state a theorem that links an agent's intentions and abilities to the eventual achievement of her intentions.

## 2. CASL

In CASL [20, 21], agents are viewed as entities with mental states, i.e., knowledge and goals, and the specifier can define the behavior of the agents in terms of these mental states. CASL combines a declarative action theory defined in the situation calculus with a rich programming/process language, ConGolog [2]. Domain dynamics and agents' mental states are specified declaratively in the theory, while system behavior is specified procedurally in ConGolog.

In CASL, a dynamic domain is represented using an action theory formulated in the situation calculus [11], a (mostly) first order language for representing dynamically changing worlds in which all changes are the result of named actions.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

AAMAS'05, July 25-29, 2005, Utrecht, Netherlands.

Copyright 2005 ACM 1-59593-094-9/05/0007 ...\$5.00.

CASL uses a theory that includes the following set of axioms:

- action precondition axioms, one per action,
- successor state axioms (SSA), one per fluent, that encode both effect and frame axioms and specify exactly when the fluent changes [14],
- initial state axioms describing what is true initially including the mental states of the agents,
- axioms identifying the agent of each action,
- unique name axioms for actions, and
- domain-independent foundational axioms describing the structure of situations [6].

Within CASL, the behavior of agents is specified using the notation of the logic programming language ConGolog [2]. A typical ConGolog program is composed of a sequence of procedure declarations, followed by a complex action. Complex actions can be composed using constructs that include primitive actions ( $a$ ), waiting for a condition ( $\phi?$ ), sequence ( $\delta_1; \delta_2$ ), nondeterministic branch ( $\delta_1 \mid \delta_2$ ), nondeterministic choice of arguments ( $\pi x. \delta$ ), conditional branching (If  $\phi$  Then  $\delta_1$  Else  $\delta_2$  EndIf), while loop (While  $\phi$  Do  $\sigma$  EndWhile), and procedure call ( $\beta(\vec{p})$ ). Intuitively,  $\pi x. \delta$  nondeterministically picks a binding for the variable  $x$  and performs the program  $\delta$  for this binding of  $x$ . ConGolog also supports nondeterministic iteration, concurrent execution with and without priorities, and interrupts. To deal with multiagent processes, primitive actions in CASL take the agent of the action as argument.

The semantics of the ConGolog process description language is defined in terms of *transitions*. Two special predicates *Final* and *Trans* are introduced, and are characterized by defining axioms for each of the above constructs, where *Final*( $\delta, s$ ) means that program  $\delta$  may legally terminate in situation  $s$ , and where *Trans*( $\delta, s, \delta', s'$ ) means that program  $\delta$  in situation  $s$  may legally execute one step, ending in situation  $s'$  with program  $\delta'$  remaining.<sup>1</sup> The overall semantics of a program is specified by the *Do* relation:

$$Do(\delta, s, s') \doteq \exists \delta'. (Trans^*(\delta, s, \delta', s') \wedge Final(\delta', s'))$$

$Do(\delta, s, s')$  holds if and only if  $s'$  can be reached by performing a sequence of transitions starting with program  $\delta$  in  $s$ , and the remaining program  $\delta'$  may legally terminate in  $s'$ . Here,  $Trans^*$  is the reflexive transitive closure of the transition relation *Trans*.

CASL incorporates a branching time temporal logic, where each situation has a linear past and a branching future. In this framework, one can write both state formulas and path formulas. A state formula  $\phi(s)$  takes a single situation as argument and is evaluated with respect to that situation. On the other hand, a path formula  $\psi(s_1, s_2)$  takes two situations as arguments and is evaluated with respect to the interval (finite path)  $[s_1, s_2]$ . We often use  $\phi$  (and  $\psi$ ) to denote a formula whose fluents may contain a placeholder constant *now* (*now* and *then*, resp.) that stands for the situation in which  $\phi$  ( $\psi$ , resp.) must hold.  $\phi(s)$  (and  $\psi(s_1, s_2)$ ) is the

<sup>1</sup>Since we have predicates that take programs as arguments, we need to encode programs as first-order terms as in [2]. For notational simplicity, we suppress this encoding and use formulae as terms directly.

formula that results from replacing *now* with  $s$  (*now* and *then* with  $s_1$  and  $s_2$ , resp.). Where the intended meaning is clear, we sometimes suppress the placeholder(s).

CASL allows the specifier to model agents in terms of their mental states by including operators to specify agents' information (i.e., their knowledge), and motivation (i.e., their goals or intentions). We usually use state formulas within the scope of knowledge, and path formulas within the scope of intentions. Following [12, 18], CASL models knowledge using a possible worlds account adapted to the situation calculus.  $K(agt, s', s)$  is used to denote that in situation  $s$ ,  $agt$  thinks that she could be in situation  $s'$ .  $s'$  is called a *K-alternative situation* for  $agt$  in  $s$ . Using  $K$ , the knowledge or belief of an agent,  $Know(agt, \phi, s)$ , is defined as  $\forall s'(K(agt, s', s) \supset \phi(s'))$ , i.e.  $agt$  knows  $\phi$  in  $s$  if  $\phi$  holds in all of  $agt$ 's  $K$ -accessible situations in  $s$ . In CASL,  $K$  is constrained to be reflexive, transitive, and euclidean in the initial situation to capture the fact that agents' knowledge is true, and that agents have positive and negative introspection. As shown in [18], these constraints then continue to hold after any sequence of actions since they are preserved by the successor state axiom for  $K$ .

Scherl and Levesque [18] showed how to capture the changes in beliefs of agents that result from actions in the successor state axiom for  $K$ . These include knowledge-producing actions that can be either binary sensing actions or non-binary sensing actions. Following [9], the information provided by a binary sensing action is specified using the predicate  $SF(a, s)$ , which holds if the action  $a$  returns the binary sensing result 1 in situation  $s$ . Similarly for non-binary sensing actions, the term  $sf f(a, s)$  is used to denote the sensing value returned by the action.

Lespérance [7] extends the SSA of  $K$  in [18] to support two variants of the *inform* communicative action, namely *informWhether* and *informRef*. Here,  $inform(inf, agt, \phi)$ ,  $informWhether(inf, agt, \psi)$ , and  $informRef(inf, agt, \theta)$  mean that  $inf$  informs  $agt$  that  $\phi$  currently holds,  $inf$  informs  $agt$  about the current truth value of  $\psi$ , and  $inf$  informs  $agt$  of who/what  $\theta$  is, respectively. The preconditions of *inform* are as follows:

$$\begin{aligned} Poss(inform(inf, agt, \phi), s) &\equiv Know(inf, \phi, s) \\ &\wedge \neg Know(inf, Know(agt, \phi, now), s). \end{aligned}$$

In other words, the agent  $inf$  can inform  $agt$  that  $\phi$ , iff  $inf$  knows that  $\phi$  currently holds, and does not believe that  $agt$  currently knows that  $\phi$ . The preconditions of *informWhether* and *informRef* are similar to that of *inform*. The SSA for  $K$  is defined as follows:

$$\begin{aligned} K(agt, s^*, do(a, s)) &\equiv \\ \exists s'. [ &K(agt, s', s) \wedge s^* = do(a, s') \wedge Poss(a, s') \wedge \\ &((BinarySensingAction(a) \wedge Agent(a) = agt) \\ &\quad \supset (SF(a, s') \equiv SF(a, s))) \wedge \\ &((NonBinarySensingAction(a) \wedge Agent(a) = agt) \\ &\quad \supset (sf f(a, s') = sf f(a, s))) \wedge \\ \forall inf, \phi. ( &a = inform(inf, agt, \phi) \supset \phi(s')) \wedge \\ \forall inf, \psi. ( &a = informWhether(inf, agt, \psi) \\ &\quad \supset (\psi(s') \equiv \psi(s))) \wedge \\ \forall inf, \theta. ( &a = informRef(inf, agt, \theta) \\ &\quad \supset (\theta(s') = \theta(s)))]. \end{aligned}$$

This says that after an action happens, every agent learns that it has happened. Moreover, if the action is a sensing action, the agent performing it acquires knowledge of the associated proposition or term. Furthermore, if the action involves someone informing *agt* that  $\phi$  holds, then *agt* knows this afterwards, and similarly for *informWhether* and *informRef*. Note that this axiom only handles knowledge expansion, not revision.

CASL also incorporates goal expansion and a limited form of goal contraction. Goals or intentions are modeled using an accessibility relation  $W$  over possible situations. The  $W$ -accessible situations for an agent are the ones where she thinks that all her goals are satisfied.  $W$ -accessible situations may include situations that the agent thinks are impossible, unlike Cohen and Levesque's [1]  $G$ -accessible worlds. But intentions are defined in terms of the more primitive  $W$  and  $K$  relations so that the intention accessible situations are  $W$ -accessible situations that are also compatible with what the agent knows, in the sense that there is a  $K$ -accessible situation in their history. This guarantees that agents' intentions are realistic, that is, agents can only intend things that they believe are possible. Thus we have:

$$\begin{aligned} \text{Int}(agt, \psi, s) &\doteq \forall s', s^*. [W(agt, s^*, s) \\ &\wedge K(agt, s', s) \wedge s' \leq s^*] \supset \psi(s', s^*). \end{aligned}$$

This means that the intentions of an agent in  $s$  are those formulas that are true for all intervals between situations  $s'$  and  $s^*$  where the situations  $s^*$  are  $W$ -accessible from  $s$  and have a  $K$ -accessible situation  $s'$  in their history. Intentions are future oriented, and any goal formula will be evaluated with respect to a finite path defined by a pair of situations, a beginning situation  $s'$  and an ending situation  $s^*$ . This formalization of goals can deal with both achievement goals and maintenance goals. An achievement goal  $\psi$  is said to be satisfied if  $\psi$  holds between *now* and *then*, i.e., if  $\text{Eventually}(\psi, \text{now}, \text{then})$ , which is defined as  $\exists s'. (\text{now} \leq s' \leq \text{then} \wedge \psi(s'))$ . In [19], Shapiro showed how positive and negative introspection of intentions can be modeled by placing some constraints on  $K$  and  $W$ . To make sure that agents' wishes and intentions are consistent,  $W$  is also constrained to be serial.

The SSA for  $W$  which handles intention change in CASL, has the same structure as a SSA for a domain dependent fluent. In the following,  $W^+(agt, a, s^*, s)$  ( $W^-(agt, a, s^*, s)$ , resp.) denotes the conditions under which  $s^*$  is added to (dropped from, resp.)  $W$  as a result of the action  $a$ :

$$\begin{aligned} W(agt, s^*, do(a, s)) &\equiv \\ W^+(agt, a, s^*, s) &\vee (W(agt, s^*, s) \wedge \neg W^-(agt, a, s^*, s)). \end{aligned}$$

An agent's intentions are expanded when it is requested something by another agent. After the  $request(req, agt, \psi)$  action, *agt* adopts the goal that  $\psi$ , unless she has a conflicting goal or is not willing to serve *req* for  $\psi$ . Therefore, this action should cause *agt* to drop any paths in  $W$  where  $\psi$  does not hold. This is handled in  $W^-$ :

$$\begin{aligned} W^-(agt, a, s^*, s) &\doteq \text{IncompRequest}(agt, a, s^*, s), \\ \text{IncompRequest}(agt, a, s^*, s) &\doteq \\ [\exists req, \psi. a = request(req, agt, \psi) \\ &\wedge \text{Serves}(agt, req, \psi, s) \wedge \neg \text{Int}(agt, \neg \psi, s) \\ &\wedge \exists s'. K(agt, s', s) \wedge s' \leq s^* \\ &\wedge \neg \psi(do(a, s'), s^*)]. \end{aligned}$$

Here, the *request* action is considered a primitive action. The preconditions of request are:

$$\text{Poss}(request(req, agt, \psi), s) \equiv \text{Int}(req, \psi, s).$$

A limited form of intention contraction is also handled in CASL. Suppose that the agent *req* requests *agt* that  $\psi$  and later decides it no longer wants this. The requester *req* can perform the action  $cancelRequest(req, agt, \psi)$ , which causes *agt* to drop the goal that  $\psi$ .  $cancelRequest$  actions are handled by determining what the  $W$  relation would have been if the corresponding *request* action had never happened. This type of goal contraction is handled in  $W^+$ , which can be defined as follows:

$$\begin{aligned} W^+(agt, a, s^*, s) &\doteq \exists s_1. W(agt, s^*, s_1) \\ &\wedge \exists a_1. do(a_1, s_1) \leq s \wedge \text{Cancels}(a, a_1) \\ &\wedge (\forall a', s'. do(a_1, s_1) < do(a', s') \leq s \supset \\ &\quad \neg W^-(agt, a', s^*, s')), \\ \text{Cancels}(a, a') &\doteq [\exists req, \psi. a' = request(req, agt, \psi) \\ &\quad \wedge a = cancelRequest(req, agt, \psi)]. \end{aligned}$$

Suppose that a  $cancelRequest$  action occurs in situation  $s$ . The  $W$  relation is first restored to the way it was before the corresponding *request* action occurred, i.e., in  $s_1$ . Then starting just after the *request*, all the actions  $a'$  that occurred in the history of  $s$  (say in situation  $s'$ ) are considered, and any situation  $s^*$  in  $W$  that satisfies  $W^-(agt, a', s^*, s')$  is removed from  $W$ . A  $cancelRequest$  action can only be executed if a corresponding *request* action has occurred in the past.

### 3. SIMPLE COOPERATIVE ABILITY

An agent cannot be expected to eventually achieve an intention just because she has that intention, and she is acting rationally. We also need to make sure that the agent is *capable* of achieving the goal in the current situation [8]. In a single agent domain, an agent's ability can roughly be defined as her knowledge of a plan that is physically and epistemically executable and whose execution achieves the goal. However, modeling multiagent ability is a more complex problem, since in this case we need to consider the agents' knowledge about each other's knowledge and intentions as well as how they choose actions, behave rationally, etc. In this section, we develop a simple model of cooperative ability of agents suitable for a limited multiagent context in the absence of exogenous actions, i.e., actions whose performance is not intended by the planning agent. In an open multiagent framework, agents' actions may interfere with each other, possibly perturbing their plans. In some cases, there are multiple strategies to achieve a common goal, and the agents may fail unless they coordinate their choice of strategy by reasoning about each other's knowledge, ability, and rational choice. Moreover, agents may have conflicting goals or intentions. To simplify, we restrict our framework by only allowing plans where the actions that the other agents must do are fully specified, i.e., action delegation is possible, but (sub)goal delegation is not. The primary agent, who is doing the planning, is constrained to know the whole plan in advance. Thus, the primary agent is allowed to get help from others, but she can only ask them to perform specific actions. Given this, we do not need to model the fact that the other agents behave rationally.

When dealing with ability, it is not enough to say that the agent is able to achieve a goal iff she has a physically executable plan, and any execution of this plan starting in the current situation achieves the goal. We should also take into account the epistemic and intentional feasibility of the plan. This is necessary as physical executability does not guarantee that the executor will not get stuck in a situation where it knows that some step must be performed, but does not know which. For example, consider the plan  $(a; \text{If } \phi \text{ Then } b \text{ Else } c \text{ EndIf}) \mid d$ , where actions  $a$ ,  $b$ ,  $c$  and  $d$  are always possible, but where the agent does not know whether  $\phi$  holds after  $a$ . If the agent follows the branch where the first action is  $a$ , she will get stuck due to incomplete knowledge. Hence, the result of deliberation should be a kind of plan where the executor will know what to do next at every step, a plan that does not itself require deliberation to interpret. To deal with this, Sardiña *et al.* [17] defined the notion of *Epistemically Feasible Deterministic Programs* (EFDPs) for single agent plans and characterized deliberation in terms of it. Note that EFDPs are deterministic, since they are the result of deliberation and their execution should not require making further choices.

Since we are dealing with cooperative multiagent ability, we also need to make sure that the cooperating agents intend to perform the requested actions when it is their turn to act. We extend the notion of EFDP to handle simple multi-agent plans as follows. A program is called an *Epistemically and Intentionally Feasible Deterministic Program* (EIFDP) in situation  $s$  for agent  $agt$ , if at each step of the program starting at  $s$ ,  $agt$  always has enough information to execute the next action in the program, or knows that the executor of the next action is another agent, and that this agent has enough information to execute this action and intends to do it. Put formally:

$$\begin{aligned}
& EIFDP(agt, \delta, s) \doteq \\
& \quad \forall \delta', s'. \text{Trans}^*(\delta, s, \delta', s') \supset LEIFDP(agt, \delta', s'), \\
& LEIFDP(agt, \delta, s) \doteq \\
& \quad \text{Know}(agt, \text{Final}(\delta, now) \wedge \\
& \quad \quad \neg \exists \delta', s'. \text{Trans}(\delta, now, \delta', s'), s) \\
& \quad \vee \exists \delta'. \text{Know}(agt, \neg \text{Final}(\delta, now) \wedge \\
& \quad \quad \text{UTrans}(\delta, now, \delta', now), s) \\
& \quad \vee \exists \delta', a. \text{Know}(agt, \neg \text{Final}(\delta, now) \wedge \\
& \quad \quad \text{Agent}(a) = agt \wedge \\
& \quad \quad \text{UTrans}(\delta, now, \delta', do(a, now)), s) \\
& \quad \vee \exists \delta', agt'. \text{Know}(agt, \neg \text{Final}(\delta, now) \wedge \\
& \quad \quad \exists a. \text{UTrans}(\delta, now, \delta', do(a, now)) \wedge \\
& \quad \quad \text{Agent}(a) = agt' \neq agt \wedge \\
& \quad \quad \text{Int}(agt', \exists s'. s' \leq then \\
& \quad \quad \quad \wedge Do(a, now, s'), now), s).
\end{aligned}$$

Thus to be an EIFDP, a program must be such that all configurations reachable from the initial program and situation, involve a *Locally Epistemically and Intentionally Feasible Deterministic Program* (LEIFDP). A program is a LEIFDP in a situation with respect to an agent, if the agent knows that the program is currently in its *Final* configuration and no further transitions are possible, or knows that she is the agent of the next action and knows what unique transition (with or without an action) it can perform next, or knows

that someone else  $agt'$  is the agent of the next action, that  $agt'$  knows what the action is and intends to do it next, and knows what unique transition the program can perform next with this action. Here,  $UTrans(\delta, s, \delta', s')$  means that the program  $\delta$  in  $s$  can perform a unique transition, which takes the agent to  $s'$  with the remaining program  $\delta'$ . Note that when it is the other agent's turn,  $agt$  does not have to know exactly what the next action is, i.e., know all the parameters of the next action. However, at every step, she must know what the remaining program is.

EIFDPs are suitable results for planning. They can always be executed successfully and since they are deterministic, they do not require further deliberation to execute. Using EIFDP, ability can be defined as follows:

$$\begin{aligned}
& Can(agt, \psi(now, then), s) \doteq \\
& \quad \exists \delta. \text{Know}(agt, EIFDP(agt, \delta, now) \wedge \\
& \quad \quad \exists s'. Do(\delta, now, s') \wedge \\
& \quad \quad \forall s'. (Do(\delta, now, s') \supset \psi(now, s')), s).
\end{aligned}$$

Thus, an agent can achieve a goal in situation  $s$ , iff she knows of a plan  $\delta$  that is an EIFDP, is executable starting at  $s$ , and any possible execution of the plan starting in the current situation brings about the goal.

We use the following as our running example (adapted from [12]) throughout the paper. Consider a world in which there is a safe with a combination lock. If the safe is locked and the correct combination is dialed, then the safe becomes unlocked. However, dialing the incorrect combination will cause the safe to explode. The agent can only dial a combination if the safe is intact, and it is not possible to change the combination of the safe. Initially, the agent  $Ag_1$  has the intention to open the safe, but does not know the combination. However, she knows that  $Ag_2$  knows it. She also knows that  $Ag_2$  is willing to serve/help her, and that  $Ag_2$  does not have the intention of not informing her of the combination of the safe. Here are some of the axioms that we use to model this domain:

$$\begin{aligned}
& sf_1) Poss(a, s) \supset [Exploded(do(a, s)) \equiv \\
& \quad \exists c, agt. (a = dial(agt, c) \wedge Comb(s) \neq c) \\
& \quad \vee Exploded(s)]. \\
& sf_2) Poss(dial(agt, c), s) \equiv \neg Exploded(s). \\
& sf_3) Agent(dial(agt, c)) = agt. \\
& sf_4) \neg Exploded(S_0).
\end{aligned}$$

The first axiom, a successor state axiom, states that the safe has exploded after doing action  $a$  iff  $a$  denotes the action of dialing the wrong combination, or if the safe has already exploded. The second axiom, a precondition axiom, states that it is possible to dial a combination for the safe in situation  $s$  iff the safe is intact in  $s$ . The third axiom is an agent axiom and defines the agent of the *dial* action. The last axiom is an initial situation axiom, and states that the safe is initially intact. From now on, we will use  $D_{safe}$  to denote the set of axioms that we use to model this safe domain (see [5] for the complete axiomatization).

Now, consider the following plan:<sup>2</sup>

$$\begin{aligned} \sigma_{safe} = & \text{requestAct}(Agt_1, Agt_2, \\ & \text{informRef}(Agt_2, Agt_1, \text{Comb}(s))); \\ & \text{informRef}(Agt_2, Agt_1, \text{Comb}(s)); \\ & \text{dial}(Agt_1, \text{Comb}(s)). \end{aligned}$$

So, the plan is that  $Agt_1$  will request  $Agt_2$  to inform her of the combination of the safe,  $Agt_2$  will inform  $Agt_1$  of the combination of the safe, and finally,  $Agt_1$  will dial the combination to open the safe. We claim that  $\sigma_{safe}$  is an EIFDP in the initial situation for  $Agt_1$ , and that  $Agt_1$  is able to achieve her intention of opening the safe in the initial situation  $S_0$ :

THEOREM 1.

- a.  $D_{safe} \models EIFDP(Agt_1, \sigma_{safe}, S_0)$ .
- b.  $D_{safe} \models Can(Agt_1, \text{Eventually}(\neg \text{Locked}), S_0)$ .

(a) holds as all configurations reached by  $\sigma_{safe}$  starting in  $S_0$  are LEIFDP. (b) holds as  $Agt_1$  knows of a plan (i.e.,  $\sigma_{safe}$ ), which she knows is an EIFDP and is executable, and knows that any execution of this plan ends up in a situation where the safe is unlocked.

#### 4. FROM INTENTIONS THAT TO INTENTIONS TO ACT

In this section, we define rational plans and extend CASL to model the role of intention and rationality in determining an agent's actions. This bridges the gap between future directed intentions and present directed ones. We also present a theorem that relates intention and ability to the eventual achievement of intended goals.

Before going further, let us discuss the communication actions that we will use in ECASL. Like in CASL, we use three primitive informative communication actions, namely, *inform*, *informWhether*, and *informRef*. However, unlike in CASL, we provide two intention transfer communication actions, *request* and *requestAct*, and these are defined in terms of *inform*.<sup>3</sup> The *request* action can be used by an agent to request another agent to achieve some state of affairs, whereas *requestAct* involves an agent's request to another agent to perform some particular complex action starting in the next situation. Formally,

$$\begin{aligned} \text{request}(req, agt, \phi) & \doteq \\ & \text{inform}(req, agt, \text{Int}(req, \phi, now)), \\ \text{requestAct}(req, agt, \delta) & \doteq \\ & \text{request}(req, agt, \exists a. \text{DoNext}(\delta, do(a, now))) \\ & \wedge \text{Agent}(\delta) = agt. \end{aligned}$$

Here  $\text{DoNext}(\delta, s)$  is an abbreviation for  $\exists s^*. s \leq s^* \leq \text{then} \wedge \text{Do}(\delta, s, s^*)$ , and  $\text{Agent}(\delta) = agt$  means that the agent of all actions in  $\delta$  is  $agt$ . In our specification, we only allow sincere requests. That is, an agent can perform a request

<sup>2</sup>*requestAct* is an abbreviation introduced in the next section; it denotes a special kind of request, namely, a request to perform an action.

<sup>3</sup>A similar account of request was presented by Herzig and Longin [4], where it is defined as inform about intentions, and the requested goals are adopted via cooperation principles.

if the request is not contradictory to her current intentions. So defining requests as informing of intentions is reasonable. However, since requests are modeled in terms of *inform*, and since we are using true belief, the account seems to be overly strict. For instance, in the safe domain,  $\sigma_{safe}$  seems like a rational plan for  $Agt_1$  in the initial situation. However, initially  $Agt_1$  does not have the intention that  $Agt_2$  informs her the combination of the safe. So we cannot show that  $\sigma_{safe}$  is rational, since it requires  $Agt_1$  to know that she has the intention before she can inform about it. One way to solve this is to relax the preconditions of *inform*. However, this can have problematic consequences, as someone could inform of something without knowing it, and this might require belief revision by the addressee. Later, we will discuss another way to avoid this problem by building commitment into plans. For now, we just assume that initially  $Agt_1$  has the intention that  $Agt_2$  informs her of the combination of the safe.

To allow the cancellation of requests, we also provide two actions, namely, *cancelRequest*, and *cancelReqAct*. Unlike CASL where *cancelRequest* is primitive, we define it using *inform*. These two actions are defined as follows:

$$\begin{aligned} \text{cancelRequest}(req, agt, \psi) & \doteq \\ & \text{inform}(req, agt, \neg \text{Int}(req, \psi, now)), \\ \text{cancelReqAct}(req, agt, \delta) & \doteq \\ & \text{cancelRequest}(req, agt, \exists s^*. s^+, prev). \\ & prev = do(\text{requestAct}(req, agt, \delta), s^+) \\ & \wedge s^+ < now \leq s^* \leq \text{then} \wedge \text{Do}(\delta, prev, s^*) \\ & \wedge \text{Agent}(\delta) = agt. \end{aligned}$$

Now let us look at what plans are *rational* for an agent. An agent that is acting rationally, should prefer some plans to others. To this end, we define an ordering on plans:

$$\begin{aligned} \succeq (agt, \delta_1, \delta_2, s) & \doteq \\ & \forall s'. K(agt, s', s) \wedge \exists s^*. \text{Do}(\delta_2, s', s^*) \wedge W(agt, s^*, s) \\ & \supset [\exists s^*. \text{Do}(\delta_1, s', s^*) \wedge W(agt, s^*, s)]. \end{aligned}$$

That is, a plan  $\delta_1$  is as good as another plan  $\delta_2$  in situation  $s$  for an agent  $agt$  iff for all  $W$ -accessible situations that can be reached by following  $\delta_2$  from a situation that is  $K$ -accessible from  $s$ , say  $s'$ , there exists a  $W$ -accessible situation that can be reached from  $s'$  by following  $\delta_1$ . In other words,  $\delta_1$  is at least as good as  $\delta_2$  if it achieves the agent's goals in all the possible situations where  $\delta_2$  does.

Using EIFDP and the  $\succeq$  relation, we next define rational plans. A plan  $\delta$  is said to be *rational* in situation  $s$  for an agent  $agt$  if the following holds:

$$\begin{aligned} \text{Rational}(agt, \delta, s) & \doteq \\ & \forall \delta'. \succeq (agt, \delta', \delta, s) \supset \succeq (agt, \delta, \delta', s) \\ & \wedge EIFDP(agt, \delta, s). \end{aligned}$$

Thus, a rational plan in a situation  $s$ , is a plan that is as good as any other plan in  $s$  and is an EIFDP in  $s$ .

For example, consider the plan  $\sigma_{safe}$ . We claim that  $\sigma_{safe}$  is as good as any other plan available to  $Agt_1$  in the initial situation, and that  $\sigma_{safe}$  is rational in the initial situation.

THEOREM 2.

- a.  $D_{safe} \models \forall \sigma. \succeq (Agt_1, \sigma_{safe}, \sigma, S_0)$ .
- b.  $D_{safe} \models \text{Rational}(Agt_1, \sigma_{safe}, S_0)$ .

Since this plan achieves  $Agt_1$ 's intention of opening the safe starting in any situation that is  $K$ -accessible to  $S_0$ , (a) holds. (b) follows from the fact that  $\sigma_{safe}$  is as good as any other plan in  $S_0$  and is an EIFDP in  $S_0$ .

In most cases, there are many rational plans (i.e., ways of achieving as many goals as possible). The decision of which plan the agent commits to is made based on pragmatic/non-logical grounds. We do not model this here. Instead, we introduce a  $commit(agt, \delta)$  action that will model the agent's committing to a particular plan  $\delta$ , more specifically, committing to executing  $\delta$  next. The action precondition axiom for the  $commit$  action is as follows:

$$Poss(commit(agt, \delta), s) \equiv \neg Int(agt, \neg DoNext(\delta, now), s).$$

That is, the agent  $agt$  can commit to a plan  $\delta$  in situation  $s$ , iff the agent currently does not have the intention that the actions in the plan do not happen next.

Next, we extend the SSA for  $W$  seen earlier to handle intention revision as a result of the agent's commitment to a rational plan. We modify  $W^-$  as follows:

$$W^-(agt, a, s^*, s) \doteq IncompRequest(agt, a, s^*, s) \vee \\ IncompCommit(agt, a, s^*, s).$$

Here,  $IncompCommit$  handles the expansion of the agent's intentions that occur when a  $commit$  action occurs. We define  $IncompCommit$  as follows:

$$IncompCommit(agt, a, s^*, s) \doteq \\ [\exists \delta. a = commit(agt, \delta) \wedge \exists s'. s' \leq s^* \wedge K(agt, s', s) \\ \wedge \neg \exists s^{**}. (s' < s^{**} \leq s^* \wedge Do(\delta, do(a, s'), s^{**}))].$$

So, after the performance of a  $commit$  action in  $s$ , a  $W$ -accessible situation  $s^*$  in  $s$  will be dropped from  $agt$ 's new set of  $W$ -accessible situations if the committed to action does not happen next over the interval between the  $W$ -accessible situation  $s^*$  and its predecessor  $s'$  that is  $K$ -accessible from the current situation  $s$ .

The definition of  $W^+$  remains unchanged. Note that if exogenous actions are allowed, agents need to revise their commitments when an exogenous action occurs by uncommitting from the currently committed plan, and committing to a new rational plan. We return to this issue in Section 5.

We now show that our formalization of intentions has some desirable properties:

THEOREM 3.

- a.  $\models \neg Int(agt, \neg \phi, s) \wedge Serves(agt, req, \phi, s) \\ \wedge Poss(request(req, agt, \phi), s) \supset \\ Int(agt, \phi, do(request(req, agt, \phi), s)).$
- b.  $\models Poss(commit(agt, \delta), s) \supset \\ Int(agt, DoNext(\delta, now), do(commit(agt, \delta), s)).$

(a) says that if an agent  $agt$  does not have the intention that not  $\phi$  in  $s$ , then she will have the intention that  $\phi$  in the situation resulting from another agent  $req$ 's request to  $agt$  that  $\phi$  in  $s$ , provided that she is willing to serve  $req$  on  $\phi$ , and that the  $request$  action is possible in  $s$ . (b) states that if an agent  $agt$  does not have the intention of not performing a complex action  $\delta$  in  $s$  (i.e. if  $commit(agt, \delta)$  is possible in  $s$ ), then she will have the intention of performing it after she commits to it.

As mentioned earlier, the problem that arises as a result of defining requests as informing of intentions can be solved by

building commitment into plans. Consider the safe example; we assumed earlier that initially  $Ag_1$  has the intention that  $Ag_2$  informs her the combination of the safe. We can now relax this constraint by considering  $\sigma_{safe}^*$  to be our new rational plan, where  $\sigma_{safe}^* = commit(Agt_1, \sigma_{safe}); \sigma_{safe}$ , i.e. the new plan is that  $Ag_1$  commits to  $\sigma_{safe}$  and then  $\sigma_{safe}$  is performed. Since  $Ag_1$  commits to  $\sigma_{safe}$  before she executes it, the SSA for  $W$  will make her adopt the intention that  $Ag_2$  informs her the combination of the safe after she requests  $Ag_2$  to do so, and thus we do not need to assume that this holds.

$commit$  provides a way to link future directed intentions and present directed ones. We next specify a generic meta-controller for an agent that arbitrarily chooses a rational plan, commits to it, and executes it. Then we can prove a theorem about the relationship between intention, ability, and the eventual achievement of an intended goal. This theorem serves as a proof of soundness of our agent theory.

The following meta-controller allows us to refer to the future histories of actions that may occur for an agent who is behaving rationally until  $\psi$  holds. Rational behavior until  $\psi$  can be defined as follows (we assume that there are no exogenous actions):

$$BehaveRationallyUntil(agt, \psi(now)) \doteq \\ \pi \delta. Rational(agt, \delta, now)?; commit(agt, \delta); \\ While \neg \psi(now) Do \\ If \exists a. Int(agt, do(a, now) \leq then, now) \wedge \\ Agent(a) = agt) Then \\ [\pi a. (Int(agt, do(a, now) \leq then, now) \wedge \\ Agent(a) = agt)?; a] \\ Else \\ \pi agt'. [Int(agt, \exists a. do(a, now) \leq then \wedge \\ Agent(a) \neq agt \wedge \\ Agent(a) = agt', now)?; \\ (\pi a'. Int(agt', do(a', now) \leq then, now)?; a')] \\ EndIf EndWhile.$$

That is, rational behavior until  $\psi$  can be defined as arbitrarily choosing a rational plan, committing to it, and then executing it as long as  $\psi$  does not hold. A rational plan can have actions by the planning agent and by other agents. When it is the planning agent's turn to act, she should perform the action that she intends to perform next; otherwise, she should wait for the other agent to act. When it is the other agent's turn, it will perform the action that it is supposed to perform, because rational plans are EIFDP, and thus the other agent must intend to do the action required by the plan. Note that we only deal with achievement goals here.

One problem with CASL is that the execution of plans is viewed from the system's perspective rather than from the agents' perspective. So, although CASL includes operators that model agents' knowledge and goals, the system behavior is simply specified as a set of concurrent processes. To deal with this problem, Lespérance [7] proposed an account of subjective plan execution in CASL; the program construct  $Subj(agt, \delta)$  ensures that  $\delta$  can be executed by  $agt$  based on her own knowledge state. We have extended this notion to deal with multiagent plans (i.e. plans with actions by agents other than the executor) and consider other agents'

intentions; see [5] for the formal details.

Next, we present our “success theorem”:<sup>4</sup>

**THEOREM 4.** *From Commitment and Ability to Eventuality*

$$\begin{aligned} \models & [\text{OInt}(agt, \text{Eventually}(\gamma, now, then), s) \\ & \wedge \text{Can}(agt, \text{Eventually}(\gamma, now, then), s) \\ & \wedge \text{Int}(agt, \text{Eventually}(\psi, now, then), s)] \supset \\ & \text{AllDo}(\text{Subj}(agt, \text{BehaveRationallyUntil}(agt, \psi)), s). \end{aligned}$$

Intuitively, if in some situation, an agent intends to achieve some goal and is able to achieve all its intentions, then the agent will eventually achieve the goal in all rational histories from that situation.  $\text{OInt}(agt, \psi, s)$  means that  $\psi$  is all the intentions that  $agt$  has in  $s$ . This construct must be used as we have to assume that the agent is able to achieve all her intentions. If this is not the case, the agent may have to choose between some of its goals and the *BehaveRationallyUntil* operator will not guarantee that a specific goal (i.e.,  $\psi$ ) will be achieved. If there are exogenous actions, then a more generic meta-controller can be defined. We discuss this in the next section.

We also have the following corollary for the safe domain:

**COROLLARY 1.**

$$\begin{aligned} D_{safe} \models & \text{AllDo}(\text{Subj}(Agt_1, \\ & \text{BehaveRationallyUntil}(Agt_1, \neg \text{Locked})), S_0). \end{aligned}$$

We have shown in Theorem 1(b) that  $Agt_1$  can achieve her intention of opening the safe in the initial situation. Moreover, the only intention of  $Agt_1$  is to open the safe. It follows from Theorem 4 that  $Agt_1$  will eventually open the safe if she behaves rationally starting in  $S_0$  (see [5] for a complete proof).

## 5. DISCUSSION AND FUTURE WORK

In this paper, we have presented a formal theory of agency that deals with simple multiagent cooperation and shows how future directed intentions and present directed ones can be related. An agent’s current rational plans depend on her current intentions. The *commit* action models how the agent’s intentions can be updated to include a commitment to a rational plan. Using this, we have formulated a planning framework for multiple cooperating and communicating agents in CASL. We specified how an agent’s future directed intentions will lead the agent to adopt a rational plan and then carry it out using the meta-controller *BehaveRationallyUntil*.

To relate agents’ intentions with their actions, Cohen and Levesque [1] required that agents eventually drop all their intentions either because they had been achieved or because they were viewed as impossible to achieve (AKA the *no infinite deferral* assumption). A similar account was presented by Rao and Georgeff [13]. We believe that this no infinite deferral assumption should be a consequence of an agent behaving rationally as specified by other axioms of the theory, rather than be imposed separately. A more intuitive account was presented in [22], where Singh showed that rather than

having it as an assumption, the no infinite deferral principle can be derived from the theory.

In [15], Sadek introduces some axioms to incorporate explicit principles of rational behaviour in his adaptation of the Cohen and Levesque framework. The application of these axioms makes it possible for an agent to build rational plans in a deductive way by inference of causal chains of intention, without needing to resort to a separate planner. From an operational point of view, agents in this framework generate plans using a backward chaining planning mechanism. Sadek uses the rational effect of communication actions as an integral part of his specification. These rational effects express the reasons which lead an agent to select an action, and are related to perlocutionary effects. However, it is not specified under what conditions the rational effects become actual effects, and one cannot reason about these conditions. Moreover, the planning mechanism in [15] is incomplete and many rational plans cannot be inferred. In [16], Sadek et al. describe how this theory is used to develop an implemented rational agent engine called ARTIMIS. This technology has been used to build natural language dialogue systems and multiagent applications. Louis [10] recently extended ARTIMIS [16] to incorporate a more general model of planning (state space planning by regression and hierarchical planning) and plan adoption. His framework is more complex than ours and uses defaults (as does Sadek’s). The approach supports multiagent plans and has been implemented. But there is no formalization of epistemically feasible plans, and no success theorem. Commitment to a plan is modelled using a special predicate rather than using the intention attitude.

Although independently motivated, our account closely resembles the one in [24], where a similar notion of commitment to actions was introduced to relate intentions and actions. However, that framework does not model rationality or provide a success theorem. There has also been related work that tries to extend agent programming languages to support declarative goals (e.g. [23]).

Our semantics of communication acts is mentalistic, in contrast to recent social commitment semantics (e.g. [3]). The public social commitment level is obviously important, but we don’t think that communication can be reduced to it. The reason agents communicate is that this serves their private goals. One must usually reason about these goals and the associated beliefs to really understand the agents’ behavior.

The theory presented here is a part of our ongoing research on the semantics of speech acts and communication in the situation calculus. In [5], we present an extended version of our framework where we allow exogenous actions. To deal with these unintended actions, an agent needs to revise the plan it is committed to whenever an exogenous action occurs. In other words, she needs to un-commit from the previously committed plan, consider the new set of rational plans, and commit to one of them. We handle the un-committing part in the SSA for  $W$ . The agents’ commitment to a new rational plan is handled using a more sophisticated meta-controller. This controller iterates the *BehaveRationallyUntil* program as long as the goal remains un-achieved and there is a plan that is rational in the current situation. In [5], we also define a notion of conditional commitment, and model some simple communication protocols using it.

<sup>4</sup>The construct *AllDo* is a strict version of *Do* that requires that all possible executions of a program terminate successfully; see [7] for a formal definition.

Our current agent theory is overly simplistic in many ways. One strict constraint is that we do not allow cooperating agents to choose how they will achieve the goals delegated to them since we assume that the planning agent knows the whole plan in advance. Only one agent is assumed to do planning. In future work, we will try to relax this restriction and to model some interaction protocols that involve multiple planning agents. It would also be interesting to try to use this formalization to implement flexible communication agents as in [16] and to develop tools to support multiagent programming that conform to ECASL.

## 6. REFERENCES

- [1] P. Cohen and H. Levesque. Intention is choice with commitment. *Artificial Intelligence*, 42(2-3):213–361, 1990.
- [2] G. De Giacomo, Y. Lespérance, and H. Levesque. Congolog, a concurrent programming language based on the situation calculus. *Artificial Intelligence*, 121:109–169, 2000.
- [3] N. Fornara, F. Vigano, and M. Colombetti. Agent communication and institutional reality. In R. van Eijk, M.-P. Huget, and F. Dignum, editors, *Agent Communication: Proc. of the AAMAS 04 Workshop on Agent Communication*, vol. 3396 of *LNAI*. Springer, 2005.
- [4] A. Herzig and D. Longin. A logic of intention with cooperation principles and with assertive speech acts as communication primitives. In *Proc. of the First International Joint Conference on AAMAS*, 2002.
- [5] S. Khan. A situation calculus account of multiagent planning, speech acts, and communication. Master’s thesis, Dept. of Computer Science, York University, Toronto, ON, Canada, 2005 (In preparation).
- [6] G. Lakemeyer and H. Levesque. Aol: A logic of acting, sensing, knowing, and only-knowing. In *Proc. of the 6th International Conference on Principles of KR&R*, pages 316–327, 1998.
- [7] Y. Lespérance. On the epistemic feasibility of plans in multiagent systems specifications. In J.-J. C. Meyer and M. Tambe, editors, *Proc. of the 8th International Workshop on ATAL*, vol. 2333 of *LNAI*, pages 69–85. Springer, 2002.
- [8] Y. Lespérance, H. Levesque, F. Lin, and R. Scherl. Ability and knowing how in the situation calculus. *Studia Logica*, 66(1):165–186, 2000.
- [9] H. Levesque. What is planning in the presence of sensing? In *Proc. of the Thirteenth National Conference on AI*, pages 1139–1146, Portland, OR, 1996.
- [10] V. Louis. *Conception et Mise en Oeuvre de Modèles Formels de Calcul de Plans d’Action Complexes par un Agent Rationnel Dialoguant*. PhD thesis, Université de Caen, Caen, France, 2002.
- [11] J. McCarthy and P. Hayes. Some philosophical problems from the standpoint of artificial intelligence. *Machine Intelligence*, 4:463–502, 1969.
- [12] R. Moore. A formal theory of knowledge and action. In J. Hobbs and R. Moore, editors, *Formal Theories of the Commonsense World*, pages 319–358. Ablex, 1985.
- [13] A. Rao and M. Georgeff. Modeling rational agents within a BDI-architecture. In R. Fikes and E. Sandewall, editors, *Proc. of the 2nd International Conference on Principles of KR&R*, pages 473–484, 1991.
- [14] R. Reiter. The frame problem in the situation calculus: A simple solution (sometimes) and a completeness result for goal regression. In V. Lifschitz, editor, *Artificial Intelligence and Mathematical Theory of Computation: Papers in the Honor of John McCarthy*. Academic Press, San Diego, CA, 1991.
- [15] D. Sadek. Communication theory = rationality principles + communicative act models. In *Proc. of the AAAI 94 Workshop on Planning for Interagent Communication*, 1994.
- [16] D. Sadek and P. Bretier. ARTIMIS: Natural dialogue meets rational agency. In *Proc. of the Fifteenth IJCAI*, pages 1030–1035, 1997.
- [17] S. Sardiña, G. De Giacomo, Y. Lespérance, and H. Levesque. On the semantics of deliberation in indilog - from theory to implementation. *Annals of Mathematics and Artificial Intelligence*, 41(2-4):259–299, 2004.
- [18] R. Scherl and H. Levesque. Knowledge, action, and the frame problem. *Artificial Intelligence*, 144(1-2), 2003.
- [19] S. Shapiro. *Specifying and Verifying Multiagent Systems Using CASL*. PhD thesis, Dept. of Computer Science, University of Toronto, Toronto, ON, Canada, 2005.
- [20] S. Shapiro and Y. Lespérance. Modeling multiagent systems with the cognitive agents specification language – a feature interaction resolution application. In C. Castelfranchi and Y. Lespérance, editors, *Intelligent Agents VII: Proc. of the 2000 Workshop on ATAL*, vol. 1986 of *LNAI*, pages 244–259. Springer, 2001.
- [21] S. Shapiro, Y. Lespérance, and H. Levesque. The cognitive agents specification language and verification environment for multiagent systems. In C. Castelfranchi and W. Johnson, editors, *Proc. of the 1st Int. Joint Conference on AAMAS*, pages 19–26, 2002.
- [22] M. Singh. *Multiagent Systems: A Theoretical Framework for Intentions, Know-How, and Communications*, vol. 799 of *LNAI*. Springer, 1994.
- [23] W. van der Hoeka, K. Hindriks, F. de Boer, and J.-J. C. Meyer. Agent programming with declarative goals. In C. Castelfranchi and Y. Lespérance, editors, *Intelligent Agents VII: Proc. of the 7th International Workshop ATAL 2000*, vol. 1986 of *LNAI*. Springer, 2000.
- [24] B. van Linder, W. van der Hoek, and C. M. J.-J. Formalising motivational attitudes of agents : On preferences, goals, and commitments. In M. Wooldridge, J. Muller, and M. Tambe, editors, *Intelligent Agents vol. II – Proc. of the 1995 Workshop on ATAL*, vol. 1037 of *LNAI*, pages 17–32. Springer, 1996.