

User Evaluated Explainable BDI Agents

Joost Broekens
Maaïke Harbers
Koen Hindriks
Karel van den Bosch
Catholijn Jonker
John-Jules Meyer

Background

- Explainable autonomous agents...why?
 - VR Training
 - Games
 - Interactive storytelling
 - Flexibility and believability
 - Agent debugging
 - Use agent reasoning-language translation
- Related to expert system, but not the same.
 - Main diff: expl. for *actions* derived from autonomous behavior and reasoning, not necessarily a conclusion of a system build for explanation

Background

- Main aim:
 - explain autonomous actions in for humans understandable way.
- Two main streams
 - Directly derived from agent behavior and reasoning
 - Benefit: Agent model is all, Drawback: quality of explanation?
 - Additional information and or knowlegde is needed
 - Benefit: fine tune explanation, Drawback: more modeling
- Our appraoch:
 - BDI Agent, directly derived.

Main question

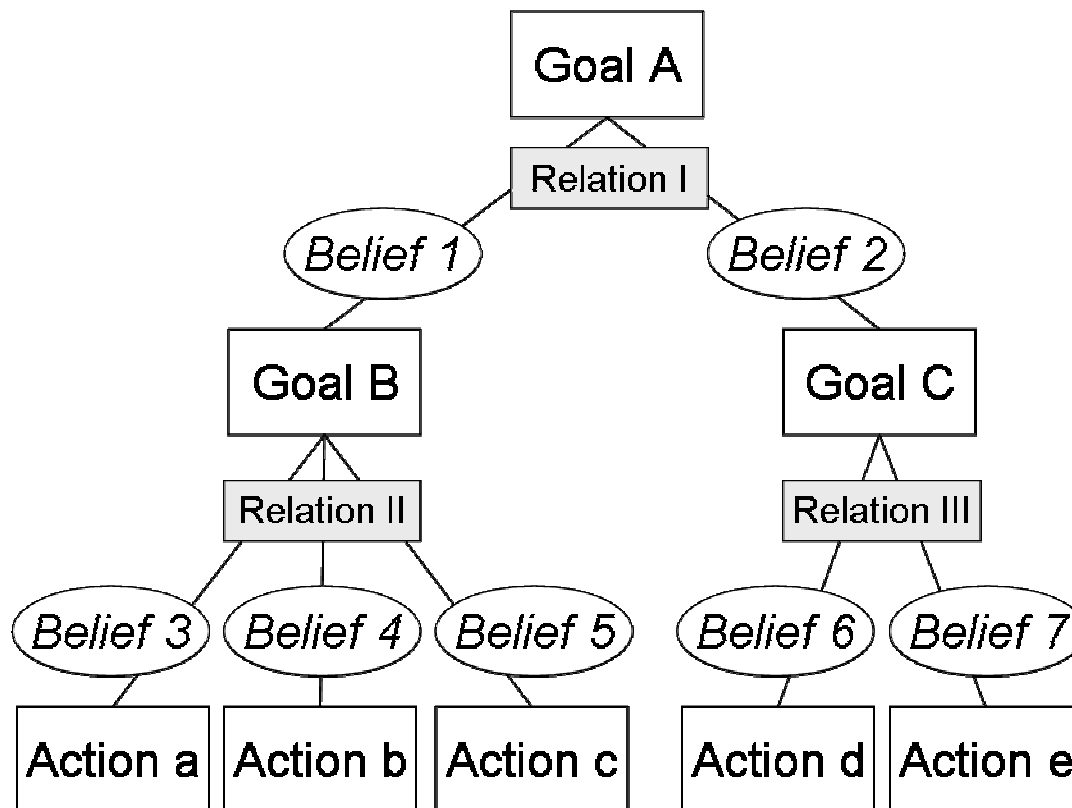
- Is it possible to generate useful explanations for an agent's actions based only on (the history of) that agent's BDI state?
- How does the usefulness depend on the type of action to be explained?
- Are there guidelines on how to do it?
 - For explainable BDI Agent development
 - For generating the explanation

"Loose" hypothesis

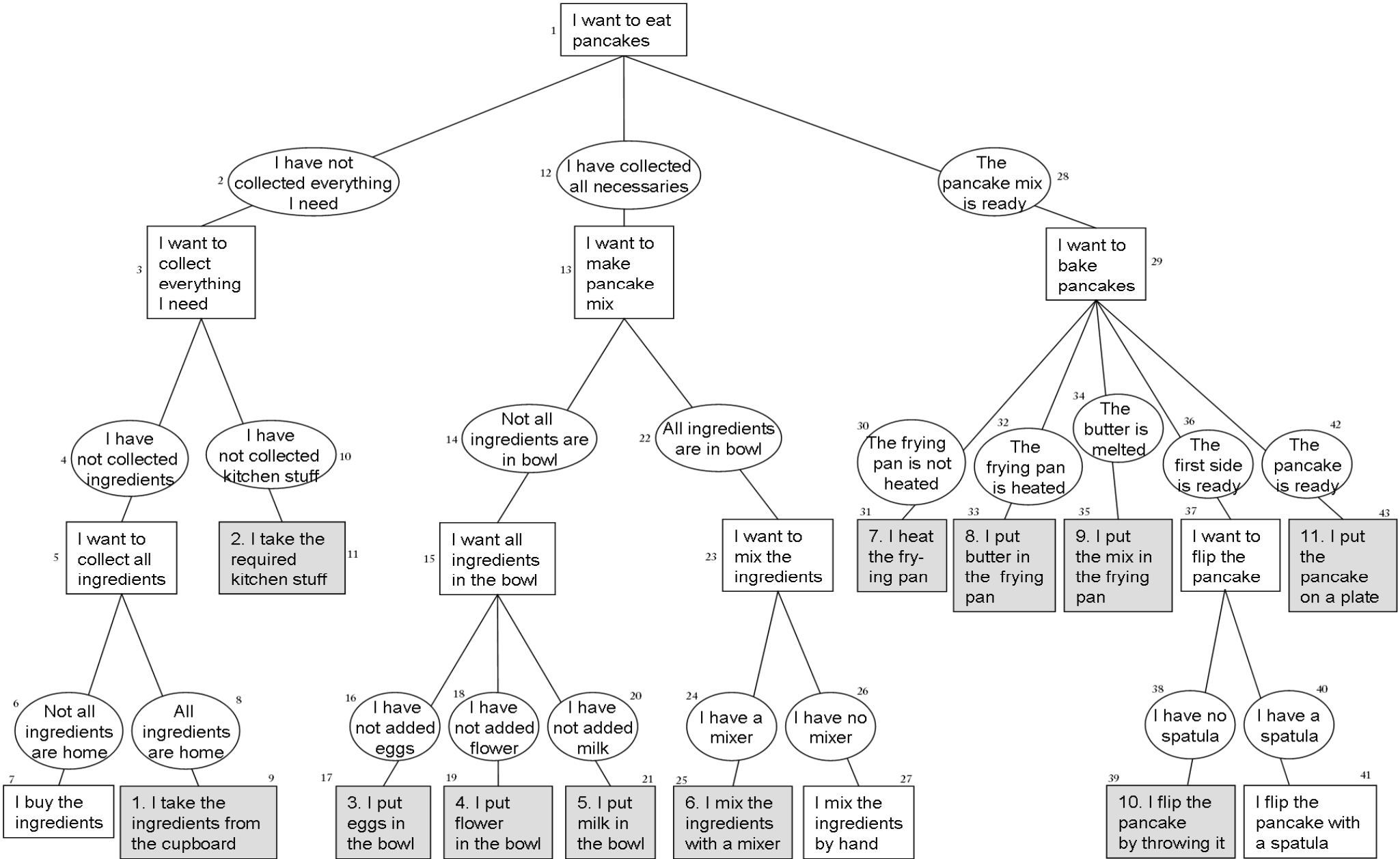
- Given a BDI agent:
 - Given an explanation mechanism:
 - Given different types of action
-
- The usefulness of an explanation mechanism depends on the type of action to be explained generated by the BDI agent.

Action types

- All: *and* relation between siblings
- Seq: *and* relation in order
- One: *xor* relation between siblings

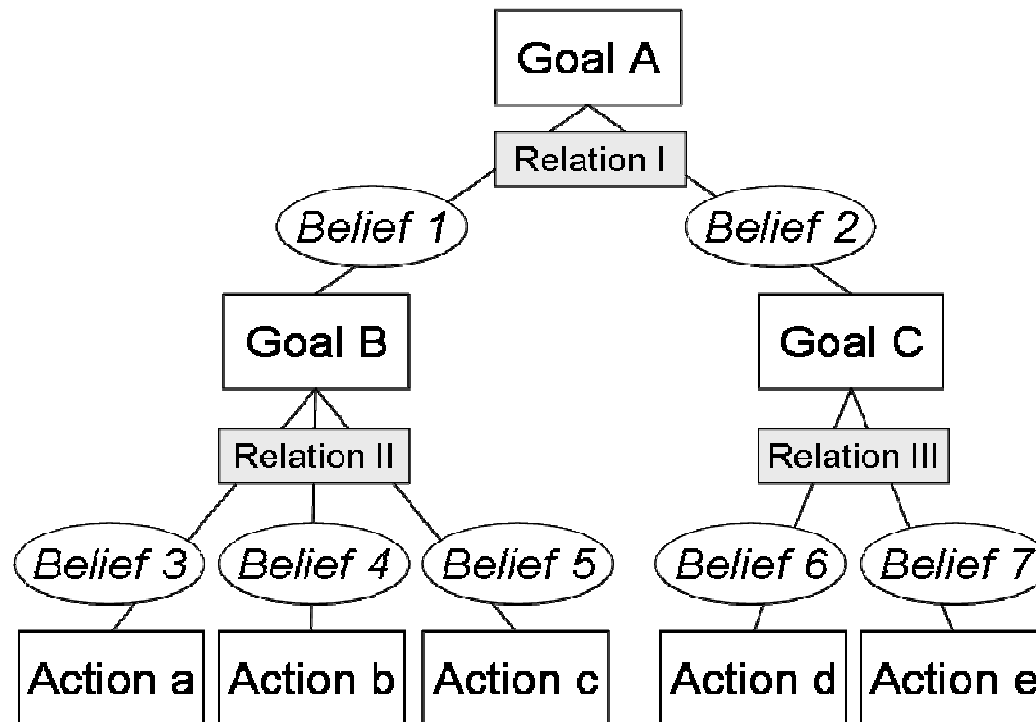


BDI Agent



Explanation mechanisms

- Alg 1: Explain action using nearest goal up
- Alg 2: Explain action using enabling belief
- Alg 3: Explain action using next goal/action to be achieved/executed



Detailed hypothesis

- The usefulness and naturalness of an explanation generated from a resulting goal-plan tree of a BDI agent depends on
 - the interaction between
 - action type (all, seq, one) and
 - explanation mechanism (goal up, enabling believe, next goal/action).

Experiment

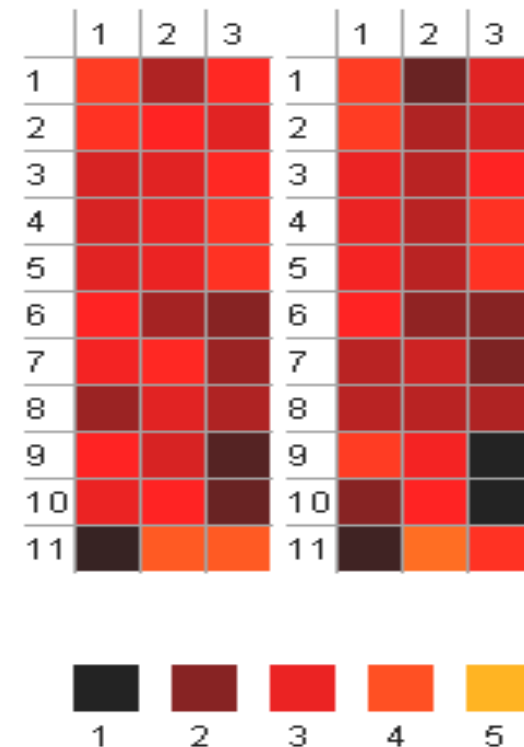
- Simple domain: cooking pancakes
- 30 users, 10 per explanation mechanism.
- 11 actions of different types.
- Each mechanism generates explanation for each of the 11 actions
- Each user rated 11 actions.
- Each user gave its own explanation based on available elements in goal-plan tree.

Example questionnaire

- (Show [questionnaire](#) for condition 1)

Rating results

- Analysis showed two important things:
 - Explanations with parent goal are slightly better on average
 - Important interaction effect between action and expl, but not as hypothesized.
 - e.g. all actions 3,4,5 better explained using next goal
- In detail:
 - left=usefulness, right=naturalness
 - x=algorithm
 - y=action



Concluding observations

- Main hypothesis confirmed
- Not as straightforward as we hypothesized
- All actions seem to need at least one extra element in addition to parent goal.
 - e.g. type one (or):
 - Enabling condition, parent's condition, and parent goal
 - e.g. for type seq:
 - Enabling condition needed

Hypothesized guidelines

- Parent goal is most important default, make it meaningful
- XOR /or choices should be modelled with meaningful subgoal (not abstract choice action) or use more elements.
- Actions starting new phase: include parent goal and enabling condition of parent goal.
- Actions in sequence: include enabling condition.