

# TOWARDS EPISTEMIC-DOXASTIC PLANNING WITH OBSERVATION AND REVISION

---

Thorsten Engesser<sup>1</sup>, Andreas Herzig<sup>1</sup>, Elise Perrotin<sup>2</sup>

KR and MAS Conventicle 2024

<sup>1</sup>IRIT, CNRS, Toulouse, France

<sup>2</sup>CRIL, CNRS, Lens, France



Planning with a **theory of mind** is a valuable skill for autonomous agents:

- Accounting for other agents with *false beliefs*.
- Planning to facilitate *coordination*.

Most existing planning formalisms support *knowledge* or *belief*, but not both.

Our approach is inspired by *lightweight* epistemic & doxastic planning approaches from the literature [Cooper et al., 2021, Muise et al., 2022].

We consider the epistemic-doxastic logic S5-EDL:

$$\varphi ::= K_i\varphi \mid B_i\varphi \mid \neg\varphi \mid \varphi \wedge \varphi$$

- Knowledge: Facts which we currently observe.
- Belief: Things we observed in the past or learn through communication.
- S5 for knowledge, KD45 for belief + interactions axioms (e.g.,  $K_i\varphi \rightarrow B_i\varphi$ ).
- The satisfiability problem of the full logic is PSPACE-complete.

We do not want to use Kripke models + DEL update models for states/actions.

⇒ Can we find something simpler?

## KNOWLEDGE ONLY: THE EPISTEMIC LOGIC OF OBSERVATIONS (EL-O)

States = Valuations over **observation atoms** [Cooper et al., 2021]:

$$\sigma ::= p \mid S_i \sigma$$

- $S_i \sigma$ : *agent i sees  $\sigma$*  (= knowing whether,  $S_i \sigma \equiv K_i \sigma \vee K_i \neg \sigma$ ).
- No negations, conjunctions or disjunctions within modal operator.
- Introspective atoms such as  $S_1 S_1 p$  are excluded (they are tautological).

Example:

$$\{p, S_1 p, S_1 S_2 p\} \models K_1 p \wedge \neg K_2 p \wedge K_1 \neg K_2 p$$

## KNOWLEDGE ONLY: THE EPISTEMIC LOGIC OF OBSERVATIONS (EL-O)

- Introspection-free observation atoms logically independent of each other.  
 $\Rightarrow$  S5-satisfiability of formulas over such atoms reduces to boolean SAT.

Unfortunately, this approach does not work with *having a belief about*:

$$B_i\varphi \not\equiv BA_i\varphi \wedge \varphi$$

4 epistemic situations:

	$\neg\sigma$	$\sigma$
$\neg S_i\sigma$	$\neg\sigma \wedge \neg K_i\neg\sigma$	$\sigma \wedge \neg K_i\sigma$
$S_i\sigma$	$\neg\sigma \wedge K_i\neg\sigma$	$\sigma \wedge K_i\sigma$

6 doxastic situations:

	$\neg\beta$	$\beta$
?	$\neg\beta \wedge \neg B_i\neg\beta \wedge \neg B_i\beta$	$\beta \wedge \neg B_i\neg\beta \wedge \neg B_i\beta$
?	$\neg\beta \wedge \neg B_i\neg\beta \wedge B_i\beta$	$\beta \wedge \neg B_i\neg\beta \wedge B_i\beta$
?	$\neg\beta \wedge B_i\neg\beta \wedge \neg B_i\beta$	$\beta \wedge B_i\neg\beta \wedge \neg B_i\beta$

# TRUE BELIEFS AND MERE BELIEFS [HERZIG AND PERROTIN, 2021]

True belief about  $\varphi$ :  $TBA_i\varphi \equiv (B_i\varphi \wedge \varphi) \vee (B_i\neg\varphi \wedge \neg\varphi)$

Mere belief about  $\varphi$ :  $MBA_i\varphi \equiv (B_i\varphi \wedge \neg K_i\varphi) \vee (B_i\neg\varphi \wedge \neg K_i\neg\varphi)$

All combinations of knowledge and belief are expressible:

e.g., assuming  $\varphi$  is true:

$i$ has <b>no belief about</b> $\varphi$	$\neg MBA_i\varphi \wedge \neg TBA_i\varphi$	$\varphi \wedge \neg B_i\varphi \wedge \neg B_i\neg\varphi$
$i$ has a <b>false belief about</b> $\varphi$	$MBA_i\varphi \wedge \neg TBA_i\varphi$	$\varphi \wedge B_i\neg\varphi \wedge \neg K_i\varphi$
$i$ has a <b>lucky belief about</b> $\varphi$	$MBA_i\varphi \wedge TBA_i\varphi$	$\varphi \wedge B_i\varphi \wedge \neg K_i\varphi$
$i$ <b>knows whether / observes</b> $\varphi$	$\neg MBA_i\varphi \wedge TBA_i\varphi$	$\varphi \wedge B_i\varphi \wedge K_i\varphi$

## A LIGHTWEIGHT FRAGMENT OF S5-EDL

We consider boolean formulas over so-called *REDA atoms*:

$$\alpha ::= p \mid \text{TBA}_i\alpha \mid \text{MBA}_i\alpha$$

- *REDA*: *repetition-free* epistemic-doxastic atoms.
  - ⇒ No negations, conjunctions or disjunctions within modal operator.
  - ⇒ Introspective atoms such as  $\text{TBA}_i\text{MBA}_i\alpha$  are excluded.
- Arbitrary conjunctions of such atoms are satisfiable.
  - ⇒ Satisfiability reduces to propositional SAT (NP-complete).

We use valuations over *REDA* atoms as states. For example:

$$\{p, \text{TBA}_i p, \neg \text{MBA}_i p, \text{TBA}_j p, \text{MBA}_j p\} \models K_i p \wedge B_j p$$

Actions have **indirect effects** conditional on agents' observations.

E.g., action of changing the truth value of  $p$ :

- Direct effect:  $\top \triangleright \pm p$
- Indirect effect:  $MBA_{ip} \triangleright \pm TBA_{ip}$ .
- Lucky beliefs become false beliefs and vice versa.
- There are additional higher-order indirect effects...



In our paper, we define the following types of actions:

- Ontic actions (changing the value of a proposition).
- Starting and stopping to observe (first and second-order).

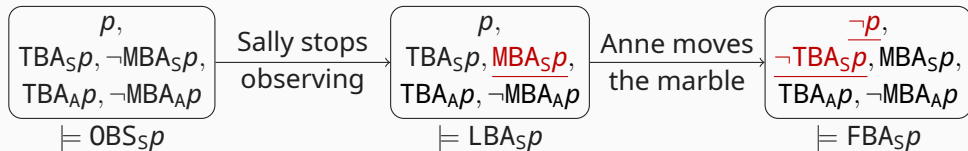
Allows us to model some first- and second-order *false-belief tasks*.

## EXAMPLE: SALLY-ANNE TASK

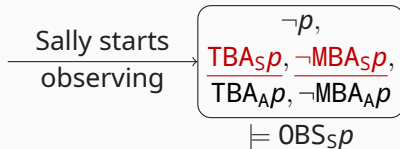
1. Two children, Sally and Anne, are in a room together.
2. Sally has a marble, which she puts into a basket.
3. Sally leaves the room to go out for a walk.
4. Anne removes the marble from the basket and puts it into a box.
5. Sally comes back into the room.

*Will Sally search for her marble in the basket or in the box?*

## EXAMPLE: SALLY-ANNE TASK



If Sally starts to observe the marble again:



We get *revision* for free!

- Satisfiability in our S5-EDL fragment reduces to propositional satisfiability.
- We define an epistemic-doxastic planning formalism.
- Planning reduces to classical planning (PSPACE-complete).

## LIMITATIONS

We define **only actions for second-order** knowledge and beliefs.

- Could be generalized to higher-order.
- Actions with second-order indirect effects are already quite complicated.


Our approach only approximates second-order observability:

$$\begin{aligned} \text{OBS}_i \text{OBS}_j p &\equiv \text{TBA}_i(\text{TBA}_j p \wedge \neg \text{MBA}_j p) \wedge \neg \text{MBA}_i(\text{TBA}_j p \wedge \neg \text{MBA}_j p) \\ &\approx \text{TBA}_i \text{TBA}_j p \wedge \text{TBA}_i \text{MBA}_j p \wedge \neg \text{MBA}_i \text{TBA}_j p \wedge \neg \text{MBA}_i \text{MBA}_j p \\ &\equiv \text{OBS}_i \text{TBA}_j p \wedge \text{OBS}_i \text{MBA}_j p \end{aligned}$$

E.g., we cannot express:


*“I know you don’t observe  $p$ , but I have no idea what you believe.”*

**THANK YOU!**

 Cooper, M. C., Herzig, A., Maffre, F., Maris, F., Perrotin, E., and Régnier, P. (2021).

**A lightweight epistemic logic and its application to planning.**

*Artif. Intell.*, 298:103437.

 Herzig, A. and Perrotin, E. (2021).

**True belief and mere belief about a proposition and the classification of epistemic-doxastic situations.**

*Filosofiska Notiser*, 8(1):103–117.

 Muise, C., Belle, V., Felli, P., McIlraith, S. A., Miller, T., Pearce, A. R., and Sonenberg, L. (2022).

**Efficient multi-agent epistemic planning: Teaching planners about nested belief.**

*Artif. Intell.*, 302:103605.