

Multi-Agent Systems

Albert-Ludwigs-Universität Freiburg



**UNI
FREIBURG**

Bernhard Nebel, Felix Lindner, and Thorsten Engesser

Summer Term 2017

- 1 Introduction
- 2 Agent-Based Simulation
- 3 Agent Architectures
- 4 Beliefs, Desires, Intentions
 - The GOAL Agent Programming Language
 - Introduction to Modal Logics
 - Epistemic Logic
 - BDI Logic
- 5 Norms and Duties
- 6 Communication and Argumentation
- 7 Coordination and Decision Making

■ Specification

- The intended behavior of a MAS can be specified using a logical specification language. The concrete program is derived from the specification (manually, in most cases).

■ Verification

- Once a program \mathcal{P} is built, one wishes to be able to proof that it behaves according to its specification φ_p , i.e., $\mathcal{P} \models \varphi_p$.

■ Agent programming

- Agents themselves can be realized deductive reasoners: What an agent knows is represented as formulae of a formal language. The agent can reason about these formulae to derive new formulae, or to determine what to do next.

Definition

Model checking is an automated technique that, given a finite-state model of a system and a formal property, systematically checks whether this property holds for (a given state in) that model.

- Model of the system \Rightarrow How the system actually behaves.
- Formal properties \Rightarrow How the system should behave.
 - Safety: something bad never happen
 - Liveness: something good eventually happens
 - Fairness: if something may happen frequently, it will happen

Definition

Runtime verification is the discipline of computer science that deals with the study, development, and application of those verification techniques that allow checking whether a run of a system under scrutiny satisfies or violates a given correctness property.

⇒ Testing using formal methods.

- **Question:** Does a given BDI agent act right (viz., according to some specified properties)?
- **Required**
 - Representation of the agent's execution.
 - Language to specify the wanted properties.
 - Algorithm to check if some given properties hold in some representation of an execution.

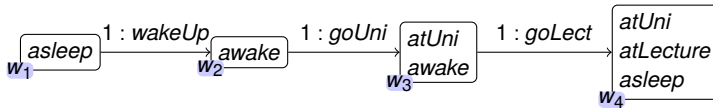
- Ingredients:
 - Action and time
 - Belief and preference
 - Definition of intention

¹The following notations are according to Meyer, Broersen, Herzig (2015). They slightly deviate from the original notations in Cohen, Levesque (1990).

A BDI Kripke model is a tuple $M = (W, R, B, P, V)$, where:

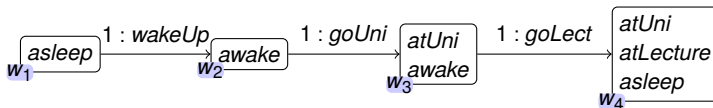
- W is a set of possible worlds.
- $R : I \times A \rightarrow W \times W$
 - Accessibility relations $R_{i:\alpha} \subseteq W \times W$ for each action $i : \alpha$.
 - (W, R) is a linear transition system.
- $B : I \rightarrow W \times W$
 - Accessibility relations $B_i \subseteq W \times W$ for each agent i .
 - Every B_i is serial, transitive, Euclidean (**KD45**).
- $P : I \rightarrow W \times W$
 - Accessibility relations $P_i \subseteq B_i \subseteq W \times W$ for each agent i .
 - Every P_i is serial (**KD**).
- $V : P \rightarrow 2^W$
 - Maps atomic propositions to their extension $V(p) \subseteq W$.

Actions: Example I



- $M, w \models \text{Happ}_{i:\alpha} \varphi$ iff. there is a $w' \in W$ s.th. $(w, w') \in R_{i:\alpha}$ and $M, w' \models \varphi$ (\Rightarrow diamond operator).
- $M, w \models \text{IfHapp}_{i:\alpha} \varphi$ iff. $M, w \models \neg \text{Happ}_{i:\alpha} \neg \varphi$ (\Rightarrow box operator).
- $M, w \models \exists \alpha \text{Happ}_{i:\alpha} \varphi$ iff. there are agent i , action type α and w' s.th. $(w, w') \in R_{i:\alpha}$ and $M, w' \models \varphi$.
- Linearity is characterised by the axiom $(\text{Happ}_{i:\alpha} \top \wedge \text{Happ}_{j:\alpha'} \varphi) \rightarrow \text{IfHapp}_{i:\alpha} \varphi$. ("If $i : \alpha$ is executable and $j : \alpha'$ brings about φ , then also $i : \alpha$ brings about φ ."

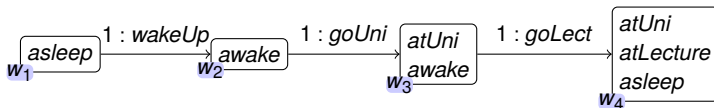
Actions: Example II



- $M, w_1 \models \text{Happ}_{1:\text{wakeUp}} \text{awake}$
- $M, w_2 \models \exists \alpha \text{Happ}_{1:\alpha} \exists \beta \text{Happ}_{1:\beta} \text{atLecture}$

- $M, w \models X\phi$ iff. $M, w' \models \phi$ for some w' s.th. $(w, w') \in R_{i:\alpha}$ for some $i : \alpha$.
- $M, w \models F\phi$ iff. $M, w \models \phi$ or $M, w \models XF\phi$.
- $M, w \models G\phi$ iff. $M, w \models \neg F\neg\phi$.
- $M, w \models \psi U\phi$ iff. $M, w \models \phi$ or $(M, w \models \psi$ and $M, w' \models \psi U\phi$) for some w' s.th. $(w, w') \in R_{i:\alpha}$ for some $i : \alpha$.

Time: Example



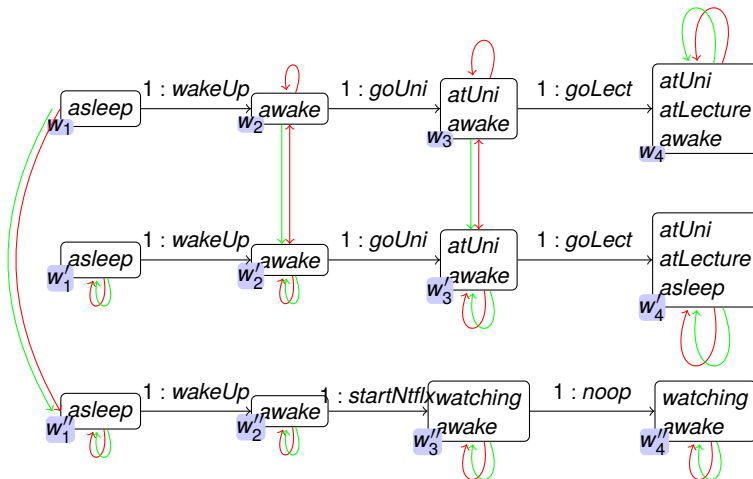
- $M, w_1 \models X(\text{awake} \text{U} \text{atLecture})$
- $M, w_1 \models \text{atSleep} \wedge X \text{F} \text{atSleep}$
- $M, w_1 \models G(\text{atSleep} \leftrightarrow \neg \text{awake})$
- $M, w_1 \models F \exists \alpha \text{Happ}_{1:\alpha} \text{atLecture}$

- $M, w \models Bel_i \varphi$ iff. for all w' s.th. $(w, w') \in B_i$: $M, w' \models \varphi$.
 - $Know_i \varphi \stackrel{\text{def}}{=} \varphi \wedge Bel_i \varphi$.
- $M, w \models Pref_i \varphi$ iff. for all w' s.th. $(w, w') \in P_i$: $M, w' \models \varphi$.
 - In the original *Pref* is called *Goal*. Some authors call it *Choice*. It is meant to be a “chosen desire” (consistent!).

Properties

- For Bel_i all properties for **KD45** operators.
- For $Pref_i$ all properties for **KD** operators.
- $\models Bel_i \varphi \rightarrow Pref_i \varphi$ (Realism)
- $\models (Pref_i \varphi \wedge Bel_i (\varphi \rightarrow \psi)) \rightarrow Pref_i \psi$.

Belief and Preference: Example



- Because of realism, all believed propositions are preferred propositions. But it only makes sense for an agent to adopt some goal φ if φ is believed to be false (compare to the GOAL programming language).

- Agent i has the **achievement goal** that φ iff i prefers that φ is eventually true and believes that φ is currently false:

$$AGoal_i \varphi \stackrel{\text{def}}{=} Pref_i F \varphi \wedge Bel_i \neg \varphi$$

Example

In the Netflix-vs.-Lecture dilemma:

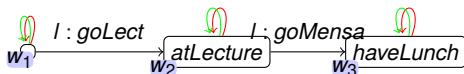
- $M, w_1 \not\models AGoal_1(\text{asleep})$
- $M, w_1 \models AGoal_1(\text{watching})$

- $\models AGoal_i \neg \phi \rightarrow \neg AGoal_i \phi$.
 - Check that $AGoal_i \neg \phi \wedge AGoal_i \phi$ is unsatisfiable, because the achievement goal that $\neg \phi$ implies to believe ϕ , and the achievement goal that ϕ implies to believe $\neg \phi$. This contradicts axiom D ($Bel_i \phi \rightarrow \neg Bel_i \neg \phi$). □
- $\not\models AGoal_i(\phi \wedge \psi) \rightarrow AGoal_i \phi \wedge AGoal_i \psi$ (for exercise).
- $\not\models AGoal_i \phi \wedge AGoal_i \psi \rightarrow AGoal_i(\phi \wedge \psi)$.
- $\not\models AGoal_i(\phi \vee \psi) \rightarrow AGoal_i \phi \vee AGoal_i \psi$.
- $\not\models AGoal_i \phi \vee AGoal_i \psi \rightarrow AGoal_i(\phi \vee \psi)$.

$$\not\models AGoal_i \varphi \wedge AGoal_i \psi \rightarrow AGoal_i (\varphi \wedge \psi)$$



“Lisa has the goal to listen to the lecture and she has the goal to have dinner” vs. “Lisa has the goal to listen to the lecture and to have dinner”

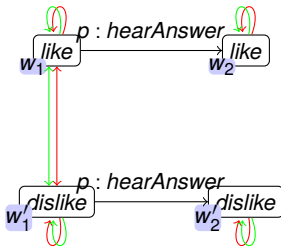


- $M, w_1 \models AGoal_i(atLecture) \wedge AGoal_i(haveLunch)$
- $M, w_1 \not\models AGoal_i(atLecture \wedge haveLunch)$

$$\not\models AGoal_i(\varphi \vee \psi) \rightarrow AGoal_i\varphi \vee AGoal_i\psi.$$

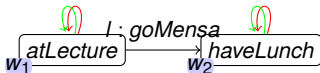


“Paul asks Lisa whether she likes him.” (Paul does not prefer any of the two possible answers.)



- $M, w_1 \models AGoal_p(Know_p like \vee Know_p dislike)$
- $M, w_1 \not\models AGoal_p(Know_p like)$
- $M, w_1 \not\models AGoal_p(Know_p dislike)$

$$\not\models AGoal_i \phi \vee AGoal_i \psi \rightarrow AGoal_i (\phi \vee \psi)$$

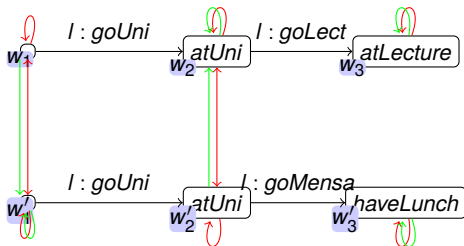


- $M, w_1 \models AGoal_1(haveLunch)$
- $M, w_1 \not\models AGoal_1(atLecture)$
 - Reason: $M, w_1 \not\models Bel_I \neg atLecture$
- $M, w_1 \not\models AGoal_1(atLecture \vee haveLunch)$
 - Reason: $M, w_1 \not\models Bel_I(\neg(atLecture \vee haveLunch))$

Achievement Goal: Too weak for Intention



- Agents can change their preferences whenever they like:
Lack of commitment!



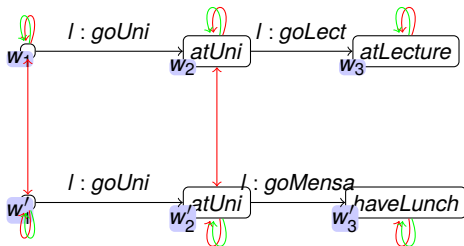
- $M, w_1 \models AGoal_1(haveLunch)$
- $M, w_2 \models \neg AGoal_1(haveLunch)$

Say a problem solver is confronted with the classic situation of a heroine, called Nell, having been tied to the tracks while a train approaches. The problem solver, called Dudley, knows that “If Nell is going to be mashed, I must remove her from the tracks.” When Dudley deduces that he must do something, he looks for, and eventually executes, a plan for doing it. This will involve finding out where Nell is, and making a navigation plan to get to her location. Assume that he knows where she is, and he is not too far away; then the fact that the plan will be carried out will be added to Dudley’s world model. Dudley must have some kind of database consistency maintainer to make sure that the plan is deleted if it is no longer necessary. Unfortunately, as soon as an apparently successful plan is added to the world model, the consistency maintainer will notice that “Nell is going to be mashed” is no longer true. But that removes any justification for the plan, so it goes too. But that means “Nell is going to be mashed” is no longer contradictory, so it comes back in. And so forth.

- Cohen & Levesque's intentions involve **commitment**. Having a commitment means having a **persistent goal**, viz., a goal the agent only abandons if s(he) comes to believe that the goal is fulfilled or unreachable. This is called **single-minded commitment**.
- Other forms of commitment:
 - **Blind commitment**: The agent maintains its intention until it is actually achieved.
 - **Open-minded commitment**: The agent maintains its intention as long as it is still believed possible. It may e.g. be rendered impossible by adapting new intentions.

- Agent i has the **persistent goal** that φ iff i has the achievement goal that φ and will keep that goal until it is either fulfilled or believed to be out of reach:

$$PGoal_i \varphi \stackrel{\text{def}}{=} AGoal_i \varphi \wedge (AGoal_i \varphi)U(Bel_i \varphi \vee Bel_i G \neg \varphi)$$



- $M, w_1 \models PGoal_i(atLecture)$

- Agent i has the **intention** that φ iff i has the persistent goal that φ and believes that (s)he can achieve φ by an action.

$$Intend_i \varphi \stackrel{\text{def}}{=} PGoal_i \varphi \wedge Bel_i F \exists \alpha Happ_{i:\alpha} \varphi$$

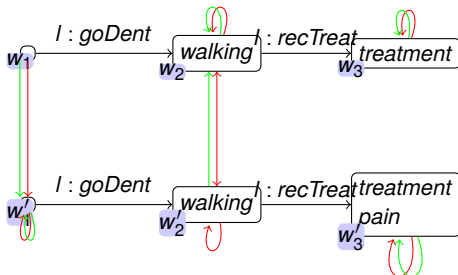
- **Intending is acting!** An agent 1 cannot intend that some other agent 2 does something. However, 1 may intend to make 2 do something.
- Viz., $Intend_1 Happ_{2:act} \top$ expands to $PGoal_1 Happ_{2:act} \top \wedge Bel_1 F \exists \alpha Happ_{1:\alpha} Happ_{2:act} \top$

- $\not\models (Intend_i \varphi \wedge Bel_i G(\varphi \rightarrow \psi)) \rightarrow Intend_i \psi.$

Proof

We provide a model for $Intend_i \varphi \wedge Bel_i G(\varphi \rightarrow \psi) \wedge \neg Intend_i \psi$:
John intends to go to the dentist. He believes that going to the dentist always implies pain. At the dentist, John gets some painkiller.

Dentist Example



- $M, w_1 \models \text{Intend}_I(\text{treatment}) \wedge \text{Bel}_I G(\text{treatment} \rightarrow \text{pain})$, but:
- $M, w_2 \not\models \text{AGoal}_I(\text{pain})$, thus:
- $M, w_1 \not\models \text{PGoal}_I(\text{pain})$, thus:
- $M, w_1 \not\models \text{Intend}_I(\text{pain})$

■ Sketch

- 1 Observe the execution of the system to be verified (e.g., log state of the environment, mental state of the agents, the agents' actions).
- 2 Represent the execution log using the semantics of Cohen & Levesque.
- 3 Model check representation against the agents' specification, e.g.:
 - $G(\text{goldNear} \rightarrow \text{Intend}(\text{hasGold}))$
 - $G(\text{Bel}(\text{goldNear}) \rightarrow \text{Intend}(\text{hasGold}))$
 - $G(\text{battLow} \rightarrow \text{Intend}(\neg \text{battLow}))$
- 4 Find time points where the specification evaluates false
⇒ Fault detection.

- We have studied an integrated logical framework that captures many aspects of agent behavior taking **belief** and **knowledge**, **preferences**, **goals**, and **intentions** into account, as well as how these mental attitudes change through **time** as progressed by **actions**.
- Next time, we'll look at another important notion, **obligations** and **permissions**, and we'll briefly discuss a **practical framework** (BOID) that deals with decision making in light of conflicts between beliefs, desires, intentions, and obligations.

- 1 Introduction
- 2 Agent-Based Simulation
- 3 Agent Architectures
- 4 Beliefs, Desires, Intentions
- 5 Norms and Duties
- 6 Communication and Argumentation
- 7 Coordination and Decision Making



M. Wooldridge, An Introduction to MultiAgent Systems, 2nd Edition, John Wiley & Sons, 2009.



Bratman, M. (1987). Intention, plans, and practical reason. Harvard University Press.



Cohen, P. R., & Levesque, H. J. (1990). Intention is choice with commitment. Artificial intelligence, 42(2-3), 213–261.



Meyer, J.-J. Ch., Broersen, J., Herzig, A. (2015). BDI Logics. In H. van Ditmarsch, J. Y. Halpern, W. van der Hoek, B. Kooi (Eds.) Handbook of Epistemic Logic. College Publications.