

## **UC Merced**

### **Proceedings of the Annual Meeting of the Cognitive Science Society**

#### **Title**

Strategic Social Planning - Looking for Willingness in Multi-Agent Domains

#### **Permalink**

<https://escholarship.org/uc/item/9n1401t2>

#### **Journal**

Proceedings of the Annual Meeting of the Cognitive Science Society, 15(0)

#### **Authors**

Miceli, Maria

Cesta, Amedeo

#### **Publication Date**

1993

Peer reviewed

# Strategic Social Planning Looking for Willingness in Multi-Agent Domains <sup>1</sup>

Maria Miceli and Amedeo Cesta  
IP-CNR  
National Research Council  
Viale Marx 15, I-00137 Rome, Italy  
*amedeo@irmkant.bitnet*

## Abstract

This paper deals with the use of social knowledge by an autonomous agent which is planning its behavior and in particular discovering it is in need of help. What we aim at showing is the possible insertion of knowledge about dependence relations in an agent architecture so that it may achieve a cognitively plausible behavior. A number of basic criteria are designed to endow our agent architecture with the ability to generate choices about social interactions and requests. Particular attention is paid to the criteria for assessing others' willingness to give help, and to the interaction of these criteria with the agent's general attitudes and skills.

## Introduction

This paper deals with the use of social knowledge by an autonomous agent that plans its behavior. In previous papers (Castelfranchi, 1990; Castelfranchi, Miceli, & Cesta, 1992), cognitive theories and formal descriptions have been developed to provide the basis for human social behavior. Our attention has been focused on *power* relations and in particular *dependence* relations. Agent's network of dependence relations is seen as a basic source of knowledge about who has the "power" needed, and what this power is for.

What we aim at showing is the possible insertion of our model of dependence in an agent architecture so that an autonomous agent may achieve a cognitively plausible behavior. Specifically, we focus on the process of planning to act on other agents in order for them to perform some needed action.

Dependence relationships among agents are viewed as a powerful interactive tool for both rational interaction and problem solving, and as a way for providing solutions to the problem of interaction and communication control. A cognitive agent who is unable to achieve a goal does not resort to somebody

else for help, nor does it limit itself to applying standard interaction protocols. On the contrary, it is likely to reason about its knowledge of dependence relations, and then act according by.

A problem-driven approach to social interaction is typical of this perspective. Here sociality (or some of its relevant aspects -- namely benevolence, or common goals) is neither given for granted nor traced back to varying degrees of overlapping among the agents' mental attitudes (which is the dominant view in Artificial Intelligence; see Grosz & Sidner, 1990; Pollack, 1990; Cohen & Levesque, 1990). On the contrary, the problem of why and under which conditions social interaction occurs is directly justified by the agents' problem solving: others are considered when help for some of the agent's problems is needed.

In Castelfranchi, Miceli, & Cesta (1992) we provided a number of basic definitions and principles for deriving a dependence relationship from previous relationships. Furthermore we focused on the relation between dependence and influencing, showing how knowledge about *x*'s dependence on *y* is likely to modify *x*'s mental state, by producing a new goal: i.e. influencing *y* to do what *x* needs. However, the goal of influencing can be a necessary but not sufficient condition for an actual influencing behavior: it would be foolish to embark in influencing without any assurance that the addressee *is likely to be influenced*. This amounts to saying that: (a) agent *x* should come to know not only which other agents it depends on and for what, but also some other relevant information about the others' mental states, including their goals as well as their beliefs about *x*'s goals, about their own duties and roles, and about their own dependence relations; (b) an agent architecture should include some reasoning capabilities about others' mental states in order to select the agents that are actually "useful".

Identifying the "right" agent is a knowledge intensive planning activity. Two aspects are critical in such a planning: (a) the need to produce a social sub-plan by using well-founded heuristics that analyze the agent's social knowledge (i.e. its knowledge about others); (b) the need to avoid producing a plan which is in conflict with the agent's other goals (e.g. to avoid producing a direct request to another agent while one is committed to demonstrate one's own skillfulness).

---

<sup>1</sup>This work has been partially supported by CNR under "Progetto Finalizzato Sistemi Informatici e Calcolo Parallelo", Grant n.104385/69/9107197 to IP-CNR. Both authors participate in the "Project for the Simulation of Social Behavior" at IP-CNR.

In this paper we briefly mention the agent architecture we are working on, then present the particular problem of reasoning about others' willingness to help. We describe the different criteria applied to identify a willing agent, and we consider the possible conflicts between one's search for others' help and one's personal goals, which can favor individual biases or preferences toward one help-seeking strategy or another.

Although our work on this architecture is still in progress, we claim that, in order to build a social planning system, we have to endow it with some key knowledge about others, their beliefs, goals, and skills. This paper explores how such knowledge can become operative.

## Agent Architecture and Dependence Theory

The agent we are currently implementing has a behavior that is strongly driven by its goals. We named it *CogAgent* to emphasize our attempt to endow it with the basic cognitive tools we have described so far. This agent's architecture has two main components: a propositional knowledge representation service, named KRAM (D'Aloisi & Castelfranchi, 1993), and a planning service. Some of the goals *CogAgent* wants to achieve are passed on to the planning component, which then works to produce a plan for them. A number of specialized reasoners contribute to the creation of the plan that guides the agent's behavior: a *memory based planner* (Cesta & Romano, 1992) builds the initial plans for achieving the goals; once a plan is created, a *resource analysis and allocation module* is called into play to identify any problem for the future action execution. In particular, for each plan action, it tests whether it is included in the action repertoire, and whether all the resources implied by the action execution are accessible to the agent. After the resource analysis has identified the need for other resources, two modules are called into play. A *social knowledge specialist* is responsible for searching social resources. This consists in checking whether there are other agents in the common world which might give *x* some help to get the work done. A *social goals conflict analyzer* is responsible for scrutinizing social plans provided by the previous specialist, in order to check for negative interactions between such a plan and other goals.

### Dependence-Based Criteria for Searching "Useful" Agents

A basic definition in our theory, whose logical formalization is drawn from Cohen and Levesque's (1990), is that of *social dependence* (Castelfranchi, Miceli, & Cesta, 1992):

$$(S-DEP\ x\ y\ a\ p) = \text{def } (GOAL\ x\ p) \wedge \\ \neg (CANDO\ x\ a) \wedge (CANDO\ y\ a) \wedge \\ ((DONE-BY\ y\ a) \supset (EVENTUALLY\ p))$$

that is: *x* depends on *y* with regard to an act *a* useful for realizing a state *p* when *p* is a goal of *x*'s and *x* is unable to realize *p* while *y* is able to do so.

Now, our starting point is *x*'s *assumed* dependence on *y*. More exactly, since *x* does not necessarily know which *y* it is dependent upon, we will have:

$$(GOAL\ x\ p) \wedge (BEL\ x\ \neg (CANDO\ x\ a)) \wedge \\ (BEL\ x\ \exists y_i (CANDO\ y_i\ a)) \wedge \\ (BEL\ x\ \exists y_i ((DONE-BY\ y_i\ a) \supset \\ (EVENTUALLY\ p)))$$

Let us provide a very simple example. Let's assume that *x* and *y*<sub>1</sub>, *y*<sub>2</sub>, ..., *y*<sub>*n*</sub> are people working in a research center, and the e-mail server of the computer network is down. (From now on, we'll refer to *x* as to "he" and to *y* as to "she".) Agent *x* has the goal *p* to make the e-mail server work (in order to send a message), but he is unable to do the (set of) action(s) *a* required to activate the server, and he assumes he is unable to do *a*, and that in the department there are some *y*s (at least one) who are able to do *a*.

In order to achieve *p*, *CogAgent x* must resort to some dependence-based meta-goals or criteria in order to use the amount of knowledge he has about others, their abilities and so on. Two basic sets of criteria have been drawn from our theory of dependence: (a) the *CANDO* criterion, for finding or retrieving information about *y*'s ability to perform the required *a*, i.e., for satisfying the predicate (*CANDO y a*); (b) the *WILL* criteria, for finding or retrieving information about *y*'s willingness to perform *a*, i.e., for satisfying the predicate (*GOAL y (DONE-BY y a)*).

Here we focus on the willingness aspects. An exhaustive presentation of the *CANDO* aspect can be found in Cesta & Miceli (1993).

### The Search for Willingness

In the following we assume that *x* has already successfully applied the *CANDO* criterion, and found out the (set of) agent(s) *y* that (*CANDO y a*). So the problem *x* has to cope with is *y*'s willingness to perform *a*. We shall propose three *WILL* criteria, that is, three possible strategies to search for an agent who has the goal of performing the required action.

#### E-WILL Criterion: Exploitation-Seeking

The most straightforward *WILL* criterion is to see whether *y* already has the goal to perform *a*, in order to take advantage of her performance. Suppose John knows that Carol will need to use the e-mail and, being able to reset the server, is likely to do the job:

John will just wait for her to do so. Now, how can  $x$  assess whether  $y$  has ( $GOAL\ y\ (DONE-BY\ y\ a)$ )?

We can apply either a classification or a performance sub-criterion. In the former case,  $x$  can see whether  $y$  belongs to a class  $Y$  of agents with ( $GOAL\ Y\ (DONE-BY\ Y\ A)$ ). For instance, John can see whether Carol belongs to the class of (competent) users of the e-mail server. We can say:

$$(BEL\ x\ (((ISA\ a\ A) \wedge (CANDO\ Y\ A) \wedge (GOAL\ Y\ (DONE-BY\ Y\ A))) \supset ((ISA\ y\ Y) \supset (GOAL\ y\ (DONE-BY\ y\ a))))))$$

The performance sub-criterion is an empirical one:  $x$  should either see  $y$  performing  $a$  or carry out some plan recognition according to which  $x$  can believe  $y$  is going to perform  $a$ . Suppose John sees Carol going to the room where there is the e-mail server after her director's request: John can easily infer that she is going to perform  $a$ .

### B-WILL Criterion: Benevolence-Seeking

In [CAS91] we define benevolence as follows:

$$(BENEVOLENT\ y\ x\ p) =_{def} (BEL\ y\ (GOAL\ x\ p)) \supset (EVENTUALLY\ (GOAL\ y\ p))$$

Thus,  $y$  is benevolent toward  $x$  for the goal  $p$  if  $y$  believes that  $x$  has goal  $p$  then  $y$  will have the same goal  $p$ .

Benevolence can be either *individualized*, or personal, when  $y$  is likely to adopt the goals of a specific individual agent, and *non individualized*, when  $y$  (usually a given set of  $ys$ ) adopts the goals of a given set of  $xs$ . In the general form of non individualized benevolence three sets of objects are mentioned: the recipients, their goals, and the benevolent agents. So, the recipient may believe he belongs to a set  $X$  which receives benevolence from a set  $Y$  of agents with regard to a given set  $P$  of goals, and that if  $y$  belongs to that set, he will receive benevolence from  $y$ :

$$(BEL\ x\ (((BENEVOLENT\ Y\ X\ P) \wedge (ISA\ x\ X) \wedge (ISA\ p\ P)) \supset ((ISA\ y\ Y) \supset (BENEVOLENT\ y\ x\ p))))$$

Non individualized benevolence seems to have two general sources: (a) *Role-tasks*, when  $x$  is mentioned as a beneficiary in the role of some  $ys$ . In such cases,  $x$  is due benevolence, and  $y$  is expected to give help to  $x$  relative to some specified set of goals. In our example, Carol might be supposed by role to reset the e-mail facility. This implies that  $x$  is supposed to have knowledge about roles and role-tasks. He should then be able to see whether the action he is in need of matches with the role-task of some other (class of) agents, as well as if his own role repertoire includes him as a beneficiary of such a role-task. (b) *General norms*: sometimes norms control whether adoption should be given and to whom, independently of the

roles structure. In some contexts, some agents might be expected to adopt some goals of other agents (for instance, on the bus, people are expected to get up and leave their seats to the disabled). A special case is represented by the norm of reciprocation:  $y$  is supposed to reciprocate, that is, to adopt some goals of those agents who have intentionally given her benefit in the past without being held to do so. This implies that *CogAgent* will check his domain knowledge for any norm mentioning a set of beneficiaries and that he is endowed with a memory of past interactions.

### D-WILL Criterion: Dependence-Seeking

There is another *WILL* criterion which looks quite widespread and crucial in natural social interaction: the search for an agent who, besides being able to perform  $a$ , is in turn dependent on  $x$  for some other action. In other words,  $x$  should find out whether it is possible to move from assumed unilateral dependence of  $x$  on  $y$  to assumed bilateral dependence of each one on the other. Bilateral dependence, in fact, would obviously allow for  $x$ 's more powerful position: it is true that, on the grounds of his assumed dependence on  $y$ ,  $x$  has the goal of influencing  $y$  to perform  $a_1$ , it should be true that also  $y$ , on the grounds of her assumed dependence on  $x$ , should have the goal of influencing  $x$  to perform some other action  $a_2$ . So,  $x$  might offer his own performance of  $a_2$  in exchange for  $y$ 's performance of  $a_1$ . This offer of exchange is in fact a form of negotiation (Zlotkin & Rosenschein, 1992) based on assumed dependence. Going back to our example, suppose that while Carol is able to reset the e-mail server, John knows ancient Greek very well and he knows that Carol doesn't and is in need of some translation from English into ancient Greek.

As one can see, situations of this sort are at the base of a great variety of social interactions, from informal exchange of favors to the more formal and socially regulated bargainings. Without claiming here to provide a detailed analysis of bilateral dependence, we just try to give an idea of the search for bilateral dependence needed to influence  $y$  to perform  $a$ .

First of all,  $x$  should see whether in his own action repertoire there is some action, that is useful towards achieving some of  $y$ 's goals, that  $x$  can perform and that  $y$  is unable to carry out.

The usual question, then, is: How can  $x$  find out  $y$ 's dependence on himself? Both experience and categorial knowledge can help also in this case. Agent  $x$  might see (or have seen)  $y$  trying to perform  $a_2$  (or more generally to pursue  $q$ ) without success; or  $x$  may know  $y$  belongs to a class of agents who  $\neg (CANDO\ Y\ A_2)$ . For instance, Carol as a computer scientist is quite unlikely to have had a classical education including the study of ancient Greek. In such cases the  $x$ 's reasoning criterion would be something like:

$(BEL x (((ISA a_2 A_2) \wedge \neg (CANDO Y A_2)) \supset ((ISA y Y) \supset \neg (CANDO y a_2))))$

In fact, when in search of others' dependence, one must go back to assessing their power, namely their *lack of power*. So,  $x$  has to apply a *CANDO* criterion of a special nature, by looking for some  $\neg (CANDO y a_2)$  which must correspond to some  $(CANDO x a_2)$ . In addition, a special problem is represented by the commensurability of  $x$ 's and  $y$ 's respective dependencies. In other words,  $x$  can not resort to *any* lack of power of  $y$ 's as a basis for his offer of exchange:  $y$ 's and  $x$ 's need for each other's help must be commensurable, as well as the costs of their respective performances. Such aspects, including many other "quantitative" features of the dependence relationships (for instance, the *degree* of dependence of one agent on another; see Castelfranchi, Miceli & Cesta, 1992) would deserve a more careful examination.

**The Other's Dependence as a Basis for Influencing Behavior.** Once assessed that  $y$  depends on  $x$  for some action  $a_2$  in view of a given goal  $q$ ,  $x$  is left with the task of influencing  $y$  to perform the action  $a_1$  in view of his own goal  $p$ . Let us try a very rough sketch of the basic steps of such an influencing behavior or, more exactly, its basic requirements. These are in fact the very reasons why the search for a dependent  $y$  turns into the search for a "useful"  $y$ , i.e. for an agent who is likely to do what  $x$  wants to.

In order to profit by  $y$ 's dependence, agent  $x$  needs for  $y$  to believe two fundamental facts: (a) First of all, quite obviously,  $y$ 's dependence on  $x$ . In fact, until  $y$  does not know about her dependence on  $x$ ,  $y$  is not likely to look for  $x$ 's help, and  $x$  has nothing to offer in exchange for his own request. If Carol does not know that John is an expert in ancient Greek, she does not know he can help her; (b) However, there is a further means-end relationship  $y$  should assume in order for  $x$  to obtain what he needs: the means-end relationship between  $y$ 's doing  $a_1$  and  $x$ 's doing  $a_2$ . Such a relationship is special in nature, in that there might be no intrinsic reasons why  $y$ 's doing  $a_1$  should be a means for  $x$ 's doing  $a_2$ . Going back to the previous example, while  $a_2$  (doing the translation from ancient Greek into English) is a means for  $q$  (to pass an exam), and  $a_1$  (resetting the e-mail server) is a means for  $p$  (using the e-mail), why should  $a_2$  be a means for  $p$ , or  $a_1$  a means for  $q$ ?

This means-end relation is a sort of artifact, a social construction on  $x$ 's part. It is John who actually creates the relation between his doing the translation and Carol's resetting the e-mail server. That is the very reason why  $x$  needs some sort of "persuasive" power over  $y$ , to make  $y$  believe that such a relation is likely to hold (Castelfranchi, Miceli & Cesta, 1992). Generally speaking such a persuasive action consists of

a communicative act, be it a more or less explicitly stated promise ( $x$  warrants he will perform  $a_2$  if  $y$  performs  $a_1$ ) or a threat ( $x$  threatens he will not perform  $a_2$  if  $y$  does not perform  $a_1$ ; or  $x$  threatens he will perform some other action  $a_3$  which will thwart some of  $y$ 's goals).

## **WILL Criteria and CogAgent's Biases**

We do not assume the *WILL* criteria are organized into a fixed sequence. *CogAgent*'s choice to prefer one criterion over another may depend on a number of factors, and to their relative weights: (a) First, the agent must consider the general information about the probabilities of success of a given criterion in a given context. For instance, a context may be given (say, a very competitive one) where benevolent behavior is quite unlikely to occur. (b) Second, *CogAgent*'s success also depends on his own skills in applying one criterion or another. This implies that the *CogAgent*'s *CANDOs* repertoire includes *social CANDOs* such as more or less detailed influencing strategies. Not surprisingly, then, there may exist very smart exploiters who are not skilled as benevolence-seekers and vice versa. (c) Third, *CogAgent* may have a personal bias or preference toward a specific criterion. Such a bias may feel the effects of the previous factors (i.e. *CogAgent*'s beliefs about either the effectiveness of a given criterion or his own ability to apply it). However, personal biases may also stem from other facts, namely other personal goals *CogAgent* wants to achieve or to defend while in search of  $y$ 's willingness to perform  $a$ . Such goals -- that we are going to outline -- may happen to be made easier by applying one criterion and/or to be threatened by applying another.

The possibility to devise different *CogAgent* "characters" starting from different interactions among the agent's goals while planning (e.g. interaction of current goals with more personal goals and consequent biases) is an important feature in our experiment. In the following we give a first account of our analysis.

## **Goal Interaction May Generate Different "Characters"**

It is worth observing that all possible characters share at least a couple of goals, due to their common condition of dependence on  $y$ :  $(GOAL x p)$  that causes *CogAgent*  $x$  to depend on  $y$ ; and  $(GOAL x (DONE-BY y a))$  that is derived from  $x$ 's assumed dependence on  $y$ . In fact, as showed in Castelfranchi, Miceli, & Cesta (1992):

$(BEL x (S-DEP x y a p)) \supset (GOAL x (DONE-BY y a))$

Moreover, since according to a postulate of rational agenthood, in order to perform an action an agent must *want* that action, one can derive that:

$(BEL\ x\ (S-DEP\ x\ y\ a\ p)) \supset$   
 $(GOAL\ x\ (GOAL\ y\ (DONE-BY\ y\ a)))$

This goal justifies the application of the *WILL* criteria.

Now, the goals  $x$  derives his assumed dependence from may come into conflict with the *other*  $x$ 's goals. More exactly, a possible plan useful for achieving  $p$  through  $(GOAL\ y\ (DONE-BY\ y\ a))$  and stemming from the application of a certain *WILL* criterion may happen to thwart some other personal goal  $q$  of  $x$ 's. In such a case,  $x$  can select some alternative plan where possible -- by applying some alternative *WILL* criterion -- which, while achieving  $(GOAL\ y\ (DONE-BY\ y\ a))$ , does not entail that  $q$  is at risk.

In the following, we will just mention a few goals which can bias  $x$  toward one *WILL* criterion or the other. This, in our view, gives an idea of the flexibility and richness needed in an agent architecture to allow for possible (individual as well as contextual) differences in goals, preferences, and sub-plans when planning for the same end-goal, i.e. to obtain others' help.

**The Exploiter.** One of the simplest ways to try and obtain for someone to do something is to let her know that you need it. This is exactly what an Exploiter wants to *avoid*. In fact, (one of) the Exploiter's goal(s) is:

$(GOAL\ x\ \neg\ (BEL\ y\ (S-DEP\ x\ y\ a\ p)))$

This may be due to a number of possible supergoals. For instance,  $x$  may want to avoid the obligation to reciprocate: if in fact  $y$  knows that  $x$  needs her help and she performs the required action, she may expect  $x$ 's help when she is in need of it. Or,  $x$  may want to save his face before  $y$  or others: to admit one's own dependence is to declare some lack of power or helplessness, and an agent who is interested in defending his image of competence and skillfulness is also interested in hiding his dependencies on others.

So, in any case,  $x$ 's goal  $q$  is that  $y$  does not come to believe that he depends on her for doing  $a$  in view of  $p$ . But, what can  $x$  do to obtain both  $p$  and  $q$ ? At least two possible scenarios are given.

If  $y$  is already performing  $a$  or is going to do so for her own personal reasons, then  $x$  just waits for this to happen, and can hide both his goal  $p$  and his consequent goal that  $(DONE-BY\ y\ a)$ .

In a different scenario, this lucky coincidence may not occur. Here is where some room is left to  $x$  to apply some influencing strategies in order to "produce"  $(GOAL\ y\ (DONE-BY\ y\ a))$ . One possible strategy is the following.

While continuing to hide  $(GOAL\ x\ p)$ ,  $x$  can declare  $(GOAL\ x\ (GOAL\ y\ (DONE-BY\ y\ a)))$ , that is, he can let  $y$  know that he wants for her to do  $a$ . What he must avoid is for  $y$  to infer  $x$ 's dependence on her from this. In fact, since

$(BEL\ x\ (S-DEP\ x\ y\ a\ p)) \supset$   
 $(GOAL\ x\ (GOAL\ y\ (DONE-BY\ y\ a)))$

$y$ , coming to believe  $(GOAL\ x\ (GOAL\ y\ (DONE-BY\ y\ a)))$ , might infer that such a goal of  $x$ 's is a consequence of his dependence on her; so,  $y$  may come to believe  $(BEL\ y\ (S-DEP\ x\ y\ a\ p))$  which would thwart the Exploiter's goal  $q$ .

In order to avoid this unwelcome consequence,  $x$  may try to persuade  $y$  that:  $a$  is a means for  $r$ , which is one of  $y$ 's goals; he is benevolent toward  $y$ , hence, knowing that  $y$  has the goal to achieve  $r$ , he also has the goal for her to achieve  $r$ ; for this reason, he also wants  $y$  to have the goal of doing  $a$  in view of  $r$ . In sum,  $x$  should try to make  $y$  believe that he wants her to do the required action *in her own interest*, in order to avoid for her to infer that the action would be in his interest instead.

**The Benevolence-Seeker.** Suppose  $x$  does not have the goal of hiding his dependence on  $y$  from  $y$  herself (because, say, either contextually or generally he is not interested in defending his image of competence and skillfulness). Conversely, he has the goal  $q$  to save a certain kind of resources, namely those implied by an offer of exchange (i.e., the influencing behavior required by an offer of exchange, and the performance of  $a_2$  in exchange for  $y$ 's performance of  $a_1$ ).

In this case,  $x$  will look for some benevolent  $y$ . In order to obtain benevolence, he will actively pursue a goal which is the exact opposite of the Exploiter's: i.e. for  $y$  to come to believe that  $x$  depends on her. In fact, according to the previous definition of benevolence, in order to "benevolently" behave,  $y$  should believe that  $(GOAL\ x\ p)$  and, plausibly, in order to actually perform some action in view of  $p$ , she must also believe that  $x$  cannot do  $a$  by himself. So, the Benevolence-Seeker must actively pursue  $(GOAL\ x\ (BEL\ y\ (S-DEP\ x\ y\ a\ p)))$ .

In view of this goal,  $x$  can show and even exaggerate his dependence, employing some influencing strategies aimed at pointing to his lack of power (e.g. by presenting himself as a needy and helpless agent).

However, it is worth observing that  $y$ 's benevolent behavior calls for some reciprocation: because of her benevolence,  $y$  is entitled by a norm of reciprocation to expect something from  $x$ . So,  $x$  cannot totally achieve his goal  $q$  of saving resources while obtaining  $y$ 's help for  $p$ . But, unlike the exchange offer, benevolence-seeking implies a form of reciprocation which is quite free from constraints: on the one hand, it should occur "sometime in the future", i.e. there are no specific temporal constraints; on the other, and even more importantly, reciprocation might coincide with a pure display of gratitude or an acknowledgement of subjection on  $x$ 's part, independent of the employment of material resources.

**The Bilateral Dependence-Seeker.** The Bilateral Dependence-Seeker (or "Exchanger", for the sake of brevity) is an agent who wants to offer something in exchange for the help he needs. Why choosing to spend resources if there might be other less expensive plans available, such as exploitation or benevolence-seeking?

The Exchanger may be guided by both some moral goal implying a negative judgement about exploitation (i.e., the goal of "not taking advantage of others") and by some "power" goal, such as that of restoring the social power he has lost when acknowledging his dependence and asking for help. In fact, as already observed while examining the *D-WILL* criterion, moving from assumed unilateral dependence of *x* on *y* to assumed bilateral dependence of each one on the other allows *x* to have a more powerful position. The Exchanger is saying to *y*: Yes, I admit I depend on you, but you depend on me. So, his goal is that of restoring the balance of power between himself and the other agent. Unlike the Benevolence-Seeker, he is mainly preoccupied with *avoiding debt and gratitude*, and consequent subjection to *y*.

### Conclusions

In this paper, we have been exploring how some social knowledge can be utilized by means of suitable heuristic strategies. This may cast a new light on *sociality*, showing how it *originates from problem solving*. Indeed, we provided a frame in which to insert and utilize knowledge about dependence, one of the most fundamental aspects of sociality. We have been showing how it can be fruitfully used to select help-givers out of a network of artificial agents. Our *CogAgent* is highly *deliberative*, essentially *goal-driven* and *knowledge-based*. These features impose a number of constraints upon his performance and applications. However, our aim is to place *CogAgent* in a realistic world where *agents are* expected to be basically *self-interested* although not necessarily hostile. In spite of the absence of some sort of pre-established harmony, such a world is not necessarily unpredictable. It can be analyzed and reasoned upon in order to take advantage of it. Our research joins a number of studies that assume cooperation in Multi-Agent Systems is not given for granted but dynamically induced during problem solving (e.g. Sycara, 1989).

Although increasing attention has been devoted to the problem solving activity in cooperative work environments, little work exists which is focused on examining ergonomic aspects, and more precisely cognitive requirements, which characterize the agents' interactions in such systems (e.g. what people really think when a colleague makes a request). In this perspective, we have identified one aspect of the problem, the help action, and a number of aspects

involved, such as searching help, as well as giving help, but also deciding to help and coping with conflicts. This paper is about the first of those aspects, searching help. It should be noted that cognitive plausibility is a natural constraint when both humans and machines are involved in a common environment. What the *CogAgent* experimental setting guarantees is a workbench where a number of idealized strategies and the knowledge involved in such tasks can be tested.

### Acknowledgements

We thank Cristiano Castelfranchi and Rosaria Conte for their thoughtful comments.

### References

- Castelfranchi, C. 1990. Social Power: A Point Missed in Multi-Agent, DAI and HCI. In Y. Demazeau, & J.P. Muller (Eds.), *Decentralized A.I.*. Amsterdam, The Netherlands: Elsevier.
- Castelfranchi, C., Miceli, M., & Cesta, A. 1992. Dependence Relations among Autonomous Agents. In E. Werner, & Y. Demazeau (Eds.), *Decentralized A.I. - 3*. Amsterdam, The Netherlands: Elsevier.
- Cohen, P.R., & Levesque, H.J. 1990. Intention Is Choice with Commitment. *Artificial Intelligence* 42: 213-261.
- Cesta, A., & Romano, G. 1992. Using Abstraction-Based Similarity to Retrieve Reuse Candidates. In J. Hendler (Ed.), *Artificial Intelligence Planning Systems: Proceedings of the First International Conference (AIPS92)*, 269-270. San Mateo, Calif.: Morgan Kaufman
- Cesta, A., & Miceli, M. 1993. In Search of Help: Strategic Social Knowledge and Plans, Technical Report IP-CNR, Rome, Italy.
- D'Aloisi, D., & Castelfranchi, C. 1993. Propositional and Terminological Knowledge Representations. *Theoretical Artificial Intelligence*, 5 (3).
- Grosz, B.J., & Sidner, C.L. 1990. Plans for Discourse. In P.R. Cohen, J. Morgan, & M.E. Pollack (Eds.), *Intentions in Communication*. Cambridge, MA: MIT Press.
- Pollack, M.E. 1990. Plans as Complex Mental Attitudes. In P.R. Cohen, J. Morgan, & M.E. Pollack (Eds.), *Intentions in Communication*. Cambridge, MA: MIT Press.
- Sycara, K. 1989. Argumentation: Planning Other Agents' Plans. In *Proceedings of IJCAI-89*, 517-522.
- Zlotkin, G., & Rosenschein, J.S. 1992. A Domain Theory for Task Oriented Negotiation. In A. Cesta, R. Conte, & M. Miceli (Eds.), *Pre-Proceedings of the Fourth European Workshop on Modeling Autonomous Agents in a Multi-Agent World*, 162-182. Rome, Italy: IP-CNR.