# Multi-Agent Systems

Albert-Ludwigs-Universität Freiburg

UNI
FREIBURG

Bernhard Nebel, Felix Lindner, and Thorsten Engesser
Winter Term 2018/19

# Recap and Overview

- Last time
  - The language of modal logics
  - Kripkean semantics for modal logics
  - Problems you are able to deal with:
    - Model Checking: Checking truth of formulas in a possible world in a Kripke model.
    - Theorem Proving: Checking validity of formulas w.r.t. a class of Kripke models.
- Today
  - The logic of knowledge (and belief)
    - Logical properties of knowledge (and belief)
    - Knowledge of groups of agents
    - Dynamics of knowledge and puzzles

# Motivation: Theory of Mind

- Sophisticated modes of social behavior require the ability to "put oneself in the position of someone else"
- Varieties of knowledge
    - Knowledge about others' knowledge:
        - First order: "John knows that the sun is shining"
        - Second order: "John knows that Mary knows that the sun is shining"
        - Third order: "John knows that Mary knows that Peter knows that the sun is shining"
        - …
    - Knowledge about one's own knowledge
        - Positive introspection: "I know that I know that the sun is shining"
        - Negative introspection: "I know that I don't know that the sun is shining"

Video: False Belief Task

# Knowledge and Belief

- Belief is the attitude of assent towards the truth of particular propositions.
- Knowledge is true justified belief (Plato).
    - Justification: Evidence, or support, for your belief. I.e., if you just claim some truth without evidence, this does not count as knowledge.

- This definition is challenged by philosophical arguments (viz., by so-called Gettier cases).
- Standard epistemic logic is much more pragmatic, though: $\approx$ knowledge is true belief.

# Operators for Knowledge and Belief

- $K_i \varphi$: Agent $i$ knows $\varphi$.
- $B_i \varphi$: Agent $i$ believes $\varphi$.

- Some authors also write $K(i)(p)$ or $iKp$ resp. $B(i)(p)$ or $iBp$..
- One distinguishes the logic of knowledge (epistemic logic) from the logic of belief (doxastic logic). We will discuss empistemic logic and refer to doxastic logic when the distinction is interesting to us.

# Epistemic Alternatives

### Def. Epistemically Accessible, Epistemic Alternative

A particular world $w'$ is epistemically accessible to an agent $i$ in world $w$ iff the set of all propositions $p$ that agent $i$ knows in $w$ are compatible with all true propositions in $w'$. All such worlds $w'$ are considered epistemic alternatives.

- We also say that the epistemic alternatives are epistemically indistinguishable to the agent.

# Knowing and Believing: Semantics

Agent $i$ knows $\varphi$ in world $w$ iff. $\varphi$ is true in all worlds $w'$ epistemically accessible from $w$ for $i$:

- $M, w \models K_i \varphi$ iff. for all $w'$ s.th. $(w, w') \in R(K_i), M, w' \models \varphi$

Agent $i$ believes $\varphi$ in world $w$ iff. $\varphi$ is true in all worlds $w'$ doxastically accessible from $w$ for $i$:

- $M, w \models B_i \varphi$ iff. for all $w'$ s.th. $(w, w') \in R(B_i), M, w' \models \varphi$
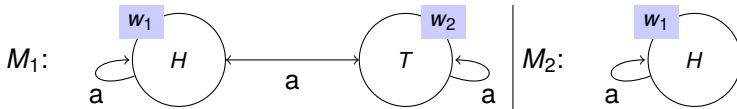
Additionally:

- $M, w \models \hat{K}_i \varphi$ iff. $M, w \models \neg K_i \neg \varphi$.
- $M, w \models \hat{B}_i \varphi$ iff. $M, w \models \neg B_i \neg \varphi$.

## Quote [Hintikka & Halonen 98]

When you know that S, you can legitimately omit from consideration all possibilities under which it is not the case that S. In other words you can restrict your attention to the situations in which it is true that S.



- $M_1, w_1 \models \neg K_a H \wedge \neg K_a T$
- $M_2, w_1 \models K_a H$

- Some sessions ago we agreed on classifying knowledge as **S5** modality, viz.,:
    - K-axiom: $K_i(\varphi \rightarrow \psi) \rightarrow (K_i \varphi \rightarrow K_i \psi)$
    - T-axiom: $K_i \varphi \rightarrow \varphi$
    - 4-axiom: $K_i \varphi \rightarrow K_i K_i \varphi$
    - 5-axiom: $\neg K_i \varphi \rightarrow K_i \neg K_i \varphi$
- T was considered inappropriate for belief. Instead, the weaker axiom D was considered appropriate:
    - D-axiom: $B_i \varphi \rightarrow \neg B_i \neg \varphi$

# Axiom K

$K_i(\varphi \rightarrow \psi) \rightarrow (K_i \varphi \rightarrow K_i \psi)$

- Objection: Logical Omniscience
  - The agent knows all implications of its knowledge.
  - The agent knows all tautologies.
- Is there a way out?
  - Not if we stick to normal modal logics. Approaches exist based on generalizations of Kripke models (impossible worlds). $\Rightarrow$ We will not go into detail. If you are interested, see Fagin et al., 1995 [2].

# Axiom N

- Necessity: If $\models \varphi$ then $\models K_i \varphi$

# Axiom T

$K_i \varphi \rightarrow \varphi$

- Uncontroversial for epistemic logic. Indeed, **KT** is considered weakest system for knowledge.

- $K_a K_b \varphi \models_{\text{KT}} K_a \varphi$, called transmissibility of knowledge. You know everything you know that others know.
    - $K_b \varphi \rightarrow \varphi$ (Axiom T)
    - $K_a(K_b \varphi \rightarrow \varphi)$ (Axiom N)
    - $K_a(K_b \varphi \rightarrow \varphi) \rightarrow (K_a K_b \varphi \rightarrow K_a \varphi)$ (K-Axiom)
    - $K_a K_b \varphi \rightarrow K_a \varphi$ (MP)  □

- Do you also believe everything you believe that others believe? ⇒Doxastic: Uncontroversial that it does not hold. Weaker axiom D instead.

$K_i\varphi \rightarrow \neg K \neg \varphi$

- Epistemic: Axiom T already entails axiom D.
- Doxastic: Considered as a substitute for T. Psychological objection: People do hold contradictory beliefs.

# Axiom 4

$K_i \varphi \rightarrow K_i K_i \varphi$

- Doxastic: Uncontroversial.
- Epistemic:
    - Objection I: So-called $KK$-regress: Agents are required to have infinitely nested knowledge.
    - Objection II: A corollary of the axiom is $\neg K_i K_i \varphi \rightarrow \neg K_i \varphi$ (counterposition). But it seems possible to have knowledge that comes into ones mind only after one gets some hints.

However, as Hintikka notes, "Knowing to know differs only in words from knowing." In system **KT4** we have $K_i K_i \varphi \rightarrow K_i \varphi$, $K_i \varphi \rightarrow K_i K_i \varphi$, hence $K_i K_i \varphi \leftrightarrow K_i \varphi$.

# Axiom 5

$\neg K_i \varphi \to K_i \neg K_i \varphi$

- Objection: A corollary is $K_i \varphi \vee K_i \neg K_i \varphi$ ($\equiv_{Def.}$ awareness). No room for ignorance! Consider John living in 17th century: According to axiom 5, John either knew Einstein's theory of relativity, or he knew that he does not know Einstein's theory of relativity. In fact, it is more appropriate to claim that John was not aware.

# Axiom B

From reflexivity (T) and Euclideanness (5) follows symmetry (B):
$\neg\varphi \to K_i\neg K_i\varphi$ (Proof: $\neg\varphi\Rightarrow_T\neg K_i\varphi\Rightarrow_5 K_i\neg K_i\varphi$ $\square$)

- Objection: What is actually true must be known to be possibly true.
- Objection: A corollary of axiom B is $\neg K_i\neg K\varphi \to \varphi$, which is the same as $\hat{K}K\varphi \to \varphi$. In words: Only true things are considered possible to be known resp. believed.

As this seems too strong to most epistemologists, 5 (and thus B) is often rejected. $\Rightarrow$Alternative Axioms proposed: 4.2, 4.3, 4.4 However, in many computer science applications, B is considered appropriate: If the agent cannot find $\varphi$ in its database, it can conclude that it knows that it does not know $\varphi$ (closed-world assumption).

$\varphi \rightarrow (\neg K_i \neg K_i \varphi \rightarrow K_i \varphi)$

- A corollary is $\neg \varphi \vee K_i \neg K_i \varphi \vee K_i \varphi$. Thus, the agent is not required to know something about things that do not hold.

- In **KT4.4**: If $\varphi$ is not the case, then $\neg K_i \varphi$ (by T), but not $K_i \neg K_i \varphi$. If $\varphi$ is the case, then either it knows $\varphi$ or it knows that it does not know $\varphi$.

$K_i(K_i\varphi \rightarrow K_i\psi) \vee K_i(K_i\psi \rightarrow K_i\varphi)$

- Variant: $K_i(\neg K_i\varphi \vee K_i\psi) \vee K_i(\neg K_i\psi \vee K_i\varphi)$
- $K_i(\neg(K_i\varphi \wedge \neg K_i\psi)) \vee K_i(\neg(K_i\psi \wedge \neg K_i\varphi))$ (De Morgan)
- The agent has introspection regarding the fact that what it knows and does not know is consistent. E.g., it knows that it cannot at the same time know that $\varphi$ and not know that $\varphi$.

$\neg K_i \neg K_i \varphi \rightarrow K_i \neg K_i \neg \varphi$

- Corollary: $K_i \hat{K}_i \neg \varphi \vee K_i \hat{K} \varphi$. Agents are only required to know whether they consider some proposition or its negation possible.

# Mixing Knowledge and Belief

Building a logic with both Knowledge and Belief modalities interacting is non-trivial. Consider interaction axioms:

- Entailment property: $K_i \varphi \rightarrow B_i \varphi$
- Positive certainty property: $B_i \varphi \rightarrow B_i K_i \varphi$
- And let B be a **KD45** modality and K a **S5** modality. Then $\neg B_i \varphi \rightarrow B_i \neg K_i \varphi$ holds (seems reasonable):
  - $\neg B_i \varphi \Rightarrow_{\textbf{Entailment}} \neg K_i \varphi \Rightarrow_{\textbf{5}} K_i \neg K_i \varphi \Rightarrow_{\textbf{Entailment}} B_i \neg K_i \varphi$ $\quad \square$

- Objection
  - Let $p$ be a proposition the agent believes, but in fact $p$ is false: $B_i p \land \neg p$
  - $\Rightarrow_{\textbf{Positive Certainty}} B_i K_i p$
  - $\Rightarrow_{\textbf{T}} \neg K_i p \Rightarrow_{\textbf{5}} K_i \neg K_i p \Rightarrow_{\textbf{Entailment}} B_i \neg K_i p$ ↯

# Mixing Knowledge and Belief

Building a logic with both Knowledge and Belief modalities interacting is non-trivial. Consider interaction axioms:

- Entailment property: $K_i \varphi \to B_i \varphi$
- Positive certainty property: $B_i \varphi \to B_i K_i \varphi$
- And let B be a **KD45** modality and K a **S5** modality. Then $\neg B_i \varphi \to B_i \neg K_i \varphi$ holds (seems reasonable):
  - $\neg B_i \varphi \Rightarrow_{\textbf{Entailment}} \neg K_i \varphi \Rightarrow_{\textbf{5}} K_i \neg K_i \varphi \Rightarrow_{\textbf{Entailment}} B_i \neg K_i \varphi \quad \square$

- Objection
  - Let $p$ be a proposition the agent believes, but in fact $p$ is false: $B_i p \wedge \neg p$
  - $\Rightarrow_{\textbf{Positive Certainty}} B_i K_i p$
  - $\Rightarrow_{\textbf{T}} \neg K_i p \Rightarrow_{\textbf{5}} K_i \neg K_i p \Rightarrow_{\textbf{Entailment}} B_i \neg K_i p \quad \lightning$

# Group Knowledge

- Extend language by three modal operators:
  - $E_G\varphi$: Everyone in the group $G$ knows $\varphi$.
  - $C_G\varphi$: It is common knowledge among the agents in $G$ that $\varphi$ is the case.
  - $D_G\varphi$: It is distributed knowledge among the agents in $G$ that $\varphi$ is the case.

- Example
  - $K_3 C_{\{1,2\}} p$: Agent 3 knows that it is common knowledge among agents 1 and 2 that $p$ is the case.

# Everyone Knows

- $M, w \models E_G \varphi$ iff. $M, w \models K_i \varphi$ for all $i \in G$.
- Write $E_G^0 \varphi$ as an abbreviation of $\varphi$, $E_G^1 \varphi$ as an abbreviation for $K_1 \varphi \wedge K_1 \varphi \wedge ...$, and let $E_G^{k+1} \varphi$ be an abbreviation for $E_G E_G^k \varphi$.

## Def. *G*-reachable in *k* steps

A world $w'$ is *G*-reachable from world $w$ in $k \geq 1$ steps iff. there exists worlds $u_0 u_1 \ldots u_k$ such that $u_0 = w$ and $u_k = w'$ and for all $j$ with $0 \leq j \leq k - 1$ there exists $i \in G$ s.th. $(u_j, u_{j+1}) \in R(K_i)$.

## Lemma

$M, w \models E_G^k \varphi$ iff. $M, w' \models \varphi$ for all $w'$ that are *G*-reachable from $w$ in $k$ steps.

# Everyone Knows

- $M, w \models E_G \varphi$ iff. $M, w \models K_i \varphi$ for all $i \in G$.
- Write $E_G^0 \varphi$ as an abbreviation of $\varphi$, $E_G^1 \varphi$ as an abbreviation for $K_1 \varphi \wedge K_1 \varphi \wedge ...$, and let $E_G^{k+1} \varphi$ be an abbreviation for $E_G E_G^k \varphi$.

### Def. $G$-reachable in $k$ steps

A world $w'$ is *G-reachable* from world $w$ in $k \geq 1$ steps iff. there exists worlds $u_0 u_1 \ldots u_k$ such that $u_0 = w$ and $u_k = w'$ and for all $j$ with $0 \leq j \leq k-1$ there exists $i \in G$ s.th. $(u_j, u_{j+1}) \in R(K_i)$.

### Lemma

$M, w \models E_G^k \varphi$ iff. $M, w' \models \varphi$ for all $w'$ that are *G-reachable* from $w$ in $k$ steps.

# Common Knowledge

- $M, w \models C_G \varphi$ iff. $M, w \models E_G^k \varphi$ for $k = 1, 2, \ldots$.

## Def. G-reachable

A world $w'$ is G-reachable from world $w$ iff. $w'$ is $G$-reachable from $w$ in $k$ steps for some $k \geq 1$.

## Lemma

$M, w \models C_G \varphi$ iff. $M, w' \models \varphi$ for all $w'$ that are $G$-reachable from $w$.

# Common Knowledge

- $M, w \models C_G \varphi$ iff. $M, w \models E_G^k \varphi$ for $k = 1, 2, \ldots$.

## Def. G-reachable

A world $w'$ is G-reachable from world $w$ iff. $w'$ is $G$-reachable from $w$ in $k$ steps for some $k \geq 1$.

## Lemma

$M, w \models C_G \varphi$ iff. $M, w' \models \varphi$ for all $w'$ that are $G$-reachable from $w$.
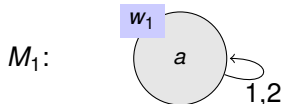
# Distributed Knowledge

## Distributed Knowledge

A group of agents $G$ has distributed knowledge of $\varphi$ iff the combined knowledge of the members of $G$ implies $\varphi$. Idea: Eliminate all worlds that some agent in $G$ considers impossible. Technically: Intersect the sets of worlds each agent in $G$ considers possible. Hence:

- $M, w \models D_G \varphi$ iff. $M, w' \models \varphi$ for all $w'$ s.th. $(w, w') \in \cap_{i \in G} K_i$
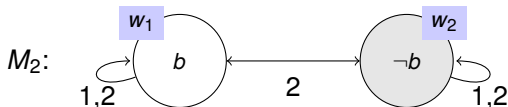
- First, it is common knowledge that the ball is in the box

$M_1$:



$$M_1, w_1 \models E_{\{1,2\}}b \wedge C_{\{1,2\}}b \wedge D_{\{1,2\}}b$$

- Afterwards, the ball is not in the box. Agent 1 knows, Agent 2 does not.

$M_2$:



$$M_2, w_2 \models \neg E_{\{1,2\}}\neg b \wedge \neg C_{\{1,2\}}\neg b \wedge D_{\{1,2\}}\neg b$$

# Summary and Outlook

- Today
    - Critical discussion of the formal properties of knowledge and belief as modeled in modal logics. **S5** seems to be quasi-standard model for knowledge in computer science, there may be good reasons to make other choices for specific modeling purpose.
    - Foundations of group knowledge: Everyone knows, Common knowledge, and distributed knowledge.
- Next time: Dynamics and Puzzles

# Literature

📄 J. Garson, 2004, https://plato.stanford.edu/entries/logic-modal/

📄 R. Fagin, J. Y. Halpern, Y. Moses, M. Y. Vardi, Reasoning about knowledge, The MIT Press, 1995.

📄 J. Symons, Logic and formal semantics for epistemology, in The Routledge Companion to Epistemology (Routledge Philosophy Companions) Duncan Pritchard and Sven Bernecker (eds.), New York: Routledge. 571–586, 2011.