Stronger Notion, mentalistic notions: anthropomorphism

- Information attitudes
  - Knowledge
  - Belief
- Pro-attitudes (guide the agent's actions)
  - Intention
  - Goal
  - Desire
  - Obligation (Shoham: Agent-Oriented programming)
  - Emotion (Joe Bates: Believable agents)

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

- What is an agent?
- Weak Notion:
  - Autonomy (operate without intervention, own control)
  - Social ability (agent communication language)
  - Reactivity (perception of the environment, action)
  - Pro–activeness (goal–directed,  $\it taking~the~initiative)$
- Self-contained, concurrently executed software process, that encapsulates some state and is able to communicate with other agents via message passing.
- Object-Based concurrent programming (Agha,86: ACTORS)

# • AGENTS:

- Theories (specification)
- Architectures (implementation)
- Languages (programming)

# • INTENTIONAL SYSTEMS

- Very complex systems: animistic, intentional explanations,
   abstraction versus mechanistic interpretations
- Information attitudes: belief, knowledge
- Pro-attitudes: desire, intention, obligation, commitment, choice, ...

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

3

- Visual representation (Maes: Interface Agents)
- Mobility
- Veracity (no communicate false information)
- Benevolence (it try to do what is asked of it)
- Rationality (it will act to achieve its goals)
- Solitary, parasite, social, selfish ...

## **SOLUTIONS:**

- Syntactic:
  - Modal language: non-truth-functional modal operators
  - Meta–language: a first–order language containing formulae of some other  $object{-}language$

$$Bel(Janine, \lceil Father(Zeus, Cronos) \rceil)$$

- Semantic:
  - Possible worlds
  - Sentencial models

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

5

Janine believes Cronos is the father of Zeus

Two problems:

- 1. Syntactic: not well-founded formula of FOL
- 2. Semantic: Suposse (Zeus = Jupiter), then Bel(Janine, Father(Jupiter, Cronos)) can not be inferred. Belief is not truth functional.

The operators are duals:

$$\Box \varphi \Leftrightarrow \neg \Diamond \neg \varphi$$

$$\Diamond \varphi \Leftrightarrow \neg \Box \neg \varphi$$

- Axiom K (Kripke):  $\Box(\varphi \Rightarrow \psi) \Rightarrow (\Box\varphi \Rightarrow \Box\psi)$
- Necessitation rule: if  $\varphi$  is valid, then  $\Box \varphi$  is valid

Four axioms:

**T** (reflexive) 
$$\Box \varphi \Rightarrow \varphi$$

**D** (serial) 
$$\Box \varphi \Rightarrow \Diamond \varphi$$

4 (transitive) 
$$\Box \varphi \Rightarrow \Box \Box \varphi$$

5 (euclidean) 
$$\Diamond \varphi \Rightarrow \Box \Diamond \varphi$$

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

7

# NORMAL MODAL LOGIC

Classic propositional logic + two operators:  $\Box$  (necessarily) and  $\Diamond$  (possibly)

Atomic propositions  $Prop = \{p, q, \ldots\}$ 

- 1. if  $p \in Prop$  then p is a formula
- 2. if  $\varphi, \psi$  are formulae, then so are  $\neg \varphi$  and  $\varphi \lor \psi$
- 3. if  $\varphi$  is a formula, then so are  $\Box \varphi$  and  $\Diamond \varphi$

Accesibility relation: what world are accesible from every other world

- $\Box \varphi$  is true if  $\varphi$  is true in every accesible world
- $\diamond \varphi$  is true if  $\varphi$  is true in at least one accesible world

**T** (knowledge)  $K_i \varphi \Rightarrow \varphi$  (distintion between knowledge and belief)

**D** (non-contradictory)  $K_i \varphi \Rightarrow \neg K_i \neg \varphi$ 

4 (positive introspection)  $K_i \varphi \Rightarrow K_i K_i \varphi$ 

5 (negative introspection)  $\neg K_i \varphi \Rightarrow K_i \neg K_i \varphi$ 

Epistemic Logic: KTD45 (idealised) knowledge

Doxastic Logic: KD45 (idealised) Belief

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

9

## EPISTEMIC LOGIC:

 $\Box \varphi$ ,  $K\varphi$ , it is known that  $\varphi$ 

Agents:  $K_i\varphi$ , i knows that  $\varphi$ 

Logical omniscience problem:

**Necessitation rule:** if  $\varphi$  is valid, then  $K_i\varphi$  is valid. An agent knows all valid formulae.

**Axiom K:**  $K_i(\varphi \Rightarrow \psi) \Rightarrow (K_i \varphi \Rightarrow K_i \psi)$ . Agent's knowledge is closed under implication.

# INTENTIONS: COHEN AND LEVESQUE

- Primary Modalities
  - (BEL x p) Agent x believes p
  - (GOAL x p) Agent x has a goal of p
  - (DONE x a) Agent x has just performed action a
  - (DOING x a) Agent x is just performing action a
- Event Sequences
  - -a; b Action a followed by action b
  - -a? Test action
- Standard future operations of temporal logic
  - $\Box$  Always
  - $\diamondsuit$  Sometimes

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

11

# THEORIES OF AGENCY

- Dynamic aspects
- $\bullet$  Relation between information attitudes and pro–attitudes
- Changes of the cognitive state over time
- Environment changes de cognitive state
- How to perform actions

```
INTENTION
```

```
(INTENT x p q) =

(PGOAL x
(DONE x

[UNTIL (DONE x a) (BEL x (DOING x a))]?; a)
q)
```

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

PERSISTENT GOALS: Agent x, Goal p, Motive q

$$(PGOAL \ x \ p \ q) =$$

- 1. (BEL  $x \neg p$ )  $\land$
- 2. (GOAL  $x \diamond p$ )  $\wedge$
- 3. (UNTIL [(BEL x p)  $\vee$  (BEL x  $\neg \Box p$ )  $\vee$  (BEL x  $\neg q$ )] (GOAL x  $\Diamond p$ ))
- 1. Agent x believes goal p is currently false
- 2. Agent x wants goal p to be eventually true
- 3. Continue until agent x believes that p is true or will never be true, or that the motivation q is no longer present

13

	OOP	AOP
Basic unit	object	agent
State	unconstrained	beliefs, commitments,
		capabilities, choices,
Computation	message passing,	message passing,
	methods	methods
Type of messages	unconstrained	inform, request, offer,
		promise, decline,
Constrains	none	honesty, consistency,
on methods		

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

# 15

# AGENT ORIENTED PROGRAMMING (AOP) $YOAV \ SHOHAM \\ STANFORD \ UNIVERSITY$

- specialization OOP
- mental state: beliefs, decisions, capabilities, obligations
- speech act theory: informing, requesting, offering, ...

17

Decision (choice) freedom central to the notion of agenthood

$$DEC_a^t \varphi \stackrel{def}{=} OBL_{a,b}^t \varphi$$

commitment to oneself

Capability at time t agent a is capable of  $\varphi$ 

$$CAN_a^t \varphi$$

$$CAN_{robot}^{5}open(door)^{8}$$

immediate version of CAN

$$ABLE_a\varphi \stackrel{def}{=} CAN_a^{time(\varphi)}\varphi$$

$$ABLE_{robot}open(door)^5 \stackrel{def}{=} CAN_{robot}^5 open(door)^5$$

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

LANGUAGE

 ${\bf Time} \ \ {\bf simple} \ \ {\bf point-based} \ \ {\bf temporal} \ \ {\bf language}$ 

$$holding(robot, cup)^t$$

Action instantaneous

$$raise\_arm(robot)^t$$

**Belief** at time t agent a believes a (recursively defined) sentence  $\varphi$ 

$$B_a^t \varphi$$

$$B_a^3 B_b^{10} like(a,b)^7$$

**Obligation** at time t agent a is obligated (committed) to agent b about  $\varphi$ 

$$OBL_{a,b}^t \varphi$$

Introspection no total introspective capabilities, obligation

- for any t, a, b,  $\varphi$ :  $OBL_{a,b}^t \varphi \equiv B_a^t OBL_{a,b}^t \varphi$
- for any t, a, b,  $\varphi$ :  $\neg OBL_{a,b}^t \varphi \equiv B_a^t \neg OBL_{a,b}^t \varphi$

Persistence of mental state mental states persist over time

- beliefs persist (until learn contradiction)
- absence of beliefs persist (until learn)
- obligations (and decisions) persist (until they are revoked)
- capabilities are fixed

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

19

## **PROPERTIES**

Internal consistency beliefs and obligations are internally consistent

- for any t, a:  $\{\varphi : B_a^t \varphi\}$  is consistent
- for any t, a:  $\{\varphi: OBL_{a,b}^t \varphi \text{ for some } b\}$  is consistent

Good faith agents commit only to what they believe themselves capable of and only if they really mean it

• for any t, a, b,  $\varphi$ :  $OBL_{a,b}^t \varphi \Rightarrow B_a^t((ABLE_a \varphi) \wedge \varphi)$ 

GENERIC INTERPRETER

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

21

# GENERIC AGENT INTERPRETER

- Basic loop
  - 1. read current messages, update mental state (belief and commitments)
  - 2. execute commitments for the current time, possibly resulting in further belief change
- message passing: addressable by name
- clock: synchronization

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

J.Puyol

23

## AGENT0

• Fact statements: atomic objective sentences

```
(t (employee (smith acme)))
(NOT (t (employee (smith acme))))
```

- Private and communicative action statements: private or communicative; conditional or unconditional
  - private

```
(DO t p-action)
```

- communicative: informing, requesting, cancelling

```
(INFORM t a fact)
(REQUEST t a action)
```

# AGENTO PROGRAM:

- Time grain
- Agent's capabilities
- Initial beliefs
- A sequence of commitment rules

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

Variables

- Commitment rules
  - conditional: message condition, logical combination of message patterns

```
(From Type Contents) (a REQUEST (DO t walk))
```

- commitment rule

25

• the mental condition part allows one to prevent commitment to incompatible actions

```
((?!time (rotate wheelbase ?degrees))
  (NOT ((CMT ?x) ?!time (service wheelbase))))
```

- belief change affect capabilities, private actions depend of mental preconditions
- belief update: examines the current commitments preconditions
- UNREQUEST: removes the commitment

Intel.ligència Artificial Distribuïda (1996)

J.Puyol

27

# AGENTO INTERPRETER, DATABASES:

# beliefs

- ullet updated when informed or as a result of a private action
- they incorporate any fact of which they are informed, retracting the contradictory atomic belief if that were previously held

## commitments

- structure: (agent action)
- ullet removed as a result of a belief change or a UNREQUEST message

# capabilities

- structure: (privateaction mntlcond)
- fixed