

# BISFAI '97

The Fifth Bar-Ilan Symposium on  
*Foundations of Artificial Intelligence*

## ABSTRACTS

June 16-18, 1997

Ramat-Gan, Israel

The Fifth Bar-Ilan Symposium on  
*Foundations of Artificial Intelligence*

June 16-18, 1997

Bar-Ilan University, Ramat-Gan, Israel

*in cooperation with*

Gelbart Institute for Mathematical Sciences  
The Leibniz Center for Research in Computer Science  
American Association for Artificial Intelligence  
Israeli Ministry of Science

**Symposium Chair**

Sarit Kraus

*Bar-Ilan University, Ramat-Gan, Israel*

**Program Co-Chairs**

D. Lehmann

*The Hebrew University, Jerusalem, Israel*

L. Joskowicz

*The Hebrew University, Jerusalem, Israel*

**Program Committee**

Y. Choueka

*Bar-Ilan U.*

I. Dagan

*Bar-Ilan U.*

R. Dechter

*UC-Irvine*

R. Feldman

*Bar-Ilan U.*

M. Golumbic

*Bar-Ilan U.*

B. Grosz

*Harvard U.*

H. Hel-Or

*Bar-Ilan U.*

J. Hendler

*U. of Maryland*

L. Joskowicz

*Hebrew U.*

D. Lehmann

*Hebrew U.*

J-J. Meyer

*Utrecht U.*

J. Pearl

*UCLA*

L. Morgenstern

*IBM TJ Watson Research*

S. Sagiv

*Hebrew U.*

E. Shamir

*Hebrew U.*

K. Sycara

*Carnegie Mellon U.*

M. Tennenholtz

*Technion*

**Local Arrangements**

A. Frank

*Bar-Ilan U.*

# A formal framework for accountable agent interactions

Carles Sierra, Pablo Noriega  
IIIA-CSIC  
Campus UAB-Bellaterra, 08193 Barcelona, Spain.  
{sierra,pablo}@iiia.csic.es

## Abstract

Agent based technologies may constitute a crucial technology for electronic commerce, but only as long as they can inspire adequate confidence to their potential users. We propose to contribute in this direction by developing a notion of "accountability" of agent interactions, which we base on two main constitutive elements: agent-mediated institutions, and shielded agents. In this paper we motivate the need for this type of accountability, and discuss examples of typical features that such accountable interactions should exhibit. We introduce a formal framework to define agent-mediated institutions and develop means for shielding agents in order to characterize accountable interactions and then show how relevant features can be identified and tested. We illustrate the applicability of these ideas in an agent-mediated electronic market place based on the traditional fish market.

Intuitively, an agent-mediated institution is the computational realization of a set of explicit enforceable restrictions imposed on a collection of dialogical agent types that concur in space and time to perform a finite repertoire of satisfiable actions.

For that characterization, we will assume that agent interactions are systematically linked to illocutions that are comprehensible to participants and refer to a basic shared ontology. We furthermore assume that all agent interactions are performed according to a pre-established *interaction protocol*, and that other explicit constraints on illocutions and their intended usage and signification can be coded into the *rules of behavior* to which participating agents are to submit. Consequently,

**Definition 1** *An Institution,  $I$ , is a 3-tuple  $I = \langle DF, PS, BR \rangle$ , where,*

1.  *$DF$  is a dialogical framework*
2.  *$PS$  is a performative structure,*
3.  *$BR$  are the rules of behavior to which participating agents are subject to.*

The *Dialogical Framework*,  $DF = \langle \text{Agents}, \text{Roles}, \text{SocialStructure}, \mathcal{CL}, \mathcal{L}, \text{Time} \rangle$  captures the intuition of *context*. It makes explicit, on one hand, the participants and their basic roles, as well as their relevant social interrelationships. On the other, it also makes explicit

the communication and object languages,  $\mathcal{CL}$ ,  $\mathcal{L}$  that will be needed for illocutions to be shared between participating agents, as well as a common notion of time to which sequencing of interactions may need to refer. Note, however, that nothing is said about the internal components of participating agents in this framework, only general rules of behavior are later on prescribed.

A *performative structure*,  $\mathcal{PS} = \langle \mathcal{S}, \text{SDG} \rangle$ , is a set of interdependent scenes. Each scene is defined as a set of agents who are each to assume a given *role*, all of the agents are subject to a common *interaction protocol*. Protocols are finite state machines where state transitions are labeled by illocutions and states have associated memory stacks (of “commitments”). We use a *Scene Dependence Graph* to establish causal and temporal co-dependencies among initial and terminal state commitments of different scenes.

$BR$  captures the notion of *intended behavior* by assigning to each participating agent a collection of schemata of the form:

$$\text{IF}(\psi_1 \wedge \dots \wedge \psi_n) \text{ THEN } \xi$$

where  $\xi \in \mathcal{CL}$ , and  $\psi_1, \dots, \psi_n \in \mathcal{L} \cup \mathcal{CL}$ .

These schemata should ideally be part of the agent’s internal theory, in order to guarantee that the agent upholds the institutional conventions. To get a handle on this issue, we introduce the notion of *governor* —or *co-agent*— to denote an agent-like entity  $c_a$  that enforces the rules of behavior for an agent-type role onto a specific agent  $a$ .

Given an agent  $a$  of role  $\text{Role}(a)$ , whose rules of behavior are  $BR(\text{Role}(a))$ ,  $c_a$  actual behavior should be consistent with those rules in the sense that for whatever illocution  $\iota(a, b, \phi, t)$ , if it is *required* by the behavior rules, it will be uttered by the agent-co-agent pair, and if  $\iota(a, b, \phi, t)$  is uttered by the agent-co-agent pair, it is not forbidden by the rules of behavior. These consistency conventions may be stated as follows:

**Definition 2** *Given an agent  $a$  of type  $\text{Role}(a)$  and subject to a set of rules of behavior  $BR(\text{Role}(a))$ , a co-agent  $c_a$  for  $a$ , is the minimal agent such that  $c_a$  upholds  $BR(\text{Role}(a))$  for  $a$ . I.e., if for every illocution  $\iota(a, b, \phi, t)$*

1.  $BR(\text{Role}(a)) \vdash \Box \iota(a, b, \phi, t) \implies T_{a \cup c_a} \vdash \iota(a, b, \phi, t)$ , and
2.  $T_{a \cup c_a} \vdash \iota(a, b, \phi, t) \implies B_{rules}(\text{Role}(a)) \not\vdash \Diamond \neg \iota(a, b, \phi, t)$

And then properly implemented. intuitively, we will have a pair of agents —an agent  $a$  and its co-agent  $c_a$  — acting as one. The co-agent  $c_a$  filters all incoming and outgoing illocutions, and in general guarantees that all rules of behavior associated with the agent’s type are actually met. Thus,  $c_a$  receives all incoming illocutions and re-sends them to  $a$  who may deliberate on them. In the meanwhile, the co-agent deliberates and prepares whatever illocutionary actions may be required by the rules of behavior, and also identifies those that may be consistent with it. When, after its deliberation,  $a$  utters an illocution, it is filtered by  $c_a$  —i.e., if the illocution is appropriate it is re-uttered by the co-agent, and if the illocution is inappropriate the co-agent inhibits it, and gives the agent an indication of failure— and the pair  $c_a$ – $a$  proceeds to a new state. However, if a triggering condition is met —e.g., if there is a time constraint for a response— and  $a$  has not been able to produce a *required answer*, then  $c_a$  should provide a default answer to guarantee compliance with the existing protocol and rules of behavior and informs the agent of its execution.