

# Creating collective intention through dialogue

Frank Dignum\*   Barbara Dunin-Kęplicz†   Rineke Verbrugge‡

\* Faculty of Mathematics and Computing Science, Technical University Eindhoven  
P.O. Box 513, 5600 MB Eindhoven, The Netherlands  
dignum@win.tue.nl

† Institute of Informatics, Warsaw University  
Banacha 2, 02-097 Warsaw, Poland  
keplicz@mimuw.edu.pl

‡ Cognitive Science and Engineering, University of Groningen  
Grote Kruisstraat 2/1, 9712 TS Groningen, The Netherlands  
rineke@tcw2.ppsw.rug.nl

**Abstract.** The process of Cooperative Problem Solving can be divided into four stages. First, finding potential team members, then forming a team followed by constructing a plan for that team. Finally, the plan is executed by the team. Traditionally, protocols like the Contract Net protocol are used for performing the first two stages of the process. In an open environment however, there can be discussion among the agents in order to form a team that can achieve the collective intention of solving the problem. For these cases fixed protocols like contract net do not suffice. In this paper we present a solution, using structured *dialogues*, with an emphasis on persuasion, that can be shown to lead to the required team formation. The dialogues are described formally using modal logics and speech acts.

## 1 Introduction

The area of Distributed Problem Solving (DPS) has been occupied for more than ten years already with solving complex problems by teams of agents. Although in this field the problem might be solved in a distributed way, the control of solving it usually lies with one agent. This means that the team of agents is either already available or can be created on demand by the controlling agent. Also, these teams of agents are collaborative by nature. Thus they are designed to participate in the team to solve the problem.

We are concerned with problems that have to be solved by a number of existing agents that are not designed to solve this particular problem together. In this case, getting the right team of agents and controlling them is of prime interest. In this setting Contract Net [24] is often proposed as a simple but effective and efficient way to distribute tasks over a number of agents in order to achieve a common goal. It basically makes use of the market mechanism of task demand and supply to match tasks with agents that are willing to perform them. The reason of the success of this approach is the fact that it uses a fixed protocol with a limited number of steps and thus is easy to implement.

In our opinion, however, this functions only in cases where the market mechanism works, namely if several agents are willing to do the task (and these can compete),

and if the tasks are all well described beforehand. We will concentrate on cases where these conditions are usually not met. Either because there is only one agent capable of performing a task and that one should be negotiated with or because the task cannot be described precisely at the beginning.

A good candidate to make conversation between agents during Cooperative Problem Solving (CPS) more flexible is Krabbe and Walton's theory of Dialogue [30]. This theory gives rules for appropriate moves within different types of dialogues. The rules direct the dialogue without completely fixing the order of the moves. These moves themselves depend on particular stages of CPS. In fact, all of them are rather complex, both from the MAS and the AI perspective. Additionally, the cooperative process takes place in a dynamic and often unpredictable environment. Our intention is to present formal means allowing realisation of some relevant forms of dialogue. These dialogue types follow a formal theory of dialogue (see [30]) and speech acts (see [27]). In this way, it will be possible to prove that in given circumstances the dialogue results in a certain outcome. This is of prime importance when constructing a MAS for automated CPS.

We base our investigation on a formal model of teamwork. Thus, four stages of CPS are distinguished according to Wooldridge and Jennings' paper [31]. We define the stages somewhat differently, however (see [14] for a discussion). The first stage is *potential recognition* in which the agent that takes the initiative tries to find out which agents are potential candidates for achieving a given overall goal and how these can be combined in a team. The second stage is *team formation*. The result of this stage is a *collective intention* among a team to achieve the overall goal to solve the problem. This is the team that will try to actually achieve the goal. The third stage is *plan formation*. Here the team divides the goal into subtasks, associates these with actions and allocates these actions to team members. In terms of motivational attitudes the end result of this stage is a *collective commitment* to perform the social plan that realizes the goal. The last stage is *plan execution* in which the team members execute the allocated actions and monitor the appropriate colleagues. If necessary a *reconfiguration* of the plan can be constructed [14].

We concentrate on the first two stages of the process and show how the collective intention resulting from team formation is built up during the dialogues. It turns out that the main type of dialogue needed is persuasion. The novelty of the presented paper lies in the combination of a theory of CPS, including agents' motivational and informational attitudes, with a theory of dialogue and speech acts. The whole process leading ultimately to the team formation is described in terms of different kinds of modal logics.

The paper is structured in the following manner. Section 2 presents the needed logical background. Section 3 discusses the individual and collective motivational attitudes that play a role during team formation. A concise typology of dialogue and speech acts is given in Section 4.2. Section 5 briefly describes an agent architecture for the first two stages of CPS. Sections 6 and 7 investigate different dialogue types during potential recognition and team formation, respectively. Finally, in section 8 conclusions and further research are discussed. This paper is a revised and extended version of [7].

## 2 Formal background

We propose the use of multi-modal logics to formalize agents' informational and motivational attitudes as well as actions they perform and their effects. In CPS, both motivational and informational attitudes are considered on the three levels: individual, social and collective. However, in the presented account concerning the phase of team formation, social and collective commitments do not play a role: they become of importance at further stages of plan generation and team action. For this reason they are not defined here. The interested reader may find them in other papers (cf. [13, 15]). For similar reasons there is no stress on actions in this paper.

### 2.1 The logical language

Individual actions and formulas are defined inductively, both with respect to a fixed finite set of agents. The basis of the induction is given in the following definition.

**Definition 1 (Language).**

The language is based on the following three sets:

- a denumerable set  $\mathcal{P}$  of *propositional symbols*;
- a finite set  $\mathcal{A}$  of *agents*, denoted by numerals  $1, 2, \dots, n$ ;
- a finite set  $\mathcal{At}$  of *atomic actions*, denoted by  $a$  or  $b$ .

Next, we will define the class of individual actions  $\mathcal{Ac}$ . They are meant to refer to agents' individual actions, and are usually represented without naming the agents.

**Definition 2 (Individual actions).**

The class  $\mathcal{Ac}$  of individual actions is defined inductively as follows:

- AC1** each atomic action  $a \in \mathcal{At}$  is an individual action;
- AC2** if  $\varphi \in \mathcal{L}$ , then  $\text{confirm } \varphi$  is an individual action; (confirmation)
- AC3** if  $\alpha_1, \alpha_2 \in \mathcal{Ac}$ , then  $\alpha_1; \alpha_2$  is an individual action; (sequential composition)
- AC4** if  $\alpha_1, \alpha_2 \in \mathcal{Ac}$ , then  $\alpha_1 \cup \alpha_2$  is an individual action; (non-deterministic choice)
- AC5** if  $\alpha \in \mathcal{Ac}$ , then  $\alpha^*$  is an individual action; (iteration)
- AC6** if  $\varphi \in \mathcal{L}$ , then  $\text{stit}(\varphi)$  is an action;

Here, in addition to the standard dynamic operators of [AC1] to [AC5], the operator  $\text{stit}$  of [AC6] stands for “sees to it that” or “brings it about that”, and has been extensively treated in [2, 28].

We inductively define a set  $\mathcal{L}$  of multi-modal formulas as follows.

**Definition 3 (Formulas).**

- F1** each atomic proposition  $p \in \mathcal{P}$  is a formula;
- F2** if  $\varphi$  and  $\psi$  are formulas, then so are  $\neg\varphi$  and  $\varphi \wedge \psi$ ;
- F3** if  $\varphi$  is a formula,  $\alpha$  is an action,  $i, j \in \mathcal{A}$ ,  $G \subseteq \mathcal{A}$ , then the following are formulas:  
**epistemic modalities**  $\text{BEL}(i, \varphi)$ ,  $\text{E-BEL}_G(\varphi)$ ,  $\text{C-BEL}_G(\varphi)$ ;

**motivational modalities**  $\text{GOAL}(i, \varphi), \text{GOAL}(i, \alpha), \text{INT}(i, \varphi), \text{INT}(i, \alpha),$   
 $\text{E-INT}_G(\varphi), \text{E-INT}_G(\alpha), \text{M-INT}_G(\varphi), \text{M-INT}_G(\alpha), \text{C-INT}_G(\varphi), \text{C-INT}_G(\alpha);$   
**temporal action modalities**  $\text{done}(i, \alpha);$   
**abilities, opportunities and willingness**  $\text{able}(i, \varphi), \text{opp}(i, \varphi), \text{willing}(i, \varphi);$   
**dynamic modalities**  $[\text{do}(i, \alpha)]\varphi.$

For the dynamic operator, we define a dual construct as follows:

$$\langle \text{do}(i, \alpha) \rangle \varphi = \neg[\text{do}(i, \alpha)]\neg\varphi;$$

In the next few subsections, all modalities defined above are given a semantics.

## 2.2 Kripke models

Each Kripke model for the language defined in the previous section consists of a set of worlds, a set of accessibility relations between worlds, and a valuation of the propositional atoms, as follows.

### Definition 4 (Kripke model).

A Kripke model is a tuple  $\mathcal{M} = (W, \{B_i : i \in \mathcal{A}\}, \{G_i : i \in \mathcal{A}\}, \{I_i : i \in \mathcal{A}\}, \{R_{i,\alpha} : i \in \mathcal{A}, \alpha \text{ is an action}\}, \text{Val}, \text{abl}, \text{op}, \text{will})$  such that the following holds:

1.  $W$  is a set of possible worlds, or states;
2. For all  $i \in \mathcal{A}$ , it holds that  $B_i, G_i, I_i \subseteq W \times W$ . They stand for the accessibility relations for each agent w.r.t. beliefs, goals, and intentions, respectively. For example,  $(w_1, w_2) \in B_i$  means that  $w_2$  is an epistemic alternative for agent  $i$  in state  $w_1$ .
3. For all  $i \in \mathcal{A}, \alpha \in \mathcal{A}c$ , it holds that  $R_{i,\alpha} \subseteq W \times W$ . They stand for the dynamic accessibility relations; for example,  $(w_1, w_2) \in R_{i,\alpha}$  means that  $w_2$  is a possible resulting state from  $w_1$  by agent  $i$  executing action  $\alpha$ . We suppose that accessibility relations for complex actions are built from those for atomic actions in the standard way, and that  $R_{i,\text{confirm}(\varphi)} = \{(w, w) \mid \mathcal{M}, w \models \varphi\}$ .
4.  $\text{Val} : \mathcal{P} \times W \rightarrow \{0, 1\}$  is the function that assigns the truth values to propositional formulas in states.
5.  $\text{abl} : \mathcal{A} \times \mathcal{L} \rightarrow \{0, 1\}$  is the ability function such that  $\text{abl}(i, \varphi) = 1$  indicates that agent  $i$  is able to achieve  $\varphi$ .
6.  $\text{op} : \mathcal{A} \times \mathcal{L} \rightarrow (W \rightarrow \{0, 1\})$  is the opportunity function such that  $\text{op}(i, \varphi)(w) = 1$  indicates that agent  $i$  has the opportunity to achieve  $\varphi$  in world  $w$ .
7.  $\text{will} : \mathcal{A} \times \mathcal{L} \rightarrow (W \rightarrow \{0, 1\})$  is the willingness function such that  $\text{will}(i, \varphi)(w) = 1$  indicates that agent  $i$  has the willingness to achieve  $\varphi$  in world  $w$ .

As to abilities, opportunities, and willingness, they are modeled in the above definition in a rather static way. It is possible to make a more refined definition, using a language that includes temporal operators. We have chosen not to do so here, because these concepts are not the main focus of this paper.

At this stage, it is possible to define the truth conditions pertaining to the propositional part of the language  $\mathcal{L}$  and to abilities, opportunities and willingness.

**Definition 5 (Semantics for the non-modal operators).**

- $\mathcal{M}, v \models p \Leftrightarrow Val(p, w) = 1$  for propositional atoms  $p$ ;
- $\mathcal{M}, v \models \neg\varphi \Leftrightarrow \mathcal{M}, v \not\models \varphi$ ;
- $\mathcal{M}, v \models \varphi \wedge \psi \Leftrightarrow \mathcal{M}, v \models \varphi$  and  $\mathcal{M}, v \models \psi$ ;
- $\mathcal{M}, v \models able(i, \varphi) \Leftrightarrow abl(i, \varphi) = 1$ ;
- $\mathcal{M}, v \models opp(i, \varphi) \Leftrightarrow op(i, \varphi)(v) = 1$ ;
- $\mathcal{M}, v \models willing(i, \varphi) \Leftrightarrow will(i, \varphi)(v) = 1$ ;

The truth conditions for formulas with dynamic operators as main modality are given in subsection 2.3; for those with epistemic main operators, the truth definitions are given in subsection 2.4; finally, for those with motivational modalities as main operators, the definitions follow in section 3.

**2.3 Dynamic logic for actions**

The valuations of complex formulas containing dynamic operators as main operator are defined as follows.

**Definition 6 (Valuation for dynamic operators).**

Let  $\varphi$  be a formula,  $\alpha \in \mathcal{Ac}$ . Then

$$\mathcal{M}, v \models [do(i, \alpha)]\varphi \Leftrightarrow \text{for all } w \text{ with } (v, w) \in R_{i, \alpha}, \mathcal{M}, w \models \varphi;$$

As usual, we have the following:

$$\mathcal{M}, v \models \langle do(i, \alpha) \rangle \varphi \Leftrightarrow \text{there exists a } w \text{ with } (v, w) \in R_{i, \alpha} \text{ and } \mathcal{M}, w \models \varphi;$$

For the dynamic logic of actions, we adapt the axiomatization PDL of propositional dynamic logic, as found in [19]:

**Definition 7 (Axioms and rules of PDL).**

- P1** All instantiations of propositional tautologies;
- P2**  $[do(i, \alpha)](\varphi \rightarrow \psi) \rightarrow ([do(i, \alpha)]\varphi \rightarrow [do(i, \alpha)]\psi)$ ; (Distribution)
- P3**  $[do(i, \text{confirm}(\varphi))]\psi \leftrightarrow (\varphi \rightarrow \psi)$ ;
- P4**  $[do(i, \alpha_1; \alpha_2)]\varphi \leftrightarrow [do(i, \alpha_1)][do(i, \alpha_2)]\varphi$ ;
- P5**  $[do(i, \alpha_1 \cup \alpha_2)]\varphi \leftrightarrow [do(i, \alpha_1)]\varphi \wedge [do(i, \alpha_2)]\varphi$ ;
- P6**  $[do(i, \alpha^*)]\varphi \rightarrow \varphi \wedge [do(i, \alpha)][do(i, \alpha^*)]\varphi$ ; (Mix)
- P7**  $(\varphi \wedge [do(i, \alpha^*)](\varphi \rightarrow [do(i, \alpha)]\varphi)) \rightarrow [do(i, \alpha^*)](\varphi)$ ; (Induction)
- PR1** From  $\varphi$  and  $\varphi \rightarrow \psi$ , derive  $\psi$ ; (Modus ponens)
- PR2** From  $\varphi$ , derive  $[do(i, \alpha)]\varphi$ . (Necessitation)

The axiom system PDL is sound and complete with respect to Kripke models with the dynamic accessibility relations  $R_{i, \alpha}$  as defined above. Its decision problem is exponential time complete, as proved by [18].

## 2.4 Beliefs

To represent beliefs, we adopt a standard  $KD45_n$ -system for  $n$  agents as explained in [17], where we take  $BEL(a, \varphi)$  to have as intended meaning “agent  $a$  believes proposition  $\varphi$ ”.  $KD45_n$  consists of the following axioms and rules for  $i = 1, \dots, n$ :

- A1** All instantiations of tautologies of the propositional calculus
- A2**  $BEL(i, \varphi) \wedge BEL(i, \varphi \rightarrow \psi) \rightarrow BEL(i, \psi)$  (Belief Distribution)
- A4**  $BEL(i, \varphi) \rightarrow BEL(i, BEL(i, \varphi))$  (Positive Introspection)
- A5**  $\neg BEL(i, \varphi) \rightarrow BEL(i, \neg BEL(i, \varphi))$  (Negative Introspection)
- A6**  $\neg BEL(i, \perp)$  (Consistency)
- R1** From  $\varphi$  and  $\varphi \rightarrow \psi$  infer  $\psi$  (Modus Ponens)
- R2** From  $\varphi$  infer  $BEL(i, \varphi)$  (Belief Generalization)

In the semantics, there are accessibility relations  $B_i$  that lead from worlds  $w$  to worlds that are consistent with agent  $i$ 's beliefs in  $w$ . Thus,  $BEL$  is defined semantically as follows:

$$w \models BEL(i, \varphi) \text{ iff } t \models \varphi \text{ for all } t \text{ such that } wB_it.$$

Note that the  $B_i$  need not be reflexive, corresponding to the fact that an agent's beliefs need not be true. On the other hand, the accessibility relations  $B_i$  are transitive, euclidean and serial. These conditions correspond to the axioms of positive and negative introspection and to the fact the agent has no inconsistent beliefs, respectively. It has been proved that  $KD45_n$  is sound and complete with respect to these semantics.

The property of negative introspection is controversial; we are agnostic about this and dropping [A5] will not have important consequences for the logical framework presented in this paper.

One can define modal operators for group beliefs. The formula  $E\text{-}BEL_G(\varphi)$  is meant to stand for “every agent in group  $G$  believes  $\varphi$ ”. It is defined semantically as  $w \models E\text{-}BEL_G(\varphi)$  iff for all  $i \in G$ ,  $w \models BEL(i, \varphi)$ , which corresponds to the following axiom:

$$\mathbf{C1} \quad E\text{-}BEL_G(\varphi) \leftrightarrow \bigwedge_{i \in G} BEL(i, \varphi)$$

A traditional way of lifting single-agent concepts to multiaгент ones is through the use of *collective belief*  $C\text{-}BEL_G(\varphi)$ . This rather strong operator is similar to the more usual one of common knowledge.  $C\text{-}BEL_G(\varphi)$  is meant to be true if everyone in  $G$  believes  $\varphi$ , everyone in  $G$  believes that everyone in  $G$  believes  $\varphi$ , etc. Let  $E\text{-}BEL_G^1(\varphi)$  be an abbreviation for  $E\text{-}BEL_G(\varphi)$ , and let  $E\text{-}BEL_G^{k+1}(\varphi)$  for  $k \geq 1$  be an abbreviation of  $E\text{-}BEL_G(E\text{-}BEL_G^k(\varphi))$ . Thus we have  $w \models C\text{-}BEL_G(\varphi)$  iff  $w \models E\text{-}BEL_G^k(\varphi)$  for all  $k \geq 1$ . Note that even collective beliefs need not be true, so  $w \models C\text{-}BEL_G(\varphi)$  need not imply  $w \models \varphi$ . Define  $t$  to be  $G$ -reachable from  $s$  if there is a path in the Kripke model from  $s$  to  $t$  along accessibility arrows  $B_i$  that are associated with members  $i$  of  $G$ . Then the following property holds (see [17]):

$$s \models C\text{-}BEL_G(\varphi) \text{ iff } t \models \varphi \text{ for all } t \text{ that are } G\text{-reachable from } s.$$

Using this property, it can be shown that the following axiom and rule can be soundly added to the union of  $KD45_n$  and [C1]:

**C2**  $C\text{-BEL}_G(\varphi) \rightarrow E\text{-BEL}_G(\varphi \wedge C\text{-BEL}_G(\varphi))$

**RC1** From  $\varphi \rightarrow E\text{-BEL}_G(\psi \wedge \varphi)$  infer  $\varphi \rightarrow C\text{-BEL}_G(\psi)$  (Induction Rule)

The resulting system is called  $KD45_n^C$ , and it is sound and complete with respect to Kripke models where all  $n$  accessibility relations are transitive, serial and euclidean [17].

Some of the ways in which individual beliefs can be generated are updating, revision, and contraction [22]. The establishment of collective beliefs among a group is more problematic. In [17] it is shown that bilateral sending of messages does not suffice to determine collective belief if communication channels may be faulty, or even if there is uncertainty whether message delivery may have been delayed. We assume that in our groups, a more general type of communication, e.g. by a kind of global announcement, can be achieved. A good reference to the problems concerning collective belief and to their possible solutions is [17, Chapter 11]. In any case, it is generally agreed that collective belief is a good *abstraction tool* to study teamwork.

### 3 Motivational attitudes

In our framework most axioms relating motivational attitudes of agents appear in two forms: one with respect to *propositions* denoted by  $\varphi$ , the other with respect to *actions* denoted by  $\alpha$ . These actions are interpreted in a generic way — we abstract from any particular form of actions: they may be complex or primitive, viewed traditionally with certain effects or with default effects [10–12], etc.

A proposition, on the other hand, reflects the particular state of affairs that an agent aims for. In other words, propositions represent the agent’s higher level goals. Again, we abstract from particular methods of achieving them; e.g. they may be realized by particular plans.

Table 1 gives the formulas appearing in this paper, together with their intended meanings.

#### 3.1 Individual goals and intentions

As to Kripke semantics, evaluation of formulas is with respect to a world  $w$ , using binary accessibility relations  $B_i$ ,  $D_i$  and  $I_i$  corresponding to each agent’s beliefs, goals (or desires), and intentions, all of which lead from a world to a world. Evaluation of formulas at worlds is defined in the obvious manner inspired by epistemic logic. Here we give only our  $n$ -agent version of the definitions for beliefs, goals and intentions, where the expression  $M, w \models \varphi$  is read as “formula  $\varphi$  is satisfied by world  $w$  in structure  $M$ ”. For  $i = 1, \dots, n$  we have:

$$M, w \models \text{GOAL}(i, \varphi) \text{ iff } \forall v \text{ with } wD_iv, M, v \models \varphi$$

GOAL( $a, \varphi$ )	agent $a$ has as a goal that $\varphi$ be true
GOAL( $a, \alpha$ )	agent $a$ has as a goal to do $\alpha$
stit( $a, \varphi$ )	agent $a$ sees to it that $\varphi$ holds
done( $a, \alpha$ )	agent $a$ has done $\alpha$ at the previous moment
INT( $a, \varphi$ )	agent $a$ has the intention to make $\varphi$ true
INT( $a, \alpha$ )	agent $a$ has the intention to do $\alpha$
E-INT $_G$ ( $\varphi$ )	every agent in group $G$ has the individual intention to make $\varphi$ true
E-INT $_G$ ( $\alpha$ )	every agent in group $G$ has the individual intention to do $\alpha$
C-INT $_G$ ( $\varphi$ )	group $G$ has the collective intention to make $\varphi$ true
C-INT $_G$ ( $\alpha$ )	group $G$ has the collective intention to do $\alpha$

**Table 1.** Formulas and their intended meaning

$$M, w \models \text{INT}(i, \varphi) \text{ iff } \forall v \text{ with } wI_iv, M, v \models \varphi$$

As for axioms: for the epistemic operator BEL the modal system  $KD45$  is used, which we adapt to  $KD45_n$  for  $n$  agents (see the previous section). For the motivational operators GOAL and INT the axioms include the system  $K$ , which we adapt for  $n$  agents to  $K_n$ . For  $i = 1, \dots, n$  the following axioms and rules are included:

**A1** All instantiations of tautologies of the propositional calculus

**R1** From  $\varphi$  and  $\varphi \rightarrow \psi$  infer  $\psi$  (Modus Ponens)

**A2<sub>D</sub>** (GOAL( $i, \varphi$ )  $\wedge$  GOAL( $i, \varphi \rightarrow \psi$ )  $\rightarrow$  GOAL( $i, \psi$ ) (Goal Distribution Axiom)

**A2<sub>I</sub>** (INT( $i, \varphi$ )  $\wedge$  INT( $i, \varphi \rightarrow \psi$ )  $\rightarrow$  INT( $i, \psi$ ) (Intention Distribution Axiom)

**R2<sub>D</sub>** From  $\varphi$  infer GOAL( $i, \varphi$ ) (Goal Generalization)

**R2<sub>I</sub>** From  $\varphi$  infer INT( $i, \varphi$ ) (Intention Generalization)

In a BDI system, an agent's activity starts from goals. As an agent may have many different objectives, its goals need not be consistent with each other. Then, the agent chooses a limited number of its goals to be intentions. We assume that they are chosen in such a way that consistency is preserved. Thus for intentions we assume, as Rao and Georgeff do, that they should be consistent. This can be formulated as follows:

**A6<sub>I</sub>**  $\neg\text{INT}(i, \perp)$  for  $i = 1, \dots, n$  (Intention Consistency Axiom)

Rao and Georgeff also add an analogous axiom for the consistency of goals. However, it was argued above that an agent's goals are not necessarily consistent with each other. Thus, we adopt the basic system  $K_n$  for goals. Nevertheless, in the presented approach other choices may be adopted without consequences for the rest of the definitions in this paper.

It is not hard to prove soundness and completeness of the basic axiom systems for goals and intentions with respect to suitable classes of models by a tableau method, and also give decidability results using a small model theorem.

### 3.2 Collective intentions

To model teamwork, individual attitudes naturally do not suffice. In other work we discussed the pairwise notion of social commitments, as well as collective notions like

collective belief, collective intention, and collective commitment [15]. In the first two stages of CPS, the essential notions are those of collective belief and collective intention.

The definition of *collective intention* is rather strong, because we focus on strictly cooperative groups. There, a necessary condition for a collective intention is that all members of the group have the associated individual intention  $\text{INT}(i, \varphi)$ . Moreover, to exclude the case of competition, all agents should *intend* all members to have the associated individual intention, as well as the intention that all members have the individual intention, and so on; we call such a mutual intention  $\text{M-INT}_G(\varphi)$ . Furthermore, all members of the group are aware of this mutual intention, that is, they have a collective belief about this ( $\text{C-BEL}_G(\text{M-INT}_G(\varphi))$ ).

In order to formalize the above two conditions,  $\text{E-INT}_G(\varphi)$  (standing for “everyone intends”) is defined by the following axiom, analogous to **C1** above:

$$\mathbf{M1} \quad \text{E-INT}_G(\varphi) \leftrightarrow \bigwedge_{i \in G} \text{INT}(i, \varphi).$$

The mutual intention  $\text{M-INT}_G(\varphi)$  is axiomatized by an axiom and rule analogous to **C2** and **RC1**:

$$\mathbf{M2} \quad \text{M-INT}_G(\varphi) \leftrightarrow \text{E-INT}_G(\varphi \wedge \text{M-INT}_G(\varphi))$$

$$\mathbf{RM1} \quad \text{From } \varphi \rightarrow \text{E-INT}_G(\psi \wedge \varphi) \text{ infer } \varphi \rightarrow \text{M-INT}_G(\psi) \text{ (Induction Rule)}$$

The system resulting from adding **M1**, **M2**, and **RM1** to  $KD_n$  is called  $KD_n^{\text{M-INT}_G}$ , and it is sound and complete with respect to Kripke models where all  $n$  accessibility relations are serial (by a proof analogous to the one for common knowledge in [17]). Finally, collective intention is defined by the following axiom:

$$\mathbf{M3} \quad \text{C-INT}_G(\varphi) \leftrightarrow (\text{M-INT}_G(\varphi) \wedge \text{C-BEL}_G(\text{M-INT}_G(\varphi)))$$

Note that this definition is different from the one given in [13, 15]; the new definition will be extensively discussed and compared with others in a forthcoming paper. Let us remark that, even though  $\text{C-INT}_G(\varphi)$  seems to be an infinite concept, a collective intention with respect to  $\varphi$  may be established in a finite number of steps: according to **RM1** and **M3**, it suffices that all agents in  $G$  intend  $\varphi \wedge \text{M-INT}_G(\varphi)$  and that this fact is announced to the whole group by which a collective belief is established.

## 4 Conversations in CPS

In this section we will briefly discuss some theory to describe conversations.

Conversations are sequences of messages between two (or more) agents. These sequences can be completely fixed as is done by the use of a fixed protocol which states exactly which message should come next. Conversations can also be seen as completely free sequences of messages. In this case the agents can decide at any moment what

will be the next message they send. We have already argued that fixed protocols are too rigid for the situation that we describe. However, complete freedom of choosing the next message would also be impractical. This would put a heavy burden on the agent that has to choose at each point in time which message it might send.

This is the reason that we choose for a form in between the two extremes sketched above, namely dialogue theory. We will first present dialogue theory, and then give a short description of speech act theory, which is used to describe the effects of utterances in dialogues between agents involved in CPS.

#### 4.1 Dialogue theory

Dialogue theory structures conversations by means of a number of dialogue rules. These rules limit the number of possible responses at each point, while not completely fixing the sequence of messages. The agents speak in turn, for example asking questions and giving replies, and take into account, at each turn, what has occurred previously in the dialogue.

Krabbe and Walton [30] provide a typology of dialogue types between two agents, with an emphasis on the persuasion dialogue. They create a *normative model*, representing the ideal way reasonable, cooperative agents participate in the type of dialogue in question. For each type of dialogue, they formulate *an initial situation*, *a primary goal*, and *a set of rules*. Below, their typology is briefly explained and adapted to the CPS. In the course of communication among agents, there often occurs a shift from one type of dialogue to another, in particular *embedding* occurs when the second dialogue is functionally related to the first one.

*A persuasion dialogue* arises from a conflict of opinions. It may be that one agent believes  $\varphi$  while some others either believe a contrary proposition  $\psi$  (where  $\varphi \wedge \psi$  is inconsistent) or just have doubt about  $\varphi$ . The goal of a persuasion dialogue is to resolve the conflict by verbal means, in such a way that a stable agreement results, corresponding to a collective informational attitude. In contrast to [30], we allow persuasion also with respect to motivational attitudes.

The initial situation of *negotiation* is a conflict of interests, together with a need for cooperation. The main goal is to make a deal. Thus, the selling and buying of goods and services, that is often described in the MAS literature, is only one of the many contexts where negotiation plays a role. Negotiation and persuasion are often not distinguished adequately.

The initial situation of *information seeking* occurs when one agent is ignorant about the truth of a certain proposition and seeks information from other agents. Two other important types of dialogue are *inquiry* and *deliberation*, but they play a role mainly during the stage of plan formation, which will be considered in a forthcoming paper. The last type of dialogue, *eristics* (verbal fighting between agents), is not relevant as we focus on teamwork.

In general, Krabbe and Walton [30] are not interested in informational and motivational attitudes of agents involved in dialogue, if these attitudes are not communicated explicitly. In contrast to them, our goal is to make the whole process of dialogue among computational agents transparent. For this reason, at each step of team formation the

agents' internal attitudes need to be established, and then updated and revised when necessary.

## 4.2 Speech acts in CPS

Austin's theory of speech acts, later refined and formalized by Searle ([26, 27]), is eminently suitable to account for the effects that a single speaker's utterance have on the mental state of the hearer. (For an interesting overview of the use of speech act theory in multiagent systems, see [29].)

Searle stated that in a speaker's utterance, the agent performs at least the following three kinds of acts:

1. the uttering of words: *utterance meaning*;
2. referring and predicating: *propositional acts*;
3. stating, questioning, commanding, promising etc.: *illocutionary acts*.

Searle characterized many types of illocutionary acts by four aspects: their propositional content, their preparatory conditions, sincerity conditions, and essential quality.

In this paper, we restrict ourselves to a small and fixed set of illocutionary acts that are relevant during the first stages of cooperative problem solving. These are *assert* (ASS), *request* (REQ), *concede* (CONCEDE), and *challenge* (CHALLENGE). For request and assert, the four characterizing aspects are defined in [26]. Thus, a request has as propositional content a future act  $\alpha$  of the hearer. As preparatory condition, the hearer must be able to do  $\alpha$  and the speaker must believe this; moreover, it shouldn't be obvious to both speaker and hearer that the hearer will do  $\alpha$  in the normal course of event of its own accord. As sincerity condition, the speaker must want the hearer to do  $\alpha$ . The essential quality of a request is that it counts as an attempt to get the hearer to do  $\alpha$ . For the characterization of the well-known illocutionary act assert, we refer the reader to [26].

Concede and challenge may be similarly defined. We will not give full formal characterizations here. Informally, a concession may be characterized as a hearer's positive reaction to another agent's assertion or request. In the first case, the conceiver should believe the other agent's assertion, but not so strongly that it can be called upon to defend it. In the second case, the concession counts as a promise to fulfil the other agent's request; for a full characterization of promises, see [26, 27]. Challenges count as negative reactions to another agent's assertion. The sincerity condition is that the challenger should not believe the propositional content of the other agent's assertion, although it may be persuaded later. Challenges follow the logical structure of the asserted proposition by pointing out a part that is disbelieved.

In addition to the three kinds of act predicated by Searle, Austin introduced the notion of the effects illocutionary acts have on the actions and attitudes of the hearer. He called such effects *perlocutionary acts*. In sections 6 and 7 we will define the perlocutionary acts resulting from the speech acts relevant during the first two stages of CPS, using dynamic logic.

## 5 Agent architecture for team formation by dialogue

In order to ensure the proper realization of team formation by dialogue, we postulate that an agent architecture should contain a number of specific modules. The heart of the system is, as usual, the **reasoning** module. When realizing the consecutive stages leading ultimately to team formation, interaction with the **planning, communication,** and **social reasoning** modules is necessary. All these modules contain a number of specific *reasoning rules* which will be introduced formally in the sequel. Each rule refers to a specific aspect of the reasoning process; very often they can be viewed as rules bridging different modules.

The **reasoning** module contains rules creating agents' individual beliefs, intentions, and goals, as well as the collective informational and motivational attitudes established during potential recognition and team formation. The **communication** module contains rules leading to speech acts, but also some auxiliary rules or mechanisms organizing the *information seeking* and *persuasion* dialogues. In this paper we abstract from the latter, rather technical, ones. The **social reasoning** module contains rules allowing to reason about other agents, while the **planning** module will be used for preplanning during the potential recognition stage. In fact, it will be more busy during higher stages of CPS, i.e. when generating a social plan and possibly during the reconfiguration process (see [14, 16]).

## 6 Potential recognition

After introducing the basic informational and motivational attitudes involved in CPS, defining speech acts and characterizing conversations, we are ready for the synthesis of all these ingredients in the formal model of dialogues, leading ultimately to team formation. As a reminder, we use Wooldridge and Jennings' formal model of teamwork([31]), obeying consecutive stages of *potential recognition; team formation*, resulting in a *collective intention* among a team to achieve the overall goal to solve the problem; *plan formation*, resulting in a *collective commitment* to perform the social plan that realizes the goal; and, finally, the last stage is *plan execution*.

Potential recognition is about finding the set of agents that may participate in the formation of the team that tries to achieve the overall goal. These agents are grouped into several potential teams with whom further discussion will follow during *team formation*. For simplicity, we assume that one agent takes the initiative to realize the overall goal, which is given.

### 6.1 The end result of potential recognition

The first task of the *initiator* is to form a partial (abstract) plan for the achievement of the overall goal. On the basis of the (type of) subgoals that it recognizes it will determine which agents might be most suited to form the team. In order to determine this match the initiator tries to find out the properties of the agents, being interested in three aspects, namely their *abilities, opportunities,* and their *willingness* to participate in team formation.

The aspect of ability concerns whether the agents can perform the right type of tasks. It does not depend on the situation, but may be viewed as an inherent property of the agent itself. The aspect of opportunity takes into account the possibilities of task performance in the present situation, involving resources and possibly other properties. The aspect of willingness considers the agents' mental attitudes towards participating towards the overall goal. Very capable agents that do not want to do the job are of no use. As a reminder, the components of the agent's suitability are represented as follows (cf. [21]):

1. the individual ability of agent  $b$  to achieve a goal  $\psi$  is denoted by  $able(b, \psi)$ ,
2. the resources available to agent  $b$  are reflected by the opportunity that agent  $b$  has to achieve  $\psi$ ; that there is such an opportunity is denoted by  $opp(b, \psi)$ .
3. the willingness of agent  $b$  to participate in team formation is denoted by  $willing(b, \varphi)$ .

Formal definitions of these notions were given in section 2.2.

## 6.2 Towards a potential of cooperation

In the preceding subsection we have described what type of information the initiating agent tries to gather in order to start team formation. Now, we will describe how the information is collected, leading to the formal outcome.

The output at this stage is the "potential for cooperation" that the initiator  $a$  sees with respect to  $\varphi$ , denoted as  $POTCOOP(a, \varphi)$ , meaning that  $\varphi$  is a goal of  $a$  ( $GOAL(a, \varphi)$ ), and that there is a group  $G$  such that  $a$  believes that  $G$  can collectively achieve  $\varphi$  ( $C-CAN_G(\varphi)$ ) and are willing to participate in team formation; and either  $a$  cannot or doesn't desire to achieve  $\varphi$  in isolation. This is expressed by the following definitions (see [14] for more discussion):

$$\begin{aligned} POTCOOP(a, \varphi) &\leftrightarrow GOAL(a, \varphi) \wedge \\ &\exists G \subseteq TBEL(a, C-CAN_G(\varphi) \wedge \forall i \in G willing(i, \varphi)) \wedge \\ &(\neg CAN(a, \varphi) \vee \neg GOAL(a, stit(\varphi))) \end{aligned}$$

$$CAN(a, \varphi) \leftrightarrow able(a, \varphi) \wedge opp(a, \varphi).$$

$POTCOOP(a, \varphi)$  is derived partly by introspection (on the agent's goal, and its lack of ability or goal to achieve it on its own. The part  $\exists G \subseteq TBEL(a, C-CAN_G(\varphi) \wedge \forall i \in G willing(i, \varphi))$  is derived from the information collected from the other individual agents. To derive  $C-CAN_G(\varphi)$  the initiator compares the information obtained about the other agents against a partial abstract plan for the overall goal  $\varphi$ . For this purpose  $\varphi$  is split into a number of subgoals  $\varphi_1, \dots, \varphi_n$ , which can be viewed as instrumental to the overall goal. Together they *realize*  $\varphi$  and are compared with the individual abilities and opportunities that the agents are believed to have:

$$\begin{aligned} C-CAN_G(\varphi) &\leftrightarrow \exists \varphi_1, \dots, \exists \varphi_n (realize(\langle \varphi_1, \dots, \varphi_n \rangle, \varphi) \wedge \\ &\forall i \leq n \exists j \in G (able(j, \varphi_i) \wedge opp(j, \varphi_i))) \end{aligned}$$

Here,  $realize(\langle \varphi_1, \dots, \varphi_n \rangle, \varphi)$  intuitively means: "if  $\varphi_1, \dots, \varphi_n$  are achieved, then  $\varphi$  holds". Technically, one may need some extra-logical, context-dependent, reasoning to show the implication. In the reasoning module of the agent architecture, there is a formal rule corresponding to: "if  $\langle \varphi_1, \dots, \varphi_n \rangle$  is a result of pre-planning to achieve  $\varphi$ , then  $realize(\langle \varphi_1, \dots, \varphi_n \rangle, \varphi)$  holds".

### 6.3 Information seeking dialogue

Ultimately, the initiator has to form beliefs about the abilities, opportunities, and willingness of the individual agents in order to derive  $POTCOOP(a, \varphi)$ . The possible strategies for organizing this information seeking part remain out of our interest. For example,  $a$  may first investigate the willingness of particular agents, and on this basis then ask the interested ones about their abilities and opportunities. Depending on the specific situation we deal with, another organization of this process may be more adequate. In any case, the questions in this stage form part of an *information seeking* dialogue. This can be done by  $a$  asking every agent about its properties and the agent responding with the requested information.

Formally this can be expressed as follows, where  $\psi$  may stand for any of the aspects  $opp(i, \varphi_i)$ ,  $able(i, \varphi_i)$ ,  $willing(i, \varphi)$ , etc.

$$\begin{aligned} & [REQ_{a,i}(\mathbf{if} \psi \mathbf{then} ASS_{i,a}(\psi) \mathbf{else} ASS_{i,a}(\neg\psi))] \\ & [ASS_{i,a}(\psi)](TRUST(a, i, \psi) \rightarrow BEL(a, \psi)) \end{aligned}$$

$$\begin{aligned} & [REQ_{a,i}(\mathbf{if} \psi \mathbf{then} ASS_{i,a}(\psi) \mathbf{else} ASS_{i,a}(\neg\psi))] \\ & [ASS_{i,a}(\neg\psi)](TRUST(a, i, \neg\psi) \rightarrow BEL(a, \neg\psi)) \end{aligned}$$

The above formulae are based on the formal theory on *speech acts* developed in [27] and [9]. They are expressed as dynamic logic formulas of the form  $[\alpha_1][\alpha_2]\psi$ , meaning that if  $\alpha_1$  is performed then always a situation arises such that if  $\alpha_2$  is performed then in the resulting state  $\psi$  will always hold. In the above case  $\alpha_1$  is the complex action  $REQ_{a,i}(\mathbf{if} \psi \mathbf{then} ASS_{i,a}(\psi) \mathbf{else} ASS_{i,a}(\neg\psi))$ , where  $REQ_{a,i}(\alpha)$  stands for agent  $a$  requesting agent  $i$  to perform the action  $\alpha$ . After this request  $i$  has three options.

1. It can simply ignore  $a$  and not answer at all.
2. It can state that it is not willing to divulge this information:  
 $ASS_{i,a}(\neg(\mathbf{if} \psi \mathbf{then} ASS_{i,a}(\psi) \mathbf{else} ASS_{i,a}(\neg\psi)))$ .
3. It can state that it does not have enough information:  
 $ASS_{i,a}(\neg(BEL(i, \psi) \wedge \neg BEL(i, \neg\psi)))$ .
4. It can either assert that  $\psi$  is the case (as described above) or that it is not, in which case  $a$  believes  $\psi$  is not true:

$$[ASS_{i,a}(\neg\psi)](TRUST(a, i, \neg\psi) \rightarrow BEL(a, \neg\psi)).$$

Of course in case 2, agent  $a$  can already derive that  $i$  is not willing to achieve  $\varphi$  as part of a team; only in case 4 will  $a$  have a resulting belief about  $\psi$ .

An important notion here is whether  $a$  trusts the other agent's answer; we use  $TRUST(a, i, \psi)$  to mean "agent  $a$  trusts agent  $i$  with respect to proposition  $\psi$ ". Trust is a complex concept that has been defined in many ways, from different perspectives (see [4] for some current work in this area). We will not try to define trust in any way in this paper, but simply use its intuitive meaning. We suppose that the module social reasoning contains rules giving a number of conditions (e.g. " $a$  has observed that agent  $i$  is generally trustworthy") implying  $TRUST(a, i, \psi)$ .

To sum up, an information seeking dialogue about all ingredients of  $POTCOOP(a, \varphi)$  takes place in the way described. In principle, the schema of all necessary questions may

be rather complex. For example, when recognizing the ability to achieve a specific subgoal  $\varphi_i$ , agent  $a$  should repeat this question for all the subgoals it distinguished to every agent, but this solution is apparently not acceptable from the AI perspective. Of course it is more effective for agent  $a$  to ask each agent  $i$  to divulge all the abilities it has with respect to achieving this set of goals. These are, however, strategic considerations, not related to the theory of dialogue.

A next strategic point is about case 1. To avoid that agent  $a$  will be waiting indefinitely for an answer, we assume that every speech act has an implicit deadline for reaction incorporated. After this deadline, the silent agent will not be considered as a potential team member anymore. The logical modeling of these types of deadlines is described in [8] and will not be pursued here.

Finally, the result of the potential recognition stage is that agent  $a$  has the belief that it is possible to form a team to achieve the overall goal or that it is not, or formally it has established either

$$\exists G \subseteq TBEL(a, C-CAN_G(\varphi) \wedge \forall i \in G \text{ willing}(i, \varphi))$$

or its negation. In some applications, like scientific research projects, the outcome of this stage is communicated to the agents who replied positively. In the context of services like travel agencies this may not be necessary.

## 7 Team formation

At the stage of potential recognition, individual properties of agents were considered. These could play a role in different types of cooperation, where some services are exchanged but there is no interest in other agent's activities, and more specifically, there is no collective intention to achieve a goal as a team. At the stage of team formation, however, the conditions needed for teamwork in a strict sense are created. The main condition of teamwork is the presence of a collective intention, as defined in subsection 3.2.

Note that this concept of teamwork requires agents that have a type of "social conscience". We do not consider a set of agents as a team if they cooperate by just achieving their own predefined part of a common goal. If agents are part of a team they should be interested in the performance of the other team members and willing to adjust their task on the basis of the needs of others.

At the beginning, the initiator has a sequence of groups in mind that could be suitable to form a team for achieving the goal. Although we do not go into details here, the organization of the groups into a sequence is important for the case that the most desirable group can not be formed.

All agents in these potential teams have expressed their willingness to participate towards the overall goal, but do not necessarily have the individual intention to contribute towards it yet. In this situation, the initiator tries to persuade them to take on the intention to achieve the overall goal as well as to act as a team with the other members. Note that task division and task allocation only come to the fore at the next stage, plan formation.

## 7.1 Persuasion dialogue

The goal of a persuasion dialogue is to establish a collective intention within a group  $G$  to reach the overall goal  $\varphi$  ( $C\text{-INT}_G(\varphi)$ , see subsection 3.2). Axiom **M3** makes evident that a crucial step for the initiator is to persuade all members of a potential team to take the overall goal as an individual intention. To establish the higher levels of the mutual intention, the initiator also persuades each member to take on the intention that all members of the potential team have the mutual intention, in order to strengthen cooperation from the start. It suffices if the initiator persuades all members of a potential team  $G$  to take on an individual intention towards  $\varphi$  ( $\text{INT}(i, \varphi)$ ) and the intention that there be a mutual intention among that team ( $\text{INT}(i, \text{M-INT}_G(\varphi))$ ). This results in  $\text{INT}(i, \varphi \wedge \text{M-INT}_G(\varphi))$  for all  $i \in G$ , or equivalently by axiom **M1**:  $\text{E-INT}_G(\varphi \wedge \text{M-INT}_G(\varphi))$ , which in turn implies by axiom **M2** that  $\text{M-INT}_G(\varphi)$ . When all the individual motivational attitudes are established within the team, the initiator broadcasts the fact  $\text{M-INT}_G(\varphi)$ , by which the necessary collective belief  $C\text{-BEL}_G(\text{M-INT}_G(\varphi))$  is established and the collective intention is in place.

This will be achieved during a persuasion dialogue, which according to [30] consists of three main stages: information exchange, rigorous persuasion and completion. In our case the information exchange already started in the potential recognition stage. Let us remind the reader that we extended the concept of persuasion to include also intentions, and not only beliefs. The final result of the team formation stage is reached when for one potential team all the persuasion dialogues have been concluded successfully.

**Information exchange** During the information exchange the agents make clear their initial stand with respect to the overall goal and to their being part of a certain team to achieve it. These issues are expressed partly in the form of intentions and beliefs. Other beliefs supporting or related to the above issues might also be exchanged already. Only when a conflict arises about these issues a persuasion dialogue has to take place. In each persuasion there are two parties or roles; the proponent (P) and the opponent (O). In our case the proponent is the initiator and the opponent the other agent.

The stand the other agent takes about the above issues are seen as its initial *concessions*. Concessions are beliefs and intentions that an agent takes on for the sake of argument, but need not be prepared to defend. The agents will also have private attitudes that may only become apparent later on during the dialogue. The stand of the initiator is determined by the initial thesis that it is prepared to defend during the dialogue. The initial conflict description consists of the set of O's initial concessions and of P's initial thesis.

**Rigorous persuasion** During the rigorous persuasion stage the agents exchange arguments to challenge or support a thesis. The following rules can be used to govern these moves, adapted from [30]:

1. Starting with O the two parties move alternately according to the rules of the game.
2. Each move consists of either a challenge, a question, a statement, a challenge or question accompanied by a statement (see [30]), or a final remark.

3. The game is highly asymmetrical. All P's statements are assertions, and called *theses*, all O's statements are called *concessions*. P is doing all the questioning and O all the challenging.
4. The initial move by O challenges P's initial thesis. It is P's goal to make O concede the thesis. P can do this by questioning O and thus bridge the gap between the initial concessions of O and the thesis, or by making an assertion to clinch the argument if acceptable or defensible in further dialogue.
5. Each move for O is to pertain to P's preceding move. If this was a question, then O has to answer it. If it was an assertion, then O has to challenge it, unless O gives up, see 6).
6. Each party may give up, using the final remark  $ASS_{a,i}(quit)$  for the initiator P, or  $ASS_{i,a}(INT(i, \varphi \wedge M-INT_G(\varphi)))$  for the other agent O. If O's concessions include P's thesis, then P can end the dialogue by the final remark:  $ASS_{a,i}(won)$ . In our system we assume to have the following rule:  $[ASS_{a,i}(won)]OBL(ASS_{i,a}(INT(i, \varphi \wedge M-INT_G(\varphi))))$ , which means that agent  $i$  is obliged to state that it has been persuaded and accepts its role in the team. This does not mean that  $i$  will actually make this assertion! Just that there is an obligation (according to the rules of the persuasion "game"). The modal operator OBL is taken from deontic logic [1].
7. All challenges have to follow logical rules. For example, a thesis  $A \wedge B$  can be challenged by challenging one of the two conjuncts. For a complete set of rules for the propositional part of the logic we refer to [30].

In the completion stage the outcome is made explicit, such that the agents either have a collective belief and/or intention or they know that they differ in opinion.

**Speech acts during persuasion** In contrast to [30], we need to monitor the agent's informational and motivational attitudes during persuasion. We are concerned with assertions and challenges (wrt. informational attitudes), and concessions and requests (wrt. both informational and motivational attitudes).

As for assertions, after a speech act of the form  $ASS_{a,i}(B)$ , agent  $i$  believes that the initiator believes that  $B$ :

$$[ASS_{a,i}(B)]BEL(i, BEL(a, B))$$

Let us assume that  $i$  has only two rules for answering an assertion  $B$ . If  $i$  does not have a belief that is inconsistent with  $B$  then  $i$  will concede (similarly as in default logic). If, on the other hand,  $i$  does have a belief to the contrary it will challenge the assertion. Formally:

$$\neg BEL(i, \neg B) \rightarrow DO(i, CONCEDE_{i,a}(B))$$

$$BEL(i, \neg B) \rightarrow DO(i, CHALLENGE_{i,a}(B))$$

where the operator  $DO(i, \alpha)$  indicates that  $\alpha$  is the next action performed by  $i$ .

The CONCEDE action with respect to informational attitudes is basically an assertion plus a possible mental update of the agent. Thus, it does not only assert the

proposition but actually believes it as well, even if it did not believe the proposition beforehand. Suppose that  $i$  did not have a contrary belief, then  $i$  concedes  $B$  by the speech act  $\text{CONCEDE}_{i,a}(B)$ . The effect of this speech act is similar to that of  $\text{ASS}$ , except that  $a$  can only assume that  $i$  believes the formula  $B$  during the dialogue and might retract it afterwards.

$$[\text{CONCEDE}_{i,a}(B)]\text{BEL}(a, \text{BEL}(i, B))$$

The **CHALLENGE** with respect to informational attitudes, on the other hand, is a combination of a denial (assertion of a belief in the negation of the proposition) and a request to prove the proposition. The exact form of the challenge depends on the logical form of the assertion [30]. For this reason, the complete effects of this speech act are quite complex to describe fully. We will give an example of a challenge in the next subsection.

With respect to motivational attitudes, the situation is different. For example, initiator  $a$  requests  $i$  to take on an intention  $\psi$  by the following speech act:

$$\text{REQ}_{a,i}(\text{CONCEDE}_{i,a}(\text{INT}(i, \psi))).$$

Again,  $i$  has only two rules for answering such a request. If  $i$  does not have an intention  $\neg\psi$  (that is inconsistent with  $\psi$ ) then  $i$  will concede. If, on the other hand  $i$  does have an intention to the contrary it will assert that it intends  $\neg\psi$ :

$$\neg\text{INT}(i, \neg\psi) \rightarrow \text{DO}(i, \text{CONCEDE}_{i,a}(\text{INT}(i, \psi)))$$

$$\text{INT}(i, \neg\psi) \rightarrow \text{DO}(i, \text{ASS}_{i,a}(\text{INT}(i, \neg\psi)))$$

For example, suppose that  $i$  did not have a contrary intention then  $i$  concedes by the speech act  $\text{CONCEDE}_{i,a}(\text{INT}(i, \psi))$ . The effect of this speech act is:

$$[\text{CONCEDE}_{i,a}(\text{INT}(i, \psi))]\text{BEL}(a, \text{INT}(i, \psi))$$

## 7.2 Team formation for the example

The running example in this paper is about team formation for achieving the following overall goal (further abbreviated as  $\varphi$ ): “to arrange a trip of three weeks to Australia for a certain famous family; the trip should satisfy (specific) constraints on costs, times, places and activities”. The initiative for teamwork is taken by travel agent  $a$ , who cannot arrange the whole trip on his own. The trip will be extensively publicized, so it has to be a success, even if circumstances change. Thus, he does not simply ask airline companies, hotels, and organizers of activities to deliver a number of fixed services. Instead, he believes that real teamwork, where all members are interested in the achievement of the overall goal, gives the best chances of a successful trip. This paper discusses the first two stages of cooperative problem solving only; plan formation and team action for this example will be described in further work.

In the travel example, the initiator tries to persuade the other agents  $i$  in the potential team to take on the intention to achieve the overall goal of organizing the journey ( $\text{INT}(i, \varphi)$ ), but also with respect to doing this as a team with the other agents

(INT( $i, M\text{-INT}_G(\varphi)$ )). To this end, the initiator exploits the theory of intention formation.

Intentions are formed on the basis of beliefs and previously formed intentions of a higher abstraction level by a number of formal rules (see [5]). For example, the built-in intention can be to obey the law, or avoid punishment. The (instrumental) belief is that driving slower than the speed limit is instrumental for obeying the law, and is its preferred way to do so. Together with the rule the new intention of driving slower than the speed limit is derived.

The general intention generation rule is represented as follows:

$$\mathbf{IG} \text{ INT}(i, \psi) \wedge \text{BEL}(i, \text{INSTR}(i, \chi, \psi)) \wedge \text{PREFER}(i, \chi, \psi) \rightarrow \text{INT}(i, \chi)$$

It states that if an agent  $i$  has an intention  $\psi$  and it believes that  $\chi$  is instrumental in achieving  $\psi$  and  $\chi$  is its preferred way of achieving  $\psi$ , then it will have the intention to achieve  $\chi$ . “ $\chi$  is instrumental in achieving  $\psi$ ” means that achieving  $\chi$  gets the agent “closer” to  $\psi$  in some abstract sense. We do not define this relation any further, but leave it as primitive.

The PREFER relation is based on an agent’s individual beliefs about the utility ordering between its goals, collected here into a finite set  $H$ . We abstract from the specific way in which the agent may compute the relative utilities, but see the literature about (qualitative) decision theory [3].

$$\begin{aligned} \text{PREFER}(i, \chi, \psi) \leftrightarrow \\ \bigwedge_{\xi \in H} (\text{BEL}(i, \text{INSTR}(i, \xi, \psi)) \rightarrow \text{BEL}(i, ut(i, \chi) \geq ut(i, \xi))) \end{aligned}$$

The mechanism sketched in subsection 7.1 can be used in our setting during persuasion. In our example, the initiator  $a$  tries to get the other agent  $i$  to concede to higher level intentions, instrumental beliefs and preferences that together with **IG** imply the intention to achieve the overall goal  $\varphi$ . To be more concrete, we could choose the higher level intention  $\psi$  to stand for “earn good money”. Here follows an example move of the initiator:

$$\text{ASS}_{a,i}(\forall j(\text{INT}(j, \psi) \rightarrow \text{INSTR}(j, \varphi, \psi))).$$

After this speech act agent  $i$  believes that the initiator believes that if an agent has the higher level intention to earn good money, then the overall intention  $\varphi$  is instrumental to this. Formally (see also [9]):

$$\begin{aligned} [\text{ASS}_{a,i}(\forall j(\text{INT}(j, \psi) \rightarrow \text{INSTR}(j, \varphi, \psi)))] \\ \text{BEL}(i, \text{BEL}(a, \forall j(\text{INT}(j, \psi) \rightarrow \text{INSTR}(j, \varphi, \psi)))) \end{aligned}$$

According to the general rule about assertions there are two possibilities for  $i$ ’s answer. Let us assume that the positive case holds, i.e.  $i$  does not have a contrary belief, so it concedes:  $\text{CONCEDE}_{i,a}(\forall j(\text{INT}(j, \psi) \rightarrow \text{INSTR}(j, \varphi, \psi)))$ . The effect of this speech act on agent  $a$  is given by the general rule:

$$\begin{aligned} [\text{CONCEDE}_{i,a}(\forall j(\text{INT}(j, \psi) \rightarrow \text{INSTR}(j, \varphi, \psi)))] \\ \text{BEL}(a, \text{BEL}(i, \forall j(\text{INT}(j, \psi) \rightarrow \text{INSTR}(j, \varphi, \psi)))) \end{aligned}$$

Now the formula is believed by both  $a$  and  $i$ . Thus, the initiator's next aim in the persuasion will be to get  $i$  to intend  $\psi$  (earn good money) by the question:

$$\text{REQ}_{a,i}(\text{CONCEDE}_{i,a}(\text{INT}(i, \psi))).$$

By the general rule,  $i$  is obliged to either concede it has the intention  $\psi$  (if it is consistent with its other intentions) or to assert that it intends its negation. After  $i$ 's response, the initiator believes  $i$ 's answer. Note that in the second case, it may be useful for  $a$  to embed a negotiation dialogue in the persuasion, in order to get  $i$  to revise some of its previous intentions. For the example, let us suppose that  $a$  is successful in persuading  $i$ .

When the initiator has persuaded agent  $i$  to take on the high level intention  $\psi$  and to believe the instrumentality of  $\varphi$  with respect to  $\psi$ , it can go on to persuade the other that  $\text{PREFER}(i, \phi, \psi)$  by the speech act:

$$\text{ASS}_{a,i}(\bigwedge_{\xi \in H} (\text{BEL}(i, \text{INSTR}(i, \xi, \psi)) \rightarrow \text{BEL}(i, \text{ut}(i, \varphi) \geq \text{ut}(i, \xi))))$$

To make the example more interesting, let us suppose that  $i$  does not yet prefer  $\varphi$  as a means to earn good money; instead it believes that  $\chi$ , arranging some less complex holidays for another family, has a higher utility than  $\varphi$ . Thus  $i$  does not concede to the initiator's speech act, but instead counters with a challenge. According to the logical structure of the definition of  $\text{PREFER}$ , this challenge is a complex speech act consisting of three consecutive steps. First  $i$  asserts the negation of  $a$ 's assertion, a conjunction of implications; then it concedes to the antecedent of the implication for a specific goal  $\chi \in H$ ; and finally it requests  $a$  to present a proof that  $\varphi$  has a better utility for  $i$  than  $\chi$ .

$$\begin{aligned} & \text{CHALLENGE}_{i,a} \\ & (\bigwedge_{\xi \in H} (\text{BEL}(i, \text{INSTR}(i, \xi, \psi)) \rightarrow \text{BEL}(i, \text{ut}(i, \varphi) \geq \text{ut}(i, \xi)))) \equiv \\ & \text{ASS}_{i,a}(\neg(\bigwedge_{\xi \in H} (\text{BEL}(i, \text{INSTR}(i, \xi, \psi)) \rightarrow \text{BEL}(i, \text{ut}(i, \varphi) \geq \text{ut}(i, \xi))))); \\ & \text{CONCEDE}_{i,a}(\text{BEL}(i, \text{INSTR}(i, \chi, \psi))); \\ & \text{REQ}_{i,a}(\text{ASS}_{a,i}(\text{PROOF}(\text{ut}(i, \varphi) \geq \text{ut}(i, \chi)))) \end{aligned}$$

As a reply,  $a$  could prove that the utility of  $\varphi$  is in fact higher than that of  $\chi$ , because it generates a lot of good publicity, which will be profitable for  $i$  in future – something of which  $i$  was not yet aware. Let us suppose that  $i$  is persuaded by the proof and indeed concedes to its new preference by the speech act:

$$\text{CONCEDE}_{i,a}(\text{PREFER}(i, \varphi, \psi)).$$

All these concessions, together with the general intention formation rule and the fact that agents are correct about their intentions, then lead to  $\text{INT}(i, \varphi)$ . For intentions with respect to cooperation with other potential team members, the process to persuade the agent to take on  $\text{INT}(i, \text{M-INT}_G(\varphi))$  is analogous.

## 8 Discussion and conclusions

In previous work [14] it was shown how all four stages of CPS result in specific motivational attitudes that can be formally described. In this paper, we have shown for the first

two stages, potential recognition and team formation, which rules govern the dialogues by which the motivational attitudes are formed, and how to represent the moves within the dialogues by formalized speech acts.

It is clear that, even though the dialogues are governed by strict rules, the reasoning needed to find an appropriate move is highly complex. This implies that the agents also have to contain complex reasoning mechanisms in order to execute the dialogues. It means that, although the result is much more flexible and refined than using a protocol like Contract Net, the process is also more time consuming. For practical cases one should carefully consider what carries more weight and choose the method of team formation accordingly.

Related work can be found in [23], who also present an agent architecture and a representation for agent communication. In the discussion they note that their own “negotiation” in fact covers a number of Walton and Krabbe’s different dialogue types. We find the more fine-grained typology to be very useful when designing agents for teamwork: one can use specific sets of rules governing each type of dialogue as well as the possible embeddings between the different types. Thus desired kinds of communication are allowed and harmful ones prevented, without completely fixing any protocol. Also [23] uses multi-context logic whereas we stick to (multi-)modal logic.

The emphasis on pre-planning (here in the stage of potential recognition) and establishing appropriate collective attitudes for teamwork is shared with Grosz and Kraus [20]. Nevertheless, the intentional component in their definition of collective plans is much weaker than our collective intention: Grosz and Kraus’ agents involved in a collective plan have individual intentions towards the overall goal and a collective belief about these intentions; intentions with respect to the other agents play a part only at the level of individual sub-actions of the collective plan. We stress, however, that team members’ intentions about their colleagues’ motivation to achieve the overall goal play an important role in keeping the team on track even if their plan has to be changed radically due to a changing environment (see also [14]).

The first issue for further research is to give a complete set of formal rules for all the types of dialogue and indicate how these are implemented through formal speech acts. This would make it possible to extend the framework to the next stages of cooperative problem solving, namely plan formation and execution.

A second issue is the further investigation of several aspects of the internal reasoning of the agents. One example is the concept of giving a proof as defence of an assertion during the rigorous persuasion. Finally, it should be investigated how actual proofs can be constructed in an efficient way to show that the end results of a dialogue are formed through the speech acts given the rules of the dialogue.

## **Acknowledgments**

We would like to thank Erik Krabbe, Alexandru Baltag, Magnus Boman, Wiebe van der Hoek, Chris Reed, and the anonymous referees for their useful comments. Barbara Dunin-Kęplicz’ research was supported by ESPRIT under the Grant CRIT-2 No. 20288 and KBN Grant 8T11C02519.

## References

1. L. Aqvist. Deontic logic. In: D. Gabbay and F. Guentner (eds.), *Handbook of Philosophical Logic*, Vol. III, Reidel, Dordrecht, 1984, pp. 605–714.
2. N. Belnap and M. Perloff. Seeing to it that: a canonical form for agentives. In: *Theoria*, vol. 54, 1988, pp. 175-199.
3. C. Boutilier. Toward a logic for qualitative decision theory. In: *Proceedings KR'94*, 1994, pp. 75-86.
4. C. Castelfranchi and Y.-H. Tan (eds.). *Trust and deception in virtual societies*, Kluwer, Dordrecht, 2000.
5. F. Dignum and R. Conte. Intentional agents and goal formation. In: M. Singh et. al.(eds.), *Intelligent Agents IV (LNAI 1365)*, Springer Verlag, 1998, pp. 231–244.
6. F. Dignum, B. Dunin-Kęplicz, and R. Verbrugge. Dialogue in team formation: a formal approach. In: F. Dignum and B. Chaib-draa (eds.), *IJCAI Workshop on Agent Communication Languages*, Stockholm, 1999, pp. 39–50.
7. F. Dignum, B. Dunin-Kęplicz, and R. Verbrugge. Agent theory for team formation by dialogue. In: *The VII Workshop on Agent Theories, Architectures and Languages*, Boston, 2000.
8. F. Dignum and R. Kuiper. Combining dynamic deontic logic and temporal logic for the specification of deadlines. In: Jr. R. Sprague (ed.), *Proceedings of thirtieth HICSS*, Wailea, Hawaii, 1997.
9. F. Dignum and H. Weigand. Communication and deontic logic. In: R. Wieringa and R. Feenstra (eds.), *Information Systems, Correctness and Reusability*, World Scientific, Singapore, 1995, pp. 242–260.
10. B. Dunin-Kęplicz and A. Radzikowska. Actions with typical effects: epistemic characterization of scenarios. In: V. Lesser (ed.), *Proc. First International Conference on Multi-Agent Systems, ICMAS'95*, MIT Press and AAAI Press, Cambridge (MA), 1995, p. 445.
11. B. Dunin-Kęplicz and A. Radzikowska. Epistemic approach to actions with typical effects. In: *Proceedings ECSQARU'95*, LNAI 946, 1995, pp. 180-189.
12. B. Dunin-Kęplicz and A. Radzikowska. Nondeterministic actions with typical effects: reasoning about scenarios. In: *Formal Models of Agents*, LNAI 1760, 1999, pp. 143-156.
13. B. Dunin-Kęplicz and R. Verbrugge. Collective commitments. In: *Proc. Second International Conference on Multi-Agent Systems, ICMAS'96*, IEEE Computer Society Press, Kyoto, 1996, pp. 56–63.
14. B. Dunin-Kęplicz and R. Verbrugge. A Reconfiguration algorithm for distributed problem solving. In: *Journal of Electronic Modeling*, vol. 22, nr 2, 2000, pp. 68–86.
15. B. Dunin-Kęplicz and R. Verbrugge. Collective motivational attitudes in cooperative problem solving. In: V. Gorodetsky et al. (eds.), *Proceedings of The First International Workshop of Central and Eastern Europe on Multi-Agent Systems (CEEMAS'99)*, St. Petersburg, 1999, pp. 22–41.
16. B. Dunin-Kęplicz and R. Verbrugge. The role of dialogue in Cooperative Problem Solving. To appear in: *Proceedings of the Workshop On Intelligent Agents for Computer Supported Cooperative Work: Technology and Risks*, Barcelona, 2000, pp.1–8.
17. R. Fagin, J.Y. Halpern, Y. Moses, and M.Y. Vardi. *Reasoning about Knowledge*. MIT Press, Cambridge (MA), 1995.
18. M.J. Fischer and R.E. Ladner. Propositional dynamic logic of regular programs, *Journal of Computer System Sci.* 18, 1979, pp. 194-211.
19. R. Goldblatt. *Logics of Time and Computation*. Stanford, CSLI Press, 1992.
20. B.J. Grosz and S. Kraus. Collaborative plans for group action. *Artificial Intelligence* 86, 1996, pp. 269-357.

21. B. van Linder, W. van der Hoek, and J.-J. Ch. Meyer. Formalising abilities and opportunities of agents. *Fundamenta Informaticae* 34, 1998, pp. 53-101.
22. B. van Linder, W. van der Hoek, and J.J.Ch. Meyer. Actions that make you change your mind. In: , A. Laux and H. Wansing (eds), *Knowledge and Belief in Philosophy and Artificial Intelligence*, Akademie Verlag, 1995, pp. 103-146.
23. S. Parsons, C. Sierra, and N. Jennings. Agents that reason and negotiate by arguing. *Journal of Logic and Computation*, 8(3) (1998) pp. 261-292.
24. T. Sandholm and V. Lesser. Issues in automated negotiation and electronic commerce: extending the contract net protocol. In: *Proceedings First International Conference on Multiagent Systems (ICMAS95)*, San Francisco, AAAI Press and MIT Press, 1995, pp. 328-335.
25. A.S. Rao and M.P. Georgeff. Modeling rational agents within a BDI architecture. In: R. Fikes and E. Sandewall (eds.), *Proceedings of Knowledge Representation and Reasoning (KR&R-91)*, San Mateo, Morgan Kaufmann, 1991, pp.473-484.
26. J. R. Searle. *Speech Acts: An Essay in the Philosophy of Language*. Cambridge, Cambridge University Press, 1969.
27. J.R. Searle and D. Vanderveken *Foundations of Illocutionary Logic*, Cambridge, Cambridge University Press, 1985.
28. K. Segerberg. Bringing it about. In: *Journal of Philosophical Logic*, vol. 18, 1989, pp. 327-347.
29. D.R. Traum. Speech acts for dialogue agents. In: M. Wooldridge and A. Rao (eds.), *Foundations of Rational Agency*. Dordrecht, Kluwer, 1999, pp. 169-201.
30. D. Walton and E. Krabbe. *Commitment in Dialogue*, SUNY Press, Albany, 1995.
31. M. Wooldridge and N.R. Jennings. Cooperative Problem Solving. *Journal of Logic and Computation* 9 (4) (1999) pp. 563-592.