

Dynamic Epistemic and Doxastic Logics

Sonja Smets, ILLC, Amsterdam

(Slides are based on joint lectures with A. Baltag, ILLC)

Financial Support Acknowledgement:

European Research Council



Netherlands Organisation for Scientific Research

PLAN OF THIS COURSE

1. **Puzzles. Logics of knowledge and belief. Epistemic and Doxastic models.**
2. **Core of Standard (“Hard”) Dynamic-Epistemic Logic:** Public and Private announcements. Event models. The Product Update Mechanism.
3. **Belief Revision:** Plausibility Models. Conditional belief. Belief Upgrades. Doxastic event models and the Action-Priority Rule.
4. **Further Topics in the last three lectures:** Iterated Belief Revision. Belief Merge. Collective Learning. Informational Cascades. Surprise Examination Paradox etc.

Relevant Textbooks and Surveys

- A. Baltag, H. P. van Ditmarsch and L.S. Moss, “Epistemic logic and information update”, in *Handbook of Philosophy of Information*, Elsevier, 2008.
- A. Baltag and S. Smets, “A Qualitative Theory of Dynamic Interactive Belief Revision”, in G. Bonanno, W. van der Hoek, M. Wooldridge (eds.), *Texts in Logic and Games*, Vol 3, pp.9-58, Amsterdam Univ Press, 2008.
- J. van Benthem, **Modal Logic for Open Minds**, CSLI Publications, Stanford, 2011.
- J. van Benthem, **Logical Dynamics of Information and Interaction**, Cambridge Univ Press, 2011.
- H. P. van Ditmarsch, W. van der Hoek and B. Kooi, **Dynamic Epistemic Logic**, Springer, 2007.

- R. Fagin, J.Y. Halpern, Y. Moses and M.Y. Vardi, **Reasoning about Knowledge**, MIT Press, Cambridge MA 1995.
- Research papers of van Benthem, Baltag and Smets (see their personal websites)

1.1 Epistemic Puzzles: Muddy Children

Suppose there are 4 children, all of them being good logicians, exactly 3 of them having dirty faces. *Each can see the faces of the others, but doesn't see his/her own face.*

The father publicly announces:

“At least one of you is dirty”.

Then the father does another paradoxical thing: *starts repeating over and over the same question* **“Do you know if you are dirty or not, and if so, which of the two?”**

After each question, the children have to *answer publicly, sincerely and simultaneously, based only on their knowledge, without taking any guesses*. No other communication is allowed and nobody can lie.

One can show that, after 2 rounds of questions and answers, **all the dirty children will come to know they are dirty!** So they give this answer in the 3rd round, after which **the clean child also comes to know she's clean**, giving the correct answer at the 4th round.

Muddy Children Puzzle continued

First Question: *What's the point of the father's first announcement ("At least one of you is dirty")?*

Apparently, this message is not informative to any of the children: the statement was already known to everybody! But the puzzle wouldn't work without it: in fact this announcement adds information to the system! The children implicitly learn some new fact, namely the fact that what each of them used to know *in private* is now *public knowledge*.

Second Question: *What's the point of the father's repeated questions?*

If the father knows that his children are good logicians, then at each step the father knows already the answer to his question,

before even asking it! However, the puzzle wouldn't work without these questions. In a way, it seems the father's questions are “*abnormal*”, in that they don't actually aim at filling a gap in father's knowledge; but instead they are part of a *Socratic strategy of teaching-through-questions*.

Third Question: *How can the children's statements of ignorance lead them to knowledge?*

Puzzle no 2: Sneaky Children

Let us modify the last example a bit.

Suppose the children are somehow rewarded for answering as quickly as possible, but they are punished for incorrect answers; thus they are interested in getting to the correct conclusion as fast as possible.

Suppose also that, after the first round of questions, two of the dirty children “cheat” on the others by secretly announcing each other that they’re dirty, while none of the others suspects this can happen.

Honest Children Always Suffer

One can easily see that the **third dirty child will be totally deceived, coming to the “logical” conclusion that... she is clean!**

So, after giving the wrong answer, she ends up by being punished for her credulity, despite her impeccable logic.

Clean Children Always Go Crazy

What happens to the clean child?

Well, **assuming she doesn't suspect any cheating, she is facing a contradiction**: two of the dirty children answered too quickly, coming to know they're dirty before they were supposed to know!

*If the third child simply updates her knowledge monotonically with this new information (and uses classical logic), then she ends up believing everything: **she goes crazy!***

The Amazon Island

This is another story encoding the same puzzle:

On the island of Amazonia, women are dominant and the law says that, if at any point a woman knows her husband is cheating on her, she must shoot him the same day at noon in the main square.

Now the queen (truthfully) tells the women: “At least one of your husbands is a cheater. Whenever somebody’s husband is cheating, all the other women know it”.

For 16 days, nothing happens. Then, in the 17th day, shootings are heard.

Question: How many husbands died?

The Dangers of Mercy

In the Amazonia version of the story, assume that again there are exactly 17 cheating husbands (out of 1 million husbands on the island), while the rest of 999.983 husbands are faithful.

Consider what happens now if *all the wives of the 17 cheating husbands secretly decide to break the Queen's rules*, by quietly sparing the lives of their husbands, even when they get to know that they are cheating.

We also assume that *all the other wives do not suspect this*: not only that *they strictly obey by the Queen's rules*, but *they believe that it is common knowledge that everybody else obeys by those same rules*.

It's easy to see that, in this case, *17 days will pass without any shooting.*
But it's also easy to show that in the 18th day, shots will be heard.

How many husbands will die in this scenario? How many of these are innocent?

Surprised Children

The students in a high-school class **know for sure that the date of the exam has been fixed in one of the five (working) days** of next week: it'll be the last week of the term, and it's got to be an exam, and only one exam.

But **they don't know in which day**.

Now the Teacher announces her students that the **exam's date will be a surprise**: i.e. *even in the evening before the exam, the students will still not be sure that the exam is tomorrow.*

Paradoxical Argumentation

Intuitively, one can prove (by backward induction, starting with Friday) that, **IF this announcement is true, then the exam cannot take place in any day of the week.**

So, using this argument, the students come to “know” that **the announcement is false**: the exam CANNOT be a surprise.

GIVEN THIS, they feel entitled to **dismiss** the announcement, and...
THEN, **surprise**: *whenever the exam will come* (say, on Tuesday), it **WILL** indeed be a *complete surprise!*

1.2. Epistemic-Doxastic Models and Logics

Epistemic Logic was first formalized by Hintikka (1962), who also sketched the first steps in formalizing doxastic logic.

They were further developed and studied by both philosophers (Parikh, Stalnaker etc.), economists (Aumann) and computer-scientists (Halpern, Vardi, Fagin etc.)

Syntax of Single-Agent Epistemic-Doxastic Logic

$$\varphi ::= p \mid \neg\varphi \mid \varphi \wedge \varphi \mid K\varphi \mid B\varphi$$

Models for Single-Agent Information

We are given a set of “possible worlds”, meant to represent **all the relevant epistemic/doxastic possibilities** in a certain situation.

EXAMPLE 1: a coin is on the table, but the (implicit) agent **doesn't know (nor believe he knows)** which face is up.

H

T

Knowledge or Belief

The **universal quantifier over the domain of possibilities** is interpreted as **knowledge, or belief**, by the implicit agent.

So we say the agent **knows, or believes**, a sentence φ if φ is **true in all the possible worlds** of the model.

The specific interpretation (knowledge or belief) depends on the context.

In the previous example, the agent doesn't know (nor believe) that the coin lies Heads up, and neither that it lies Tails up.

Learning: Update

EXAMPLE 2:

Suppose now **the agent looks at the coin and he sees it's Heads up.**

The model of the new situation is now:

H

Only one epistemic possibility has survived: the agent now **knows/believes that the coin lies Heads up.**

Update as World Elimination

In general, **updating corresponds to world elimination:**

an **update with a sentence φ** is simply the operation of **deleting all the non- φ possibilities**

After the update, the worlds not satisfying φ are no longer possible: the actual world is known not to be among them.

Truth and Reality

But is φ “**really**” true (in the “real” world), apart from the agent’s knowledge or beliefs?

For this, we need to specify which of the possible worlds is **is the actual world**.

Real World

Suppose that, in the original situation (before learning), the coin lied Heads up indeed (though the agent didn't know, or believe, this).

We represent this situation by marking **the actual (“real” state of the) world with a red star:**

* H

T

Mistaken Updates

But what if the real world is not among the “possible” ones? What if the agent’s sight was so bad that she only **thought** she saw the coin lying Heads up, when **in fact it lied Tails up**?

After the “update”, her epistemically-possible worlds are just

H

but we **cannot** mark the actual world here, since it **doesn’t belong to the agent’s model!**

False Beliefs

Clearly, in this case, the model only represents the agent's beliefs, but NOT her "knowledge" (in any meaningful sense): the agent believes that the coin lies Heads up, but this is wrong!

Knowledge is usually assumed to be truthful, but in this case the agent's belief is false.

But still, **how can we talk about "truth" in a model** in which the actual world is not represented?!

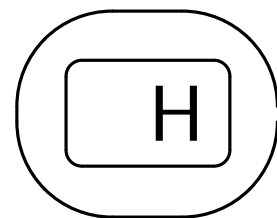
Third-person Models

The solution is to go beyond the agent's own model, by taking an “objective” (third-person) perspective: the real possibility is always in the model, even if the agent believes it to be impossible.

To point out which worlds are **believed to be possible** by the agent we **encircle them**: these worlds form the “**sphere of beliefs**”.

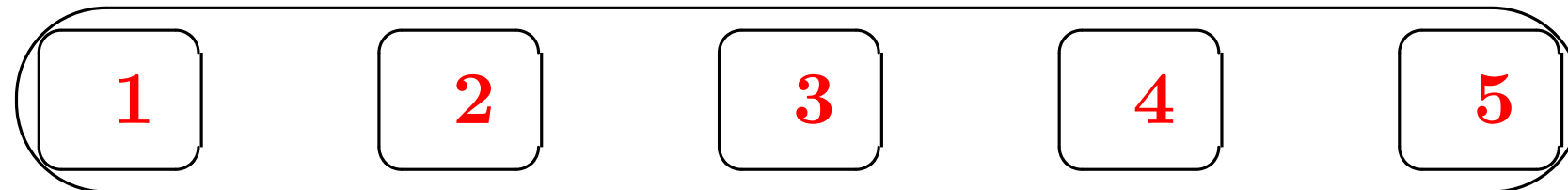
“*Belief*” now **quantifies ONLY** over the worlds in this sphere, while “*knowledge*” still quantifies over **ALL** possible worlds.

EXAMPLE 3:



Example 4

In the Surprise Exam story, a possible initial situation (BEFORE the Teacher's announcement) might be given by:



where i means that: the exam takes place in the i -th (working) day of the week.

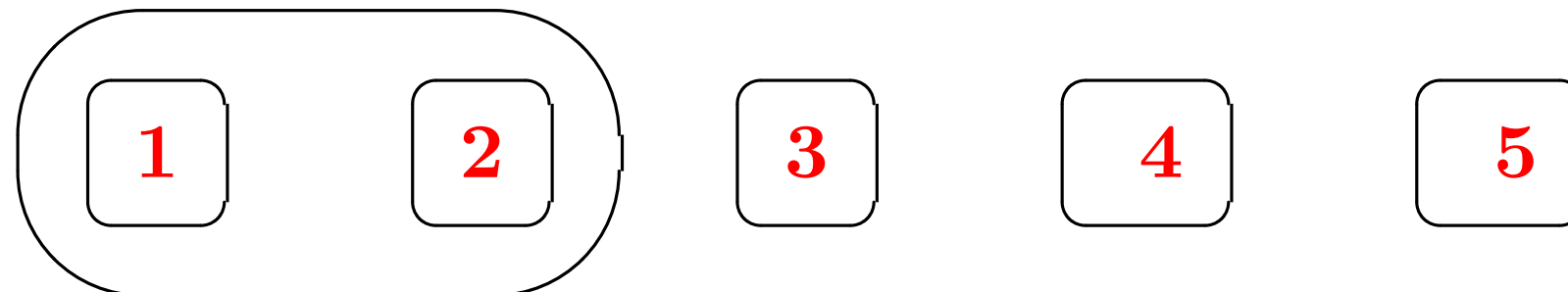
This encodes an initial situation in which the student **knows that there will be an exam** in (exactly) one of the days, but he **doesn't know the day**, and moreover he **doesn't have any special belief about this**: he considers all days as being *possible*.

We are not told when will the exam take place: *no red star*.

Beliefs

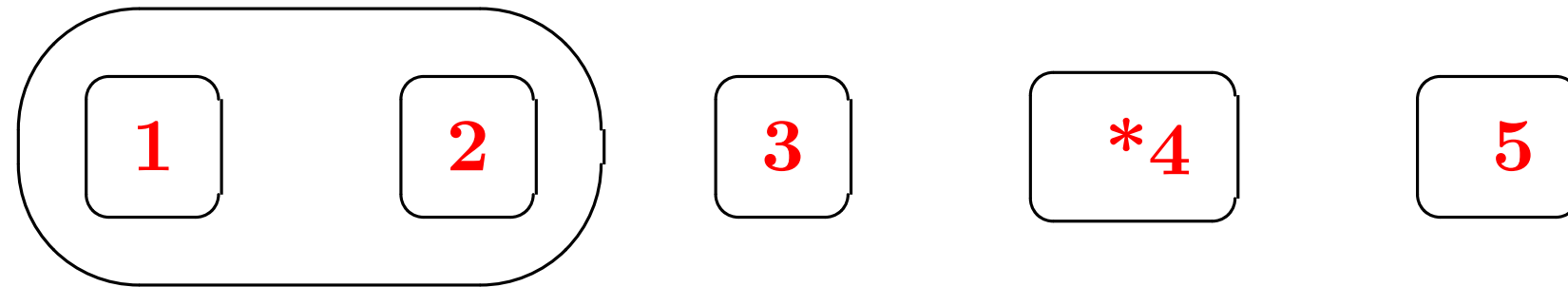
EXAMPLE 5:

If however, the Student **believes** (for some reason or another) that the exam will take place either Monday or Tuesday, then the correct representation is:



Again, we are not told when is the exam, so no red star.

However, if we are told that the exam is in fact on Thursday (though the student still doesn't know this), then the model is:



In this model, **some of the student's beliefs are false**, since the real world does NOT belong to his "sphere of beliefs".

Simple Models for Knowledge and Belief

For a set Φ of *facts*, a **(single-agent, pointed) epistemic-doxastic model** is a structure:

$$\mathbf{S} = (S, S_0, \|\cdot\|, s_*) , \quad \text{consisting of:}$$

1. A set S of “*possible worlds*” (or possible “states of the world”, also known as “*ontic states*”). S defines the agent’s **epistemic state**: these are the states that are “*epistemically possible*”.
2. A *non-empty* subset $S_0 \subseteq S, S_0 \neq \emptyset$, called the “*sphere of beliefs*”, or the agent’s **doxastic state**: these are the states that “doxastically possible”.
3. A map $\|\cdot\| : \Phi \rightarrow \mathcal{P}(S)$, called the **valuation**, assigning to each $p \in \Phi$ a set $\|p\|_S$ of states.

4. A designated world $s_* \in S$, called the “**actual world**”.

Interpretation

- The epistemic state S gives us an (implicit) agent's **state of knowledge**: he **knows** the real world belongs to S , but *cannot distinguish* between the states in S , so cannot know which of them is the real one.
- The doxastic state S_0 gives us the agent's **state of belief**: he **believes** that the real world belongs to S_0 , but his beliefs are consistent with *any* world in S_0 .
- The valuation tells us **which ontic facts hold in which world**: we say that p is **true** at s if $s \in \llbracket p \rrbracket$.
- The actual world s_* gives us the “**real state**” of the world: what really is the case.

Truth

For any world w in a model \mathbf{S} and any sentence φ , we write

$$w \models_{\mathbf{S}} \varphi$$

if φ is **true** in the world w .

When the model \mathbf{S} is fixed, we skip the subscript and simply write

$$w \models \varphi.$$

For **atomic sentences**, this is given by the **valuation map**:

$$w \models p \text{ iff } w \in \|p\|,$$

while for other propositional formulas is given by **the usual truth clauses**:

$$w \models \neg\varphi \text{ iff } w \not\models \varphi,$$

$$w \models \varphi \wedge \psi \text{ iff } w \models \varphi \text{ and } w \models \psi.$$

Disjunction, Conditional, Biconditional: We take $\varphi \vee \psi$ to be just an abbreviation for $\neg(\neg\varphi \wedge \neg\psi)$, $\varphi \Rightarrow \psi$ to be just an abbreviation for $\neg\varphi \vee \psi$, and $\varphi \Leftrightarrow \psi$ to be an abbreviation for $(\varphi \Rightarrow \psi) \wedge (\psi \Rightarrow \varphi)$.

As a consequence, we have e.g:

$$w \models_{\mathbf{s}} \varphi \vee \psi \text{ iff either } w \models_{\mathbf{s}} \varphi \text{ or } w \models_{\mathbf{s}} \psi$$

etc.

Interpretation Map

We can extend the valuation $\|p\|_{\mathbf{s}}$ to an **interpretation map** $\|\varphi\|_{\mathbf{s}}$ for all propositional formulas φ :

$$\|\varphi\|_{\mathbf{s}} := \{w \in S : w \models_{\mathbf{s}} \varphi\}.$$

Obviously, this has the property that

$$\|\neg\varphi\|_{\mathbf{s}} = S \setminus \|\varphi\|_{\mathbf{s}},$$

$$\|\varphi \wedge \psi\|_{\mathbf{s}} = \|\varphi\|_{\mathbf{s}} \cap \|\psi\|_{\mathbf{s}},$$

$$\|\varphi \vee \psi\|_{\mathbf{s}} = \|\varphi\|_{\mathbf{s}} \cup \|\psi\|_{\mathbf{s}}.$$

We now want to extend the interpretation to all the sentences in doxastic-epistemic logic.

Knowledge and Belief

Knowledge is defined as “**truth in all epistemically possible worlds**”, while **belief** is “**truth in all doxastically possible worlds**”

Formally:

$$w \models K\varphi \text{ iff } t \models \varphi \text{ for all } t \in S,$$
$$w \models B\varphi \text{ iff } t \models \varphi \text{ for all } t \in S_0.$$

Validity and Satisfiability

A sentence is **valid** over epistemic-doxastic models if it is true at every state in every epistemic-doxastic model.

A sentence is **satisfiable** (over epistemic-doxastic models) if it is true some state in some epistemic-doxastic model.

Consequences

For every sentence φ, ψ etc, the following are valid over epistemic-doxastic models:

1. **Veracity of Knowledge:**

$$K\varphi \Rightarrow \varphi$$

2. **Positive Introspection of Knowledge:**

$$K\varphi \Rightarrow KK\varphi$$

3. **Negative Introspection of Knowledge:**

$$\neg K\varphi \Rightarrow K\neg K\varphi$$

4. **Consistency of Belief:**

$$\neg B(\varphi \wedge \neg\varphi)$$

5. Positive Introspection of Belief:

$$B\varphi \Rightarrow BB\varphi$$

6. Negative Introspection of Belief:

$$\neg B\varphi \Rightarrow B\neg B\varphi$$

7. Strong Positive Introspection of Belief:

$$B\varphi \Rightarrow KB\varphi$$

8. Strong Negative Introspection of Belief:

$$\neg B\varphi \Rightarrow K\neg B\varphi$$

9. Knowledge implies Belief:

$$K\varphi \Rightarrow B\varphi$$

Epistemic-Doxastic Logic: Sound and Complete Proof System

In fact, a **sound and complete proof system for single-agent epistemic-doxastic logic** can be obtained by taking as axioms: **validities (1)-(4) and (7)-(9) above**, together with **all propositional tautologies**, as well as “**Kripke’s axioms**” for **knowledge and belief**

$$K(\varphi \Rightarrow \psi) \Rightarrow (K\varphi \Rightarrow K\psi),$$

$$B(\varphi \Rightarrow \psi) \Rightarrow (B\varphi \Rightarrow B\psi),$$

and together with following *inference rules*:

Modus Ponens: From φ and $\varphi \Rightarrow \psi$ infer ψ .

Necessitation: From φ infer $K\varphi$.

Generalization

Many philosophers deny that knowledge is introspective, and some philosophers deny that belief is introspective. In particular, both common usage and Platonic dialogues suggest that people **may believe they know things that they don't actually know.**

Other of the above validities may also be debatable: e.g. some “crazy” agents may have inconsistent beliefs.

So it is convenient to have a more general semantics, in which the above principles do not necessarily hold, so that one can pick whichever principles one considers true.

Kripke Semantics: multi-agent

For a set Φ of *facts* and a finite set \mathcal{A} of *agents*, a **Φ -Kripke model** is a structure

$$\mathbf{S} = (S, \overset{\mathcal{A}}{\rightarrow}, \|\cdot\|, s_*)$$

consisting of

1. a set S of “*possible worlds*”
2. a family of **binary accessibility relations** $\overset{a}{\rightarrow} \subseteq S \times S$, one for each agent $a \in \mathcal{A}$
3. and a *valuation* $\|\cdot\| : \Phi \rightarrow \mathcal{P}(S)$, assigning to each $p \in \Phi$ a set $\|p\|_{\mathbf{S}}$ of worlds
4. a designated world s_* : the “*actual*” one.

- The valuation is also called a *truth map*. It is meant to express the *factual content* of a given world.
- The arrows (accessibility relations) \xrightarrow{A} express the agents' uncertainty between various worlds.
- A Kripke model is called a **state model** whenever we think of its “worlds” as *possible states*. In this case, the elements $p \in \Phi$ are called *atomic sentences*, being meant to represent **basic “ontic” (non-epistemic) facts**, which may hold or not at a given state.

Satisfaction Relation

Write $s \models_{\mathbf{S}} \varphi$ for the **satisfaction relation**: φ is true at world s in model \mathbf{S} . This is defined inductively:

$$s \models_{\mathbf{S}} p \text{ iff } s \in \llbracket p \rrbracket_{\mathbf{S}}$$

$$s \models_{\mathbf{S}} \neg\varphi \text{ iff } s \not\models_{\mathbf{S}} \varphi$$

$$s \models_{\mathbf{S}} \varphi \wedge \psi \text{ iff } s \models_{\mathbf{S}} \varphi \text{ and } s \models_{\mathbf{S}} \psi$$

Extending the Truth Map

Equivalently, this allows us to *extend the truth map* $\|\varphi\|_{\mathbf{s}}$ to *all* propositional formulas, by putting:

$$\|\varphi\|_{\mathbf{s}} := \{s \in S : s \models_{\mathbf{s}} \varphi\}.$$

Obviously, this has the property that

$$\|\neg\varphi\|_{\mathbf{s}} = S \setminus \|\varphi\|_{\mathbf{s}},$$

$$\|\varphi \wedge \psi\|_{\mathbf{s}} = \|\varphi\|_{\mathbf{s}} \cap \|\psi\|_{\mathbf{s}},$$

$$\|\varphi \vee \psi\|_{\mathbf{s}} = \|\varphi\|_{\mathbf{s}} \cup \|\psi\|_{\mathbf{s}}.$$

Any *new* propositional operator $A(\varphi_1, \dots, \varphi_n)$ is “*defined*” by *extending the truth map* to define $\|A(\varphi_1, \dots, \varphi_n)\|_{\mathbf{s}}$, i.e. by *giving a defining inductive clause for satisfaction* $s \models A(\varphi_1, \dots, \varphi_n)$.

Modalities

For every sentence φ , we can define a sentence $\Box_a\varphi$ by (universally) quantifying over accessible worlds:

$$s \models_{\mathbf{S}} \Box_a\varphi \text{ iff } t \models_{\mathbf{S}} \varphi \text{ for all } t \text{ such that } s \xrightarrow{a} t.$$

Its *existential dual*

$$\Diamond_a\varphi := \neg\Box_a\neg\varphi$$

denotes a sense of “**epistemic/doxastic possibility**”.

Kripke Models for Knowledge and Belief

From now on, new notational convention:

In a context when we interpret a modality $\Box_a\varphi$ as **knowledge**, we use the notation $K_a\varphi$ instead, and we denote by \sim_a the underlying binary accessibility relation.

When we interpret the modality $\Box_a\varphi$ as **belief**, we use the notation $B_a\varphi$ instead, and we use the (long arrow) notation \xrightarrow{a} for the underlying binary doxastic accessibility relation.

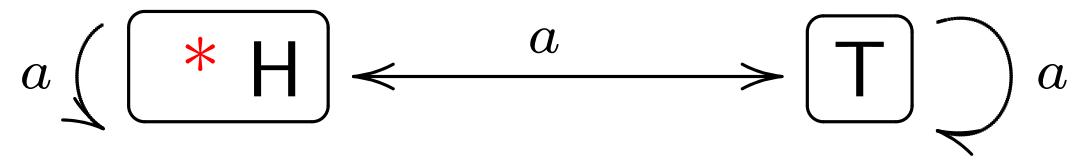
So a **Kripke model for knowledge AND belief** is of the form

$$(S, \{\sim_a\}_{a \in \mathcal{A}}, \{\xrightarrow{a}\}_{a \in \mathcal{A}}, \|\cdot\|, s_*)$$

with K_a interpreted as the modality $[\sim_a]$ for the epistemic relation \sim_a ,
and B_a as the modality $[\xrightarrow{a}]$ for the doxastic relation \xrightarrow{a} .

Coin example again: knowledge

The (single) agent's **knowledge** in the concealed coin scenario can now be represented as:



The arrows represent the **epistemic relation** \sim_a , which captures the agent's **uncertainty** about the state the world. An arrow from state s to state t means that, if s were the real state, then the agent *wouldn't distinguish it* from state t : *for all he knows, the real state might be t .*

Knowledge properties

The fact that K_a in this model satisfied our validities (1)-(3) is now reflected in the fact that \sim_a is an **equivalence relation** in this model:

- The **Veracity** (known as axiom **T** in modal logic) $K_a\varphi \Rightarrow \varphi$ corresponds to the **reflexivity** of the relation \sim_a .
- **Positive Introspection** (known as axiom **4** in modal logic) $K_a\varphi \Rightarrow K_aK_a\varphi$ corresponds to the **transitivity** of the relation \sim_a .
- **Negative Introspection** (known as axiom **5** in modal logic) $\neg K_a\varphi \Rightarrow K_a\neg K_a\varphi$ corresponds to **Euclideaness** of the relation \sim_a :

if $s \sim_a t$ and $s \sim_a w$ then $t \sim_a w$.

In the context of the other two, Euclideaness is equivalent to **symmetry**:

if $s \sim_a t$ then $t \sim_a s$.

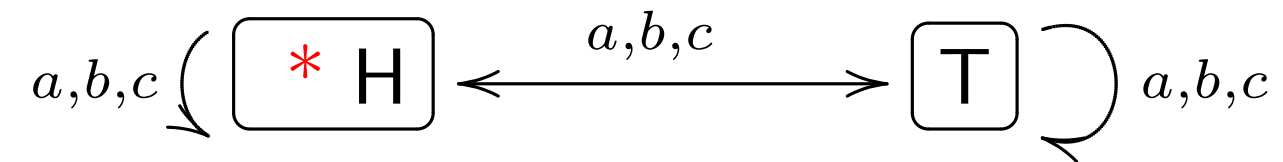
Epistemic Models

An **epistemic model** (or **S5-model**) is a Kripke model in which all the accessibility relations are **equivalence relations**, i.e. **reflexive**, **transitive** and **symmetric**

(or equivalently: **reflexive**, **transitive** and **Euclidean**).

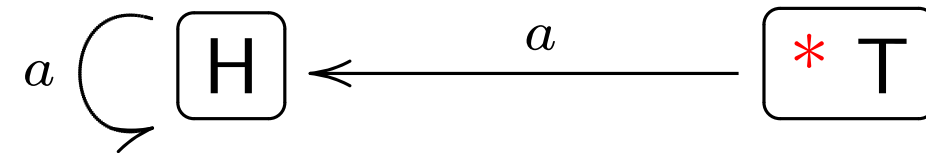
Multi-agent scenario: the concealed coin

Two players a , b and a referee c play a game. In front of everybody, the referee throws a fair coin, catching it in his palm and fully covering it, before anybody (including himself) can see on which side the coin has landed.



Coin example again: beliefs

The (single) agent's beliefs after the mistaken update are now representable as:



In both worlds (i.e. irrespective of what world is the real one), the agent a believes that the coin lies Heads up.

Belief properties

The fact that belief in this model satisfied our validities (4)-(6) is now reflected in the fact that the doxastic accessibility in the above model has the following properties:

- **Consistency** of beliefs (known as axiom **D** in modal logic)
 $\neg B_a(\varphi \wedge \neg\varphi)$ corresponds to the **seriality** of the relation \xrightarrow{a} :

$$\forall s \exists t \text{ such that } s \xrightarrow{a} t.$$

- *Positive Introspection for Beliefs* (axiom **4**) $B_a\varphi \Rightarrow B_a B_a\varphi$
corresponds to the **transitivity** of the relation \xrightarrow{a} .
- *Negative Introspection for Beliefs* (axiom **5**) $\neg B_a\varphi \Rightarrow B_a \neg B_a\varphi$
corresponds to **Euclideaness** of the relation \xrightarrow{a} .

Doxastic Models

A **doxastic model** (or *KD45-model*) is a Φ -Kripke model satisfying the following properties:

- **(D) Seriality**: for every s there exists some t such that $s \xrightarrow{a} t$;
- **(4) Transitivity**: If $s \xrightarrow{a} t$ and $t \xrightarrow{a} w$ then $s \xrightarrow{a} w$
- **(5) Euclideaness** : If $s \xrightarrow{a} t$ and $s \xrightarrow{a} w$ then $t \xrightarrow{a} w$

Putting together in the same structure the belief arrows \xrightarrow{a} from the previous example with the knowledge arrows from before, now denoted by \sim_a ,

we obtain a **Kripke model for both knowledge AND belief.**

Properties connecting Knowledge and Belief

The fact that knowledge and belief in this model satisfied our validities (7)-(9) is now reflected in the fact that the accessibility relations \xrightarrow{a} in the above model have the following properties:

- **Strong Positive Introspection** of beliefs $B_a\varphi \Rightarrow K_a B_a\varphi$ corresponds to

if $s \sim_a t$ and $t \xrightarrow{a} w$ then $s \xrightarrow{a} w$.

- **Strong Negative Introspection** of beliefs $\neg B_a\varphi \Rightarrow K_a \neg B_a\varphi$ corresponds to

if $s \sim_a t$ and $s \xrightarrow{a} w$ then $t \xrightarrow{a} w$.

- **Knowledge Implies Beliefs** $K_a\varphi \Rightarrow B_a\varphi$ corresponds to

if $s \xrightarrow{a} t$ then $s \sim_a t$.

Epistemic-Doxastic Kripke Models

A Kripke model satisfying all the above conditions on the relations \sim_a and \xrightarrow{a} is called an **epistemic-doxastic Kripke model**.

There are two important observations to be made about these models:

- first, *they are completely equivalent to our simple, sphere-based epistemic-doxastic models;*
- second, *the epistemic relation is completely determined by the doxastic relation.*

Equivalence of (single agent) Models

EXERCISE: *For every epistemic-doxastic model $\mathbf{S} = (S, S_0, \|\cdot\|, s_*)$ there exists a doxastic-epistemic Kripke model $\mathbf{S}' = (S, \sim, \longrightarrow, \|\cdot\|, s_*)$ (having the same set of worlds S , same valuation $\|\cdot\|$ and same real world s_*), such that **the same sentences of doxastic-epistemic logic are true at the real world s in model \mathbf{S} as in model \mathbf{S}' :***

$$s \models_{\mathbf{S}} \varphi \text{ iff } s \models_{\mathbf{S}'} \varphi,$$

for every sentence φ .

Conversely, *for every doxastic-epistemic Kripke model $\mathbf{S}' = (S, \sim, \longrightarrow, \|\cdot\|, s_*)$ there exist a doxastic-epistemic model $\mathbf{S} = (S, S_0, \|\cdot\|, s_*)$ such that, for every sentence φ , we have:*

$$s \models_{\mathbf{S}} \varphi \text{ iff } s \models_{\mathbf{S}'} \varphi.$$

Doxastic Relations Uniquely Determine Epistemic Ones

EXERCISE:

Given a doxastic Kripke model $(S, \rightarrow, \|\cdot\|, s_*)$ (i.e. one in which the accessibility relation \rightarrow is serial, transitive and Euclidean), there is a unique relation $\sim \subseteq S \times S$ such that $(S, \sim, \rightarrow, \|\cdot\|, s_*)$ is a doxastic-epistemic Kripke model.

This means that, to encode an epistemic-doxastic model as a Kripke model, we only need to draw the arrows for the doxastic relation.

S4 Models for weak types of knowledge

But, we can see that, in the setting of Kripke models, **the properties specific to “epistemic-doxastic models” are NOT automatically satisfied.**

So Kripke semantics is *more general* than the “sphere semantics”.

In fact, one can use Kripke semantics to interpret various *weaker notions of “knowledge”*, e.g. a type of knowledge that is *truthful* (factive) and *positively introspective*, but *NOT necessarily negative introspective*.

An *S4-model for knowledge* is a Kripke model satisfying only *reflexivity and transitivity* (but not necessarily symmetry or Euclideaness).

Kripke Models for Non-Standard Notions of Belief

Similarly, by *dropping the corresponding semantic conditions*, one can use Kripke models to represent **non-introspective beliefs**, or even **inconsistent beliefs**.

The Problem of Logical Omniscience

However, it is easy to see that any Kripke modality \Box for an accessibility relation still validates **Kripke's axiom**

$$(K) \quad \Box(\varphi \Rightarrow \psi) \Rightarrow (\Box\varphi \Rightarrow \psi),$$

and still satisfies the **Necessitation Rule**:

if φ is valid, then $\Box\varphi$ is valid.

So, if we interpret the modality as “knowledge” or “belief”, then **every logical validity is known/believed**, and similarly **every logical entailment between two propositions is known/believed**.

This means that Kripke semantics can only model “**ideal**” reasoners, who may have *limited access to external truths*, but have *unlimited inference powers*.

1.3. Common Attitudes

- Distributed Knowledge (belief)
- Common Knowledge (belief)
- Everyone Knows (believes)

Scenario: Distributed Knowledge

Agents: Alice, Bob, Charles and Eve

Suppose Alice would like to know with whom did Bob go out for dinner. But Alice only knows he went out with one of his two friends, Charles or Eve (but not both: they can't stand each other).

Suppose that in fact Bob went out with Eve; so Charles obviously know that Bob didn't go out with him.

If Alice and Charles could put their knowledge *together*, they would find out that Bob went out with Eve. So Alice gives a phone call to Charles, they chat and find out.

Before the chat, none of them knew that Bob has gone out with Eve, but this fact was *distributed knowledge* between the two of them: putting their knowledge together was enough to ensure the knowledge of this new fact.

“Distributed” Modalities

The sentence $D\Box\varphi$ is obtained by quantifying over all worlds that are **simultaneously accessible** by **all** arrows (from a given world):

$s \models_{\mathfrak{S}} D\Box\varphi$ iff $t \models_{\mathfrak{S}} \varphi$ for every t such that $s \xrightarrow{a} t$ holds for all $a \in \mathcal{A}$.

In other words, $D\Box$ is the Kripke modality corresponding to the **intersection of all epistemic relations** $\bigcap_{a \in \mathcal{A}} \xrightarrow{a}$

- When the relations \xrightarrow{a} are *reflexive* (corresponding to some form of “*knowledge*”), $D\Box\varphi$ may be interpreted as **distributed knowledge** (in which case we use the notation $Dk\varphi$ instead).
- When the relations \xrightarrow{a} represent *beliefs*, one can also interpret $D\Box$ as **“distributed belief”** $Db\varphi$, but in this case it might actually be *false*.

Distributed Knowledge Within a Group

Distributed knowledge can also be considered in a restricted form:
distributed knowledge within a given (sub)group $G \subseteq \mathcal{A}$.

The definition is the same, except we restrict the intersection of the arrows within the group G :

$s \models_{\mathbf{s}} D \Box_G \varphi$ iff $t \models_{\mathbf{s}} \varphi$ for every t such that $s \xrightarrow{a} t$ holds for all $a \in G$.

Distributed knowledge captures the **implicit (or “virtual”) knowledge of the group G** : what the agents in G *could come to know* if they would *pool together all their private knowledge*.

Distributed Knowledge: Axiomatization

It is easy to see that each of the semantic properties (reflexivity, transitivity, Euclideaness) corresponding to logical postulates usually attributed to “knowledge” (Veracity, Positive Introspection, Negative Introspection) holds for the intersection relation $\bigcap_{a \in G} \xrightarrow{a}$ whenever it holds for each of the arrows \xrightarrow{a} (for each $a \in G$).

Thus, a **complete axiomatization of epistemic logic with distributed knowledge** is given by: your favorite axioms and rules for multi-agent epistemic logic (i.e. a subset of the axioms and rules of multi-agent *S5*); the corresponding axioms and rules for Distributed Knowledge (corresponding to the same subset of *S5*); the axiom

$$\Box_a \varphi \Rightarrow D \Box_G \varphi \quad , \quad \text{for every } a \in G .$$

“Everybody knows...”

Suppose that, in fact, *everybody knows the road rules in France*.

For instance, everybody knows that a red light means “stop” and a green light means “go”. And suppose *everybody respects the rules that (s)he knows*.

Question: Is this enough for you to feel safe, as a driver?

Answer: NO.

Why? Think about it!

Common Knowledge

Suppose the road rules (and the fact they are respected) are *common knowledge*: everybody knows (and respects) the rules, and everybody knows that everybody knows (and respects) the rules, and... etc.

Now, you can drive safely!

Another Example: The Coordinated Attack

Two divisions of the same army, commanded by general A and general B , are camped on two hilltops overlooking a valley. In the valley awaits the enemy (C).

It is clear that **if both divisions attack simultaneously they will win, while if only 1 division attacks it will be defeated.**

So neither general will attack unless he is absolutely sure that the other will attack with him. General A sends a messenger to general B to coordinate a simultaneous attack, by conveying the message “attack at dawn”. But it is possible that the messenger would be captured by the enemy. Fortunately, on this particular night, everything goes smooth.

How long it will take them to coordinate an attack?

Well, B cannot attack at dawn, after receiving the message, since he's still not sure that A knows he received the message; indeed, A might think it possible the messenger was captured, in which case A will not attack at dawn, since he'll fear B won't attack. So B has to send another messenger to A to confirm the receipt of the first message (an 'acknowledgment'). After receiving it, A knows that B got the first message. But he still cannot attack, since he's not sure B will: for all that B knows, *his messenger* might have been captured (in which case A wouldn't know the first message was received). So A has to send back to B another messenger, confirming the receipt of the previous acknowledgment.

This goes forever, without achieving any coordination: even if no messenger is captured, one can show that no finite number of successful deliveries of "acknowledgments to acknowledgments" can allow the generals to attack!

“Common” Modalities

The sentence $C\Box\varphi$ is obtained by quantifying over all worlds that are accessible by any concatenations of arrows:

$s \models_{\mathbf{S}} C\Box\varphi$ iff $t \models_{\mathbf{S}} \varphi$ for every t and every a finite chain
(of length $n \geq 0$) of the form $s = s_0 \xrightarrow{a_1} s_1 \xrightarrow{a_2} s_2 \cdots \xrightarrow{a_n} s_n = t$.

$C\Box\varphi$ may be interpreted as **common knowledge** (in which case we use the notation $Ck\varphi$ instead) or **common belief** (in which case we use $Cb\varphi$ instead), depending on the context.

Common Knowledge Within a Group

Common knowledge (or belief) can also be considered in a restricted form, as **common knowledge within a given (sub)group** $G \subseteq \mathcal{A}$. Here we restrict the concatenated arrows to arrows within the group G :

$s \models_{\mathbf{S}} C \square_G \varphi$ iff $t \models_{\mathbf{S}} \varphi$ for every t and every a finite chain of the form $s = s_0 \xrightarrow{a_1} s_1 \xrightarrow{a_2} s_2 \cdots \xrightarrow{a_n} s_n = t$, with $a_1, \dots, a_n \in G$.

Full common knowledge/belief $C \square$ (as previously defined) corresponds to the case that G is the set \mathcal{A} of **all** agents:

$$C \square \varphi = C \square_{\mathcal{A}} \varphi$$

Common Knowledge as an Infinite Conjunction

If we make the abbreviation

$$E_G\varphi := \bigwedge_{a \in G} \Box_a \varphi$$

(“*everybody knows φ* ”), then we can easily check that:

$s \models_{\mathbf{S}} C\Box_G\varphi$ iff s satisfies **all** the (infinitely many) sentences

$$\varphi, E_G\varphi, E_GE_G\varphi, E_GE_GE_G\varphi, \dots$$

In *this sense*, we can say that $E_G\varphi$ is equivalent to the “infinite conjunction”

$$\varphi \wedge E_G\varphi \wedge E_GE_G\varphi \wedge \dots$$

However, the most used modal-epistemic languages are *finitary*, so that $C\Box$ cannot be defined as the (impossible to form) infinite conjunction.

Instead, $C\Box$ is usually taken as a **primitive** operator (introduced via the above semantic clause).

2.1. Logics of public and private announcements

PAL (the logic of public announcements) was first formalized (including Reduction Laws) by Plaza (1989) and independently by Gerbrandy and Groeneveld (1997).

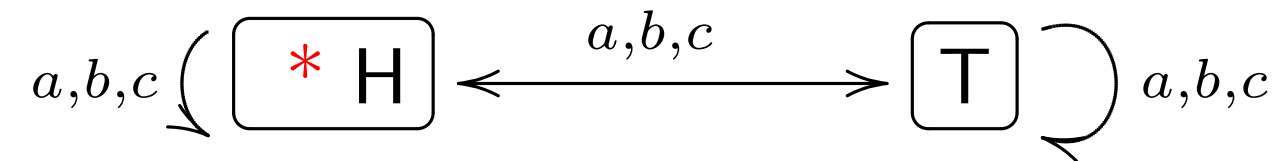
The problem of **completely axiomatizing PAL in the presence of the common knowledge operator** was first solved by Baltag, Moss and Solecki (1998).

A logic for “**secret (fully private) announcements**” was first proposed by Gerbrandy (1999).

A logic for “**private, but legal, announcements**” (what we will call “*fair-game announcements*”) was developed by H. van Ditmarsch (2000).

Recall multi-agent scenario: the concealed coin

Two players a , b and a referee c play a game. In front of everybody, the referee throws a fair coin, catching it in his palm and fully covering it, before anybody (including himself) can see on which side the coin has landed.

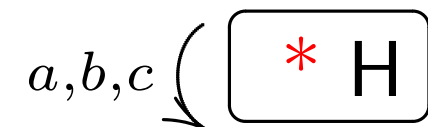


Scenario 2: The coin revealed

The referee c opens his palm and shows the face of the coin to everybody (to the public, composed of a and b , but also to himself): they **all see** it's Heads up, and **they all see that the others see it** etc.

So this is a **“public announcement” that the coin lies Heads up.**

We denote this event by $!H$. Intuitively, after the announcement, we have common knowledge of H , so the model of the new situation is:



Public Announcements are (Joint) Updates!

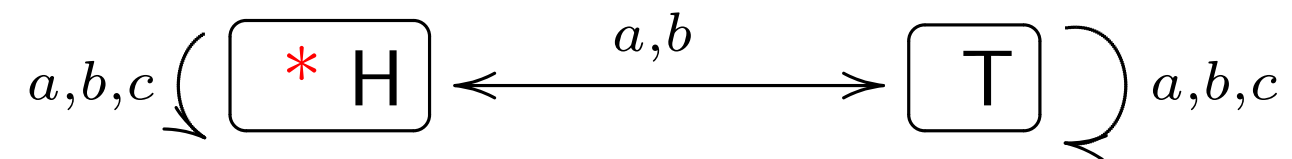
But this is just the result of **updating with H**: deleting all the non-H-worlds.

So, in the multi-agent case, **updating captures public announcements**.

From now on, we denote by $!\varphi$ the operation of deleting the non- φ worlds, and call it **public announcement with φ** , or **joint update with φ** .

Scenario 3: ‘Legal’ Private Viewing

Instead of Scenario 2: in front of everybody, the referee (c) uncovers the coin, so that (they all see that) **he, and only he, can see the upper face**. This changes the initial model to



Now, c **knows** the real state. E.g. if it's Heads, he knows it, and disregards the possibility of Tails. a and b don't know the real state, but *they know that c knows it*.

c 's viewing of the coin is a “legal”, non-deceitful, although **private action**.

Fair-Game Announcements

Equivalently: in front of everybody, an announcement of the upper face of the coin is made, but in such a way that (it is common knowledge that) only c hears it.

Such announcements (first modeled by H. van Ditmarsch) are called **fair-game announcements**:

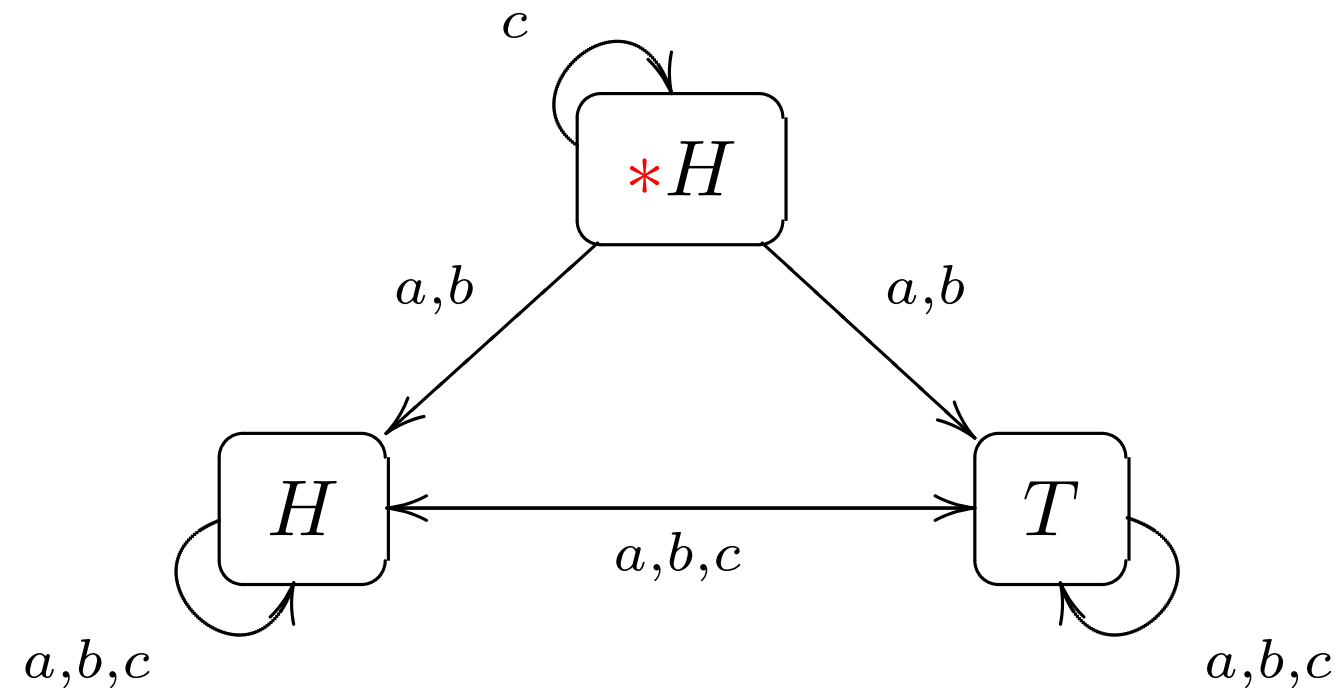
they can be thought of as “legal moves” in a fair game: nobody is cheating, all players are aware of the possibility of this move, but only some of the players (usually the one who makes the move) can see the actual move. The others know the range of possible moves at that moment, and they know that the “insider” knows his move, but they don’t necessarily know the move.

Scenario 4: Cheating

Suppose that, after Scenario 1, the referee c has **taken a peek at the coin**, before covering it. **Nobody has noticed this**. Indeed, let's assume that c **knows that a and b did not suspect anything**.

This is an instance of **cheating**: a private viewing which is “illegal”, in the sense that it is deceitful for a and b . Now, a and b think that nobody knows on which side the coin is lying. But they are wrong!

The Model after Cheating



We indicated the *real world* here. In the actual world (above), a and b think that the only possibilities are the worlds below. That is, they *do not even consider the “real” world as a possibility.*

Such models in which we indicate the *real world* are called **pointed models**.

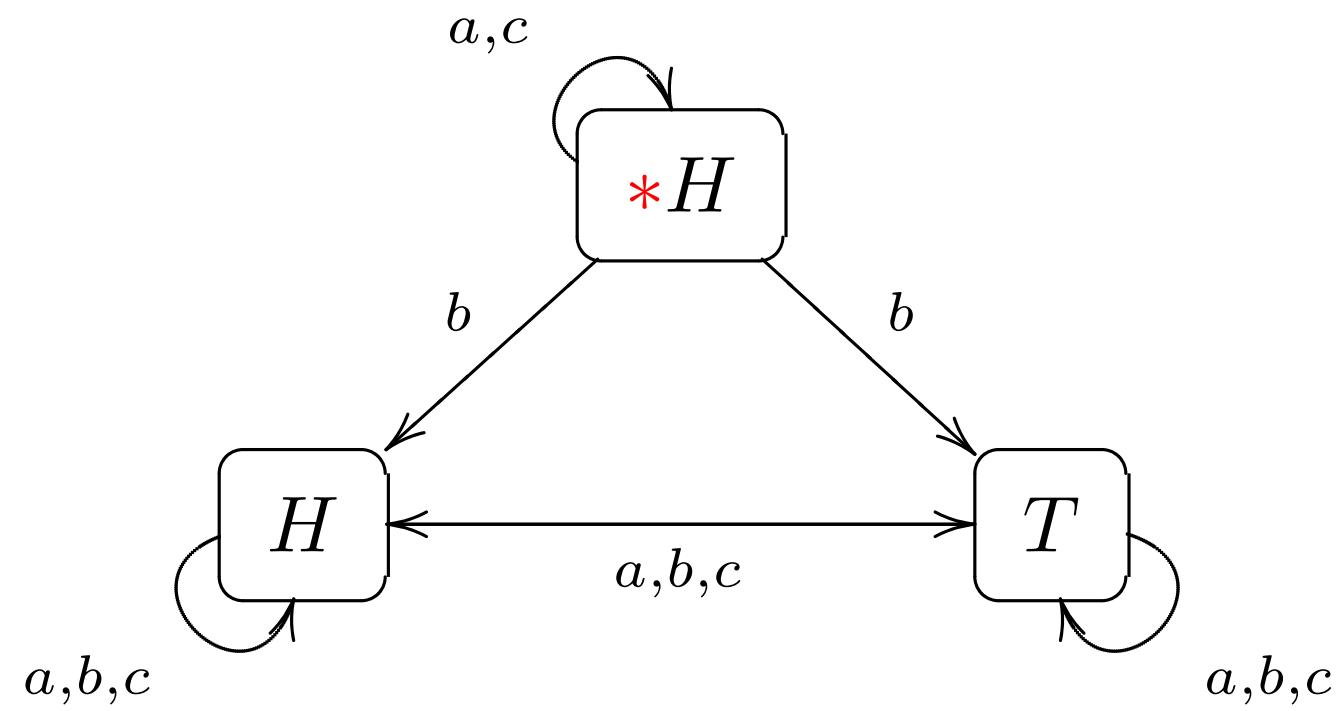
Scenario 5: Secret Communication

After cheating (Scenario 4), c engages in another "illegal" action: **he secretly sends an email to his friend a , informing her that the coin is Heads up.**

Suppose the delivery and the secrecy of the message are guaranteed: so a and c have common knowledge that H, and that b doesn't know they know this.

Indeed, b is completely fooled: he doesn't suspect that c could have taken a peek, nor that he could have been engaged in secret communication.

The model is



Private Announcements

Both of the above actions were examples of completely **private announcements**

$!_G\varphi$

of a sentence φ to a group G of agents: in the first case $G = \{c\}$, in the second case $G = \{a, c\}$.

The “insiders” (in G) know what’s going on, the “outsiders” don’t suspect anything.

Scenario 5': Wiretapping?

In Scenario 5', everything goes on as in Scenario 5, except that in the meantime b is **secretely breaking** into c 's email account (or **wiretapping** his phone) and reading c 's secret message.

Nobody suspects this illegal attack on c 's privacy. So both c and a think their secret communication is really secret and unsuspected by b : **the deceivers are deceived.**

What is the model of the situation after this action?!

Things are getting rather complicated!

Scenario 6

This starts right after Scenario 2, when it was common knowledge that c knew the face. c attempts to send a secret message to a announcing that H is the case.

- c is convinced the communication channel is fully secure and reliable; moreover, he thinks that b doesn't even suspect this secret communication is going on.
- In fact, unknown and unsuspected by c , the message is *intercepted, stopped and read* by b . As a result, *it never makes it to a* , and in fact a never knows or suspects any of this.
- As for b , he *knows* all of the above: not only now he knows the message, but he knows that he “fooled” everybody, in the way described above.

The Update Problem

We need to find a *general method* to solve all the above problems, i.e. to compute all these different kinds of updates.

2.2. “Standard DEL”

- studies the **multi-agent information flow of “hard information”** (irrevocable, absolutely certain, fully introspective “knowledge”) as well as “soft”, but essentially un-revisable, information (“beliefs” that change monotonically, but are never overturned);
- gives an answer to the Update Problem, based on the BMS (Baltag, Moss and Solecki 98) setting: **logics of epistemic actions**;
- it arose from generalizing previous work on logics for public/private announcements.
- this dynamics is **essentially monotonic** (no belief revision!), though *it can model very complex forms of communication*.

Models for 'Events'

- Until now, our Kripke models capture only *epistemic situations*, i.e. they only contain *static* information: they all are *state models*. We can thus represent the *result* of each of our Scenarios, but not what is actually going on.
- Our scenarios involve various *types of changes* that may affect agents' beliefs or state of knowledge: a public announcement, a 'legal' (non-deceitful) act of private learning, 'illegal' (unsuspected) private learning etc.
- We want to use now Kripke models to represent such types of *epistemic events*, in a way that is similar to the representations we have for epistemic states.

Event Models

An **event model** (or “*action model*”)

$$\Sigma = (\Sigma, \xrightarrow{A}, pre)$$

is just like an Kripke model,

except that its elements are now called **actions** (or “*simple events*”)
and instead of the valuation we have a **precondition map** pre ,
associating a sentence pre_σ to each action σ .

Epistemic/Doxastic Event Models

An event model is **epistemic**, or respectively a **doxastic**, event model if it satisfies the S5, or respectively the KD45, conditions.

Interpretation

We call the simple events $\sigma \in \Sigma$: *deterministic actions* of a particularly simple kind, i.e. they do not change the “facts” of the world, but only the agents’ beliefs/knowledge. In other words, they are “**purely epistemic**” actions.

For $\sigma \in \Sigma$, we interpret pre_σ as giving the **precondition** of the action σ : this is a sentence that is true in a world iff σ can be performed. In a sense, pre_σ gives the implicit information carried by σ .

Finally, the *accessibility relations* express the agents’ **knowledge/beliefs about the current action taking place.**

The Product Update

Given a state model $\mathbf{S} = (S, \xrightarrow{A}, \|\cdot\|)$ and an action model $\Sigma = (\Sigma, \xrightarrow{A}, pre)$, we define their *update product*

$$\mathbf{S} \otimes \Sigma = (S \otimes \Sigma, \xrightarrow{A}, \|\cdot\|)$$

to be a new state model, given by:

1. $S \otimes \Sigma$ is

$$\{(s, \sigma) \in S \times \Sigma : s \models_{\mathbf{S}} pre_{\sigma}\}.$$

2. $(s, \sigma) \xrightarrow{A} (s', \sigma')$ iff $s \xrightarrow{A} s'$ and $\sigma \xrightarrow{A} \sigma'$.

3. $\|p\|_{\mathbf{S} \otimes \Sigma} = \{(s, \sigma) \in S \otimes \Sigma : s \in \|p\|_{\mathbf{S}}\}.$

Product of Pointed Models

As before, we can consider **pointed event models**, if we want to specify the **actual event** taking place.

Naturally, if initially the actual state was s and then the actual event is σ , then the actual output-state is (s, σ) .

Interpretation

The **product arrows encode the idea that**: two output-states are indistinguishable iff they are the result of indistinguishable actions performed on indistinguishable input-states.

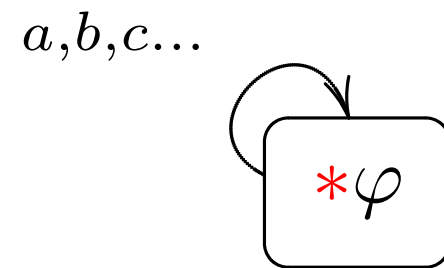
This comprises two intuitions:

1. “*No Miracles*”: knowledge can only be gained from (the epistemic appearance of) actions;
2. “*Perfect Recall*”: once gained, knowledge is never lost.

The fact that the valuation is the same as on the input-state tells us that these actions are **purely epistemic**.

Examples: Public Announcement

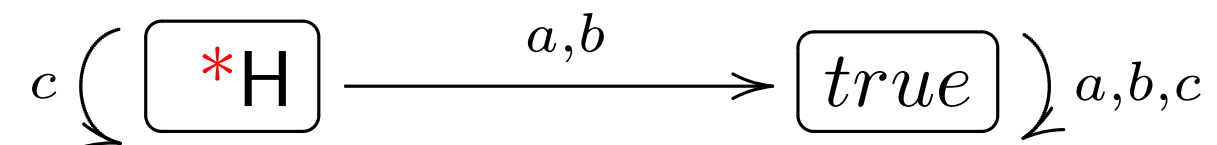
The event model $\Sigma_{!\varphi}$ for public announcement $!\varphi$ consists of a single action, with precondition φ and reflexive arrows:



EXERCISE: Check that, for every state model \mathbf{S} , $\mathbf{S} \otimes \Sigma_{!\varphi}$ is indeed the result of deleting all non- φ worlds from \mathbf{S} .

More Examples: Taking a Peek

The action in Scenario 4: c takes a peek at the coin and sees the Head is up, without anybody noticing.

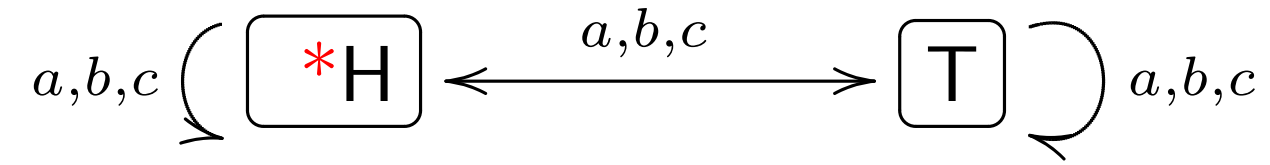


There are two actions in this model:

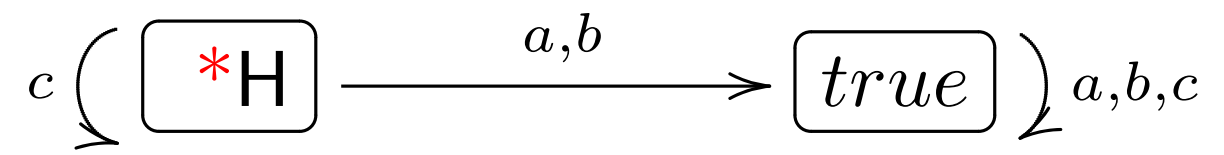
- 1) the real event (on the left) is the **cheating action** of c “taking a peek”.
- 2) The action on the right is the apparent action *skip*, having any tautological sentence *true* as its precondition: this is the action in which **nothing happens**. *This is what the outsiders (a and b) think it is going on: nothing, really.*

The Product Update

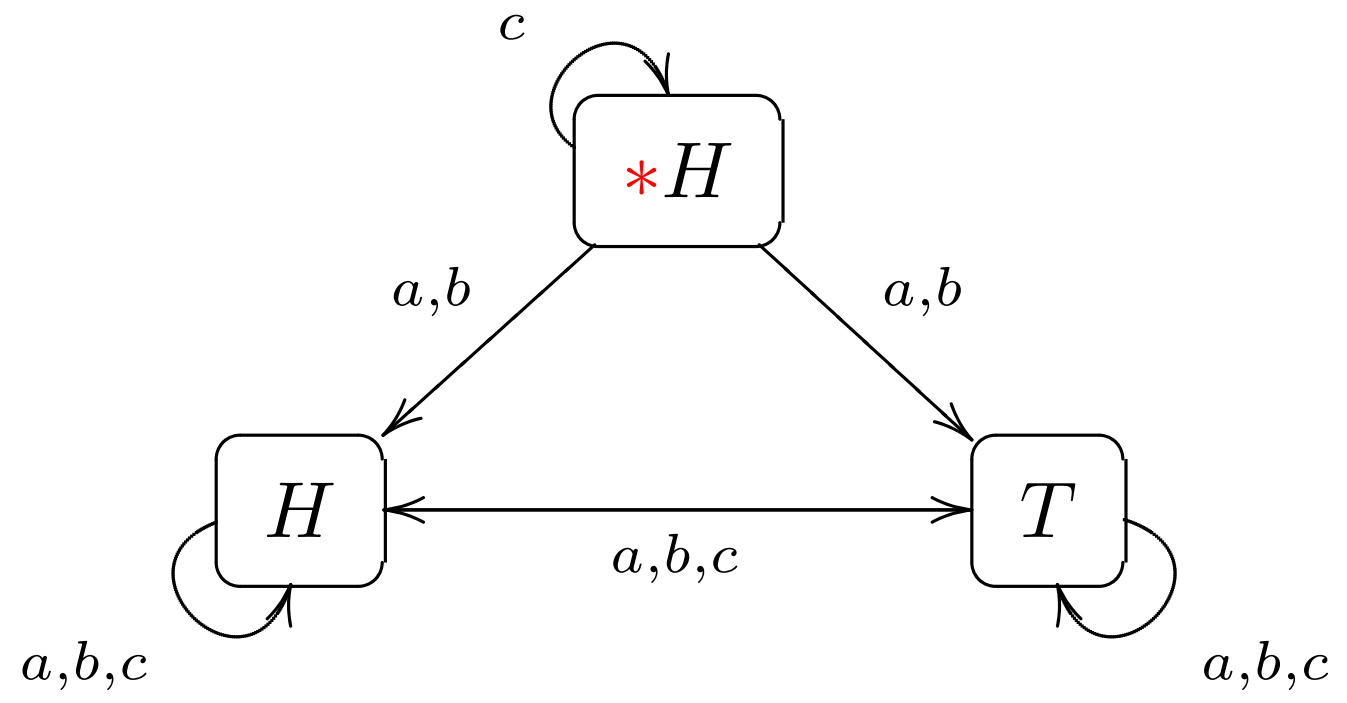
We can now check that the product of



and

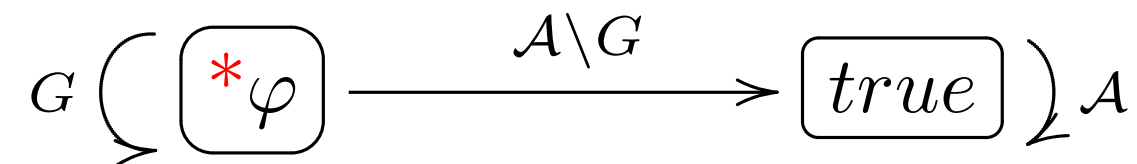


is indeed what intuitively should be:



Private Announcements

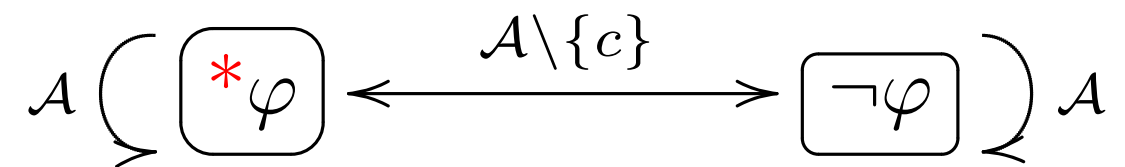
More generally, a fully **private announcement** $!_G\varphi$ of φ to a subgroup G is described by the action on the left in the event model



This subsumes both taking a peak (Example 4) and the secret communication in Example 5.

Fair-Game Announcements

The following event model represents the situation in which *it is common knowledge that an agent c privately learns whether φ or $\neg\varphi$ is the case*:

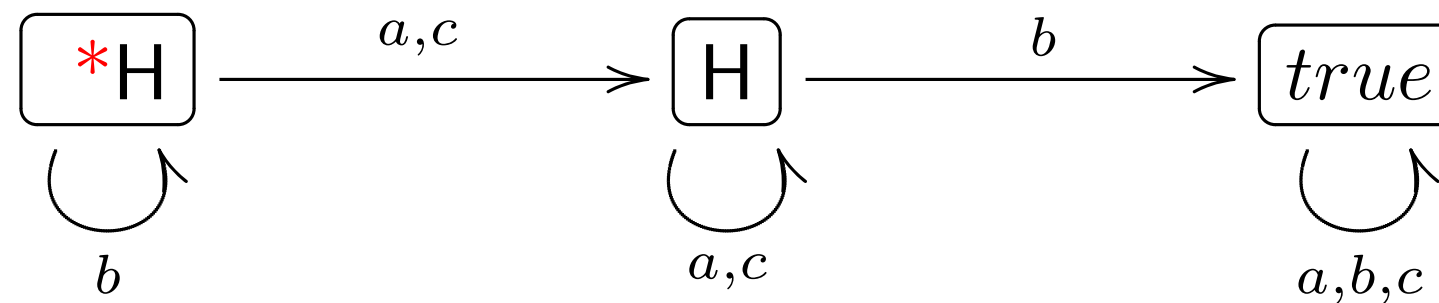


This is a **“fair-game announcement”** $Fair_c \varphi$.

The case $\varphi := H$ represents the action in Example 3 (“legal viewing” of the coin by c).

Solving Scenario 5': Wiretapping

Recall Scenario 5: the supposedly secret message from c to a is secretly intercepted by b . This is an instance of a *private announcements with (secret) interception by a group of outsiders*.



Dynamic Modalities

For any action $\sigma \in \Sigma$, we can consider the corresponding **dynamic modality** $[\sigma]\varphi$.

This is a property *of the original model*, expressing the fact that, if action σ happens, then φ will come to be true after that.

We can easily define the epistemic proposition $[\sigma]\varphi$ by:

$$s \models_{\mathbf{S}} [\sigma]\varphi \text{ iff } (s, \sigma) \in \mathbf{S} \otimes \Sigma \text{ implies } (s, \sigma) \models_{\mathbf{S} \otimes \Sigma} \varphi$$

Appearance

For any agent a and any action $\sigma \in \Sigma$, we define the **appearance of action σ to a** , denoted by σ_a , as:

$$\sigma_a = \{\sigma' \in \Sigma : \sigma \xrightarrow{a} \sigma'\}$$

When σ happens, it appears to a as if either one of the actions $\sigma' \in \sigma_a$ is happening.

Examples

$$(!\varphi)_a = \{!\varphi\} \text{ for all } a \in \mathcal{A},$$

$$(!G\varphi)_a = \{!G\varphi\} \text{ for all insiders } a \in G,$$

$$(!G\varphi)_a = \{skip\} = \{!(true)\} \text{ for all outsiders } a \notin G,$$

$$(Fair_a\varphi)_a = \{Fair_a\varphi\}$$

$$(Fair_a\varphi)_b = \{Fair_a\varphi, Fair_a\neg\varphi\} \text{ for } b \neq a.$$

Reduction Laws

If $\sigma \in \Sigma$ is a simple epistemic action, then we have the following properties (or “axioms”):

- *Preservation of “Facts”*. For all atomic $p \in \Phi$:

$$[\sigma]p \iff pre_\sigma \Rightarrow p$$

- *Partial Functionality*:

$$[\sigma]\neg\varphi \iff pre_\sigma \Rightarrow \neg[\sigma]\varphi$$

- *Normality*:

$$[\sigma](\varphi \wedge \psi) \iff [\sigma]\varphi \wedge [\sigma]\psi$$

Here, \square can be *either knowledge K or belief B* , depending on whether the model is doxastic or epistemic.

- “Action-Knowledge Axiom”:

$$[\sigma]\Box_a\varphi \iff pre_\sigma \Rightarrow \bigwedge_{\sigma' \in \sigma_a} \Box_a[\sigma']\varphi$$

This Action-Knowledge Axiom helps us to *compute the state of knowledge/belief* of an agent *after* an event, in terms of the agent’s *initial state of knowledge or belief* and of the event’s *appearance* to the agent.

Instances of Action-Knowledge Axiom

If $a \in G$, $b \notin G$, $c \neq a$, then:

$$[!\theta]B_a\varphi \iff \theta \Rightarrow B_a[!\theta]\varphi$$

$$[!_G\theta]B_a\varphi \iff \theta \Rightarrow B_a[!_G\theta]\varphi$$

$$[!_G\theta]B_b\varphi \iff \theta \Rightarrow B_b\varphi$$

$$[Fair_a\theta]B_a\varphi \iff \theta \Rightarrow B_a[Fair_a\theta]\varphi$$

$$[Fair_a\theta]B_c\varphi \iff \theta \Rightarrow B_c([\![Fair_a\theta]\!]\varphi \wedge [\![Fair_a\neg\theta]\!]\varphi)$$

EXERCISES

- Solve Scenario 5', by computing the update product of the state model obtained in Scenario 4 with the event model that we saw.
- Solve Scenario 6 using update product.
- Solve the Muddy Children puzzle, using repeated updates. Encode the conclusion of the puzzle in a DEL sentence.