

Formalisation de la confiance sociale et mise en oeuvre dans un système multi-agent

Projet ANR ForTrust (07-11) www.irit.fr/ForTrust

IRIT-LILaC, Toulouse
EMSE-G2I, Saint-Etienne
ISTC-CNR, Rome

October 17, 2008

Context and Motivations

Castelfranchi & Falcone's theory of trust

A logical framework of trust

A logical framework of reputation

Agent Prototype

Context and Motivations

Trust: “everything is based on it”

- ▶ examples:

- ▶ “Is information from this website correct?”
- ▶ “I can book this hotel via two different sites, and prices are different; which one should I choose?”
- ▶ “Should I take this email serious?”
- ▶ “Will this webservice sell my data without consent?”
- ▶ “Will this intelligent device do what I asked for? (and not more than that?)”
- ▶ ...

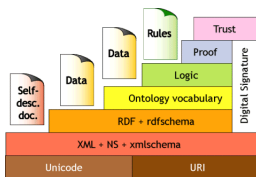
Background: existing trust management systems

“... address only narrow subsets of the overall trust management problem and often do so in a manner that is appropriate to only one application” [Blaze, Feigenbaum & Lacy, 96]

- ▶ essentially still true today
 - ▶ focus: static form of trust
 - ⇒ certificates
 - ⇒ authentication, access control
 - ▶ missing: cognitive aspect, dynamics, versatility

Human2Machine vs. Machine2Machine

- ▶ Trust for Human to Machine interactions (cf. T. Berners-Lee's cake model of the Internet)



- ▶ Trust for Machine to Machine interactions
 - ▶ Peer-to-peer systems
 - ▶ Business to Business, automatic negotiation
 - ▶ Field of **Multi-Agent Systems** (MAS)

Background: existing models of reputation

- ▶ Wide variety of multi-agent reputation models: SPORAS [Zacharia, 99], Regret [Sabater,02], socio-cognitive model [Conte & Paolucci, 02], FIRE [Huyns, 04], LIAR [Muller & Vercouter, 06], RepAge [Sabater & Paolucci, 06], ...
- ▶ varying focus \Rightarrow difficult to compare
 - ▶ reputation
 - ▶ reliability
 - ▶ sincerity
 - ▶ quality of service
 - ▶ ...
- ▶ semi-formal
- ▶ no consensus

The ForTrust project (ANR SETIN 07-11)

- ▶ statement
 - ▶ many implemented systems
 - ▶ many theoretical models
 - ▶ quantitative
 - ▶ no logical account
- ▶ ForTrust project goals
 - ▶ cognitive model of trust by Castelfranchi&Falcone [C&F]
 - ▶ formalization in modal logic
 - ▶ establish relation with reputation
 - ▶ implementation in MAS
 - ▶ application to security systems

The ForTrust project ANR SETIN (07-11)

- ▶ statement
 - ▶ many implemented systems
 - ▶ many theoretical models
 - ▶ quantitative
 - ▶ no logical account
- ▶ ForTrust project goals
 - ▶ cognitive model of trust by Castelfranchi&Falcone [C&F]
 - ▶ formalization in modal logic
 - ▶ establish relation with reputation
 - ▶ implementation in MAS
 - ▶ application to security systems

Castelfranchi & Falcone's theory of trust

General properties of trust

- ▶ related to truster's goals
 - ▶ I trust you only if you can help me to achieve my goals or to solve tasks which are important for me.
- ▶ involves complex evaluation of the trustee
 - ▶ The decision to trust you results from the evaluation of your capabilities to solve the task, your powers, your willingness,...

Trust as a subjective attitude in a social context

- ▶ trust based on objective information
 - ▶ series of observations (e.g. it has been observed for a long time that this payment protocol is secure)
 - ▶ proof (e.g. it has been proved that this cryptographic technology requires at least such amount of resources to be broken)
- ▶ trust based on social relationships
 - ▶ reputation: the certificate delivered by such institution has the reputation to guarantee safe payments
 - ▶ commitment: such internet bookshop commits itself to deliver books in less than 2 days
 - ▶ inter-individual relationships: I trust this service because my friend trusts it

C&F's informal definition

“ i trusts j to do α in order to achieve φ ”

- ⇒ defined from four more primitive concepts
 - ▶ goal, ability, power and intention
 - ▶ no formal definition of these concepts in [C&F]
- ⇒ here: represented as formulas of a modal logic of **belief**, **preference**, **time**, and **action**

A logical framework of trust

Occurrent trust vs. dispositional trust

Two perspectives on trustee's action α :

- ▶ truster believes trustee is going to do α *here and now*
⇒ *occurrent trust*
- ▶ truster believes trustee is going to do α *whenever some conditions are satisfied.*
⇒ *dispositional trust*

(cf. occurrent belief vs. dispositional belief, Searle 95)

Occurrent trust: formal definition

Agent i trusts agent j in an occurrent manner to do α in order to achieve φ iff

- ▶ i has the goal that φ ;
- ▶ i believes that j has the occurrent capability to do α
- ▶ i believes that j has the occurrent power to ensure φ by doing α
- ▶ i believes that j has the occurrent intention to do α

$$\begin{aligned} \text{OccTrust}(i, j, \alpha, \varphi) \stackrel{\text{def}}{=} & \text{Pref}_i \text{Eventually } \varphi \wedge \\ & \text{Bel}_i \neg \text{After}_{j:\alpha} \perp \wedge \\ & \text{Bel}_i \text{After}_{j:\alpha} \varphi \wedge \\ & \text{Bel}_i \text{Pref}_j \text{Does}_{j:\alpha} \top \end{aligned}$$

Dispositional trust: formal definition

Agent i trusts agent j in a dispositional manner to do α in order to achieve φ iff

- ▶ i has potentially the goal that φ ;
- ▶ i believes that j has the conditional capability to do α
- ▶ i believes that j has the conditional power to ensure φ by doing α
- ▶ i believes that j has the conditional intention to do α

$$\begin{aligned} \text{DispTrust}(i, j, \alpha, \varphi) &\stackrel{\text{def}}{=} \\ &\neg \text{Bel}_i \text{Henceforth} \neg \text{Pref}_j \text{Eventually } \varphi \wedge \\ &\text{Bel}_i \text{Henceforth} (\kappa_{\text{OccCap}}(j, \alpha) \rightarrow \neg \text{After}_{j:\alpha} \perp) \wedge \\ &\text{Bel}_i \text{Henceforth} (\kappa_{\text{OccPower}}(j, \alpha, \varphi) \rightarrow \text{After}_{j:\alpha} \varphi) \wedge \\ &\text{Bel}_i \text{Henceforth} (\kappa_{\text{OccIntends}}(j, \alpha) \rightarrow \text{Pref}_j \text{Does}_{j:\alpha} \top) \end{aligned}$$

From dispositional to occurrent trust

Theorem

Suppose Bel_i and Henceforth are normal modal operators, and suppose Henceforth obeys the axiom

$\text{Henceforth } \varphi \rightarrow \varphi$.

Then

$$\vdash \left(\begin{array}{l} \text{DispTrust}(i, j, \alpha, \varphi) \\ \wedge \text{Pref}_j \text{Eventually } \varphi \\ \wedge \text{Bel}_i \kappa \text{OccCap}(j, \alpha) \\ \wedge \text{Bel}_i \kappa \text{OccPower}(j, \alpha, \varphi) \\ \wedge \text{Bel}_i \kappa \text{OccIntends}(j, \alpha) \end{array} \right) \rightarrow \text{OccTrust}(i, j, \alpha, \varphi)$$

Two kinds of occurrent trust: trust in the action vs. trust in the inaction

- ▶ truster believes that trustee is in condition to *further* the achievement of his goals, and he will *do* that
⇒ *trust in the action*
- ▶ truster believes that trustee is in condition to *endanger* the achievement of his goals, but he will *refrain* from doing that.
⇒ *trust in the inaction*

(cf. *doing* vs. *refraining* (or *forbearing*), Von Wright 63, Belnap et al. 01)

Trust in the action vs. trust in the inaction: formal definitions

$$\text{OccTrust}(i, j, \alpha, \varphi) \stackrel{\text{def}}{=} \text{Pref}_i \text{Eventually } \varphi \wedge \\ \text{Bel}_i \neg \text{After}_{j:\alpha} \perp \wedge \\ \text{Bel}_i \text{After}_{j:\alpha} \varphi \wedge \\ \text{Bel}_i \text{Pref}_j \text{Does}_{j:\alpha} \top$$

$$\text{OccTrust}(i, j, \sim\alpha, \varphi) \stackrel{\text{def}}{=} \text{Pref}_i \text{Eventually } \varphi \wedge \\ \text{Bel}_i \neg \text{After}_{j:\alpha} \perp \wedge \\ \text{Bel}_i \text{After}_{j:\alpha} \neg \varphi \wedge \\ \text{Bel}_i \neg \text{Pref}_j \text{Does}_{j:\alpha} \top$$

A logical framework of reputation

Formal definition of reputation: the building blocks

$Rep(I, j, \alpha, \varphi)$ = “ j has reputation to do α in order to achieve φ in group I ”

- ▶ also to be defined from :
 - ▶ j 's capability,
 - ▶ j 's willingness,
 - ▶ j 's power.
- ▶ formally builded from:
 - ▶ *group belief* of the evaluating agents
(vs individual belief in the case of trust)
 - ▶ *group goals* of the evaluating agents
(vs individual goals in the case of trust)

⇒ To be defined: group beliefs, group goals

Group beliefs and group goals

$\text{GroupBelief}_I \varphi$ = “it is public in the group I that φ ”

- ▶ group belief \neq mutual belief
(group belief does not imply individual belief
[Tuomela 95, 02; Gaudou et al. 06, 08])

$\text{GroupPref}_I \varphi$ = “the group I *prefers* that φ ”

- ▶ group goals: broad sense (including norms, standards, values ...)
- ▶ does not imply individual preferences

Formal definition of reputation

$Rep(I, j, \alpha, \varphi) \stackrel{\text{def}}{=}$

$\neg \text{GroupBelief}_I \text{ Henceforth } \neg \text{GroupPref}_I \text{ Eventually } \varphi \wedge$
 $\text{GroupBelief}_I \text{ Henceforth } (\kappa_{\text{OccCap}}(j, \alpha) \rightarrow \neg \text{After}_{j:\alpha} \perp) \wedge$
 $\text{GroupBelief}_I \text{ Henceforth } (\kappa_{\text{OccPower}}(j, \alpha, \varphi) \rightarrow \text{After}_{j:\alpha} \varphi) \wedge$
 $\text{GroupBelief}_I \text{ Henceforth } (\kappa_{\text{OccIntends}}(j, \alpha) \rightarrow \text{Pref}_j \text{ Does}_{j:\alpha} \top)$

Remember: $DispTrust(i, j, \alpha, \varphi) \stackrel{\text{def}}{=}$

$\neg \text{Bel}_i \text{ Henceforth } \neg \text{Pref}_i \text{ Eventually } \varphi \wedge$
 $\text{Bel}_i \text{ Henceforth } (\kappa_{\text{OccCap}}(j, \alpha) \rightarrow \neg \text{After}_{j:\alpha} \perp) \wedge$
 $\text{Bel}_i \text{ Henceforth } (\kappa_{\text{OccPower}}(j, \alpha, \varphi) \rightarrow \text{After}_{j:\alpha} \varphi) \wedge$
 $\text{Bel}_i \text{ Henceforth } (\kappa_{\text{OccIntends}}(j, \alpha) \rightarrow \text{Pref}_j \text{ Does}_{j:\alpha} \top)$

Agent Prototype

Implementation of a prototype

- ▶ Implementation of an agent reasoning about trust using our model
- ▶ Logical formalization used as a specification
- ▶ Operational prototype of the trust framework
 - ▶ in a BDI language: Jason
 - ▶ deployed on the “standard” testbed: ART
- ▶ Reputation framework not yet included

Application to the ART Testbed

- ▶ Definition of ART's Occurent Capability

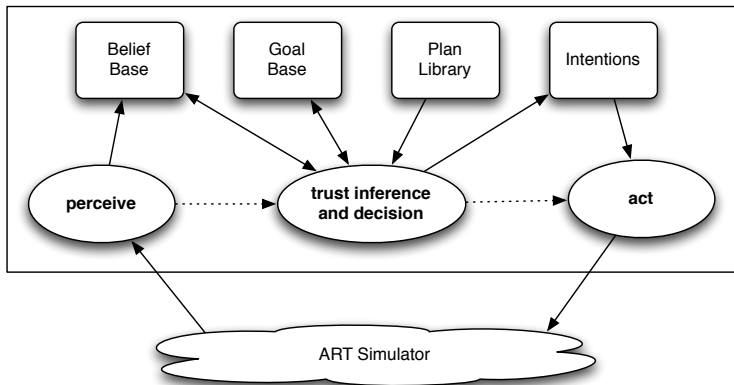
$$\begin{aligned} & \textit{Believes}(i, \textit{OccAct}(j, \alpha), x) \\ & \leftarrow \textit{Believes}(i, \textit{opinions_count}(j, \textit{asked}, \textit{got}), 1) \wedge \\ & \quad \textit{asked} > 0 \wedge x = \frac{\textit{got}}{\textit{asked}} \wedge x > \epsilon \end{aligned} \tag{1}$$

- ▶ Definition of ART's Occurent Power

$$\begin{aligned} & \textit{Believes}(i, \textit{OccPower}(j, \alpha, \varphi), y) \\ & \leftarrow \textit{Believes}(i, \textit{sincere}(j), 1) \wedge \\ & \quad \textit{Believes}(i, \textit{painting}(e, p), 1) \wedge \\ & \quad y = \textit{image}_t(j, e) \wedge y > \delta \end{aligned} \tag{2}$$

- ▶ ...

Agent architecture in Jason



Legend



.....▶ control flux

————▶ data flux

Trust evaluation in Jason

```
trust(J, Action, Goal)[strength(C)] : - // trust inference
    .intend(Goal) & // I have the goal
    does(J, Action)[strength(X)] & // J is capable and intend
    after(J, Action, Goal)[strength(Y)] & // he has the power
    C = math.min(X, Y). // strength of the trust

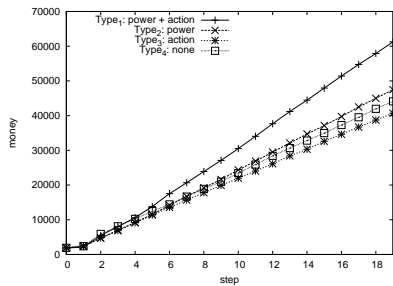
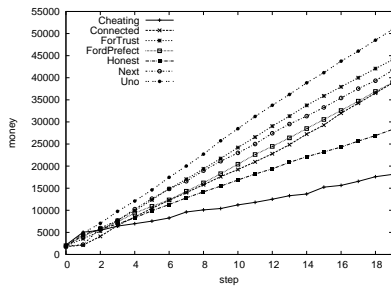
+painting(Era, P) ←!good_eval(P).

does(J, appraise(P))[strength(X)] : -
    opinions_count(J, Asked, Provided) &
    Asked > 0 &
    X = Provided/Asked & X > 0.9.

after(J, appraise(P), _)[strength(Y)] : -
    sincere(J) & painting(Era, P) &
    image(J, Era, Y) & Y > 0.5.
```

(3)

Experiments



Conclusion & Future Works

- ▶ Study of the components of C & F theory
- ▶ Proposition of a logical formalization
 - ▶ compliant to the properties of C & F
 - ▶ extended to reputation
- ▶ A first implemented prototype
 - ▶ logical formalization used as specifications
 - ▶ use of MAS language (Jason) and platform (ART)
- ▶ Future works
 - ▶ improvement and generalization of the implementation
 - ▶ application to security systems