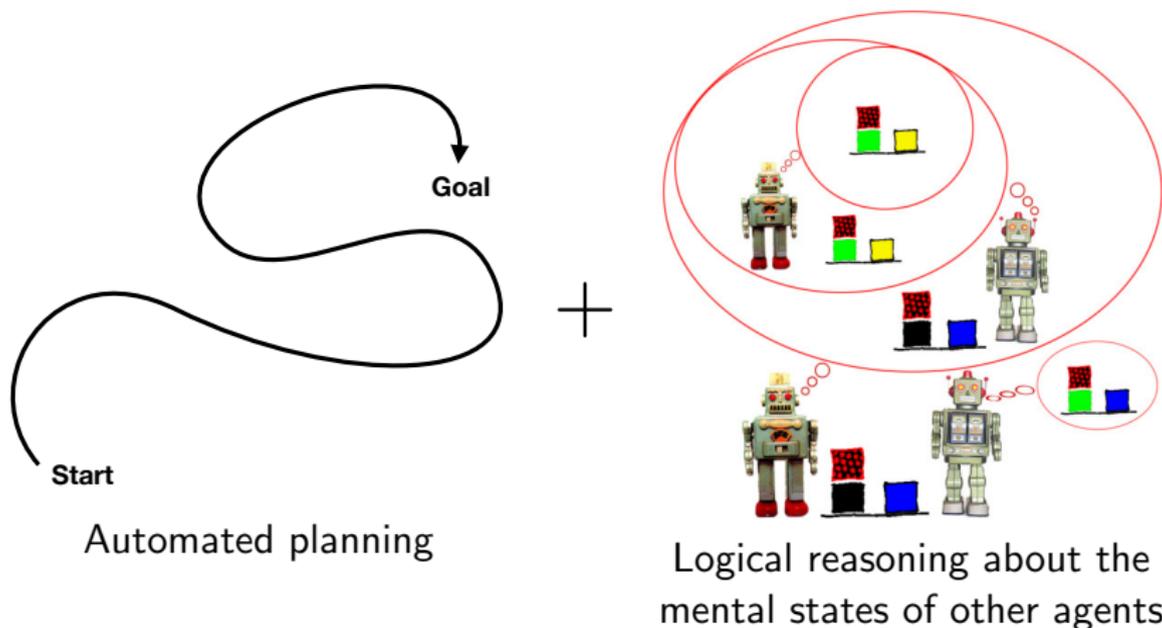
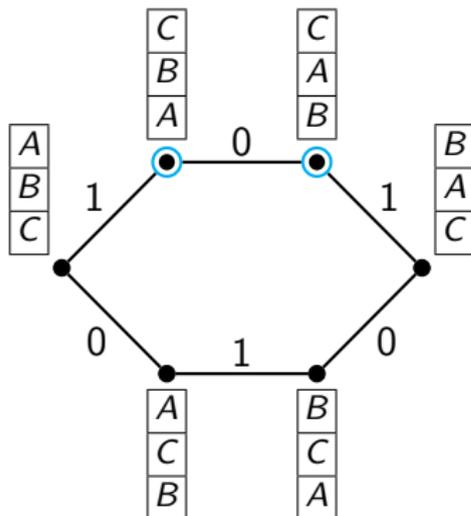
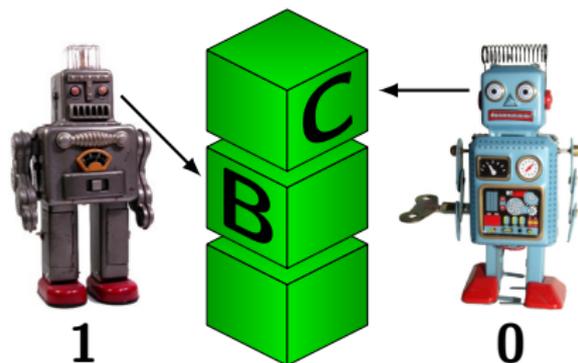
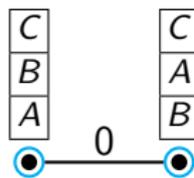
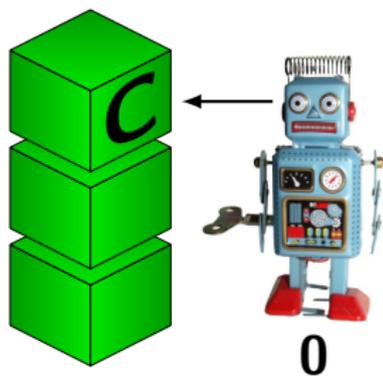




**Epistemic planning** =  
automated *planning* + (dynamic) epistemic *logic*

**Goal:** To compute plans that can take the mental states of other agents into account.





**Epistemic states:** Multi-pointed, finite epistemic models of multi-agent S5.

**Designated states:** ● (those considered possible by planning agent).

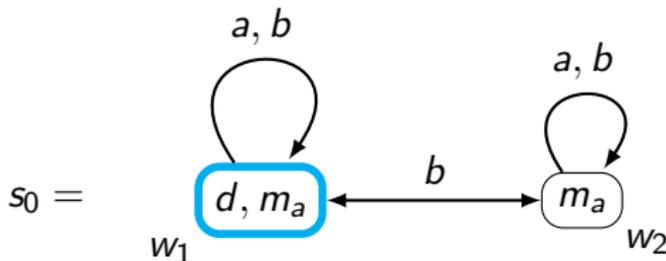
# The coordinated attack problem in dynamic epistemic logic (DEL)

Two generals (agents),  $a$  and  $b$ . They want to coordinate an attack, and only win if they attack simultaneously.

$d$ : “general  $a$  will attack at dawn”.

$m_i$ : the messenger is at general  $i$  (for  $i = a, b$ ).

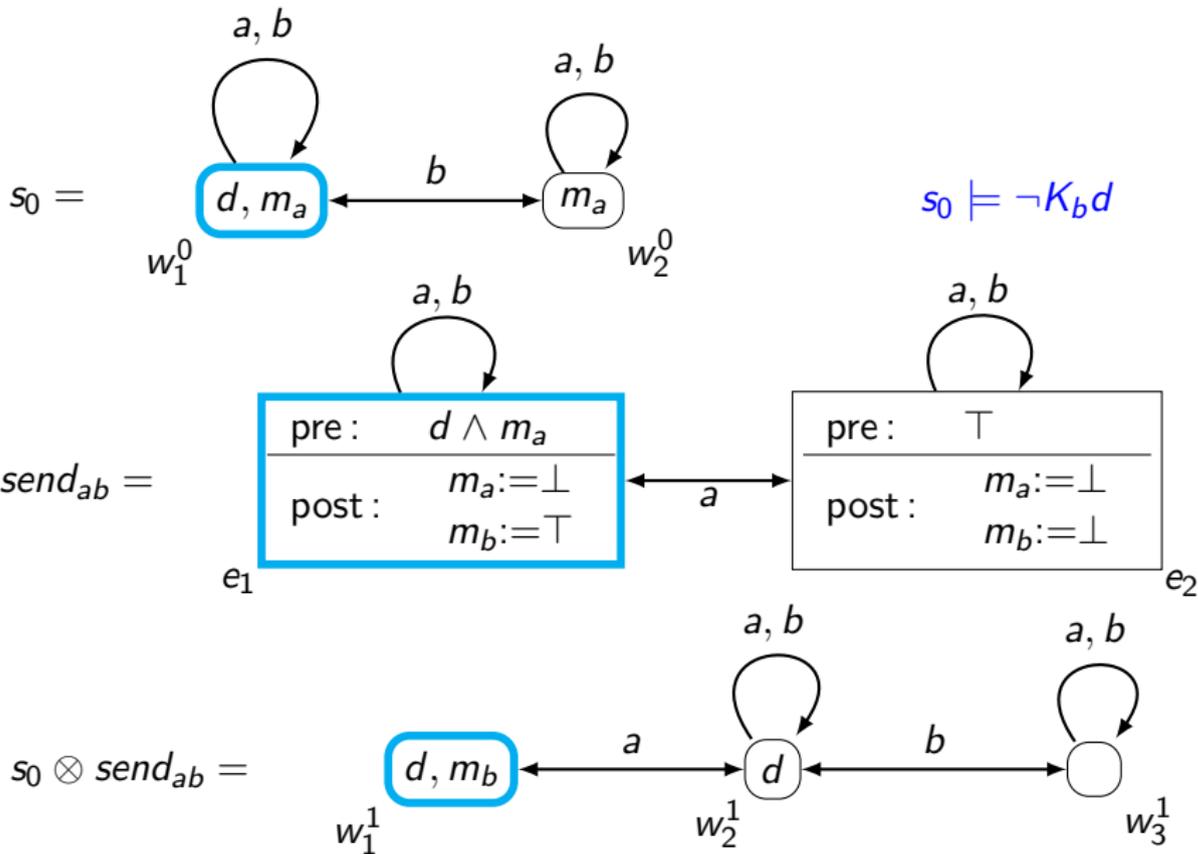
Initial **epistemic state**:



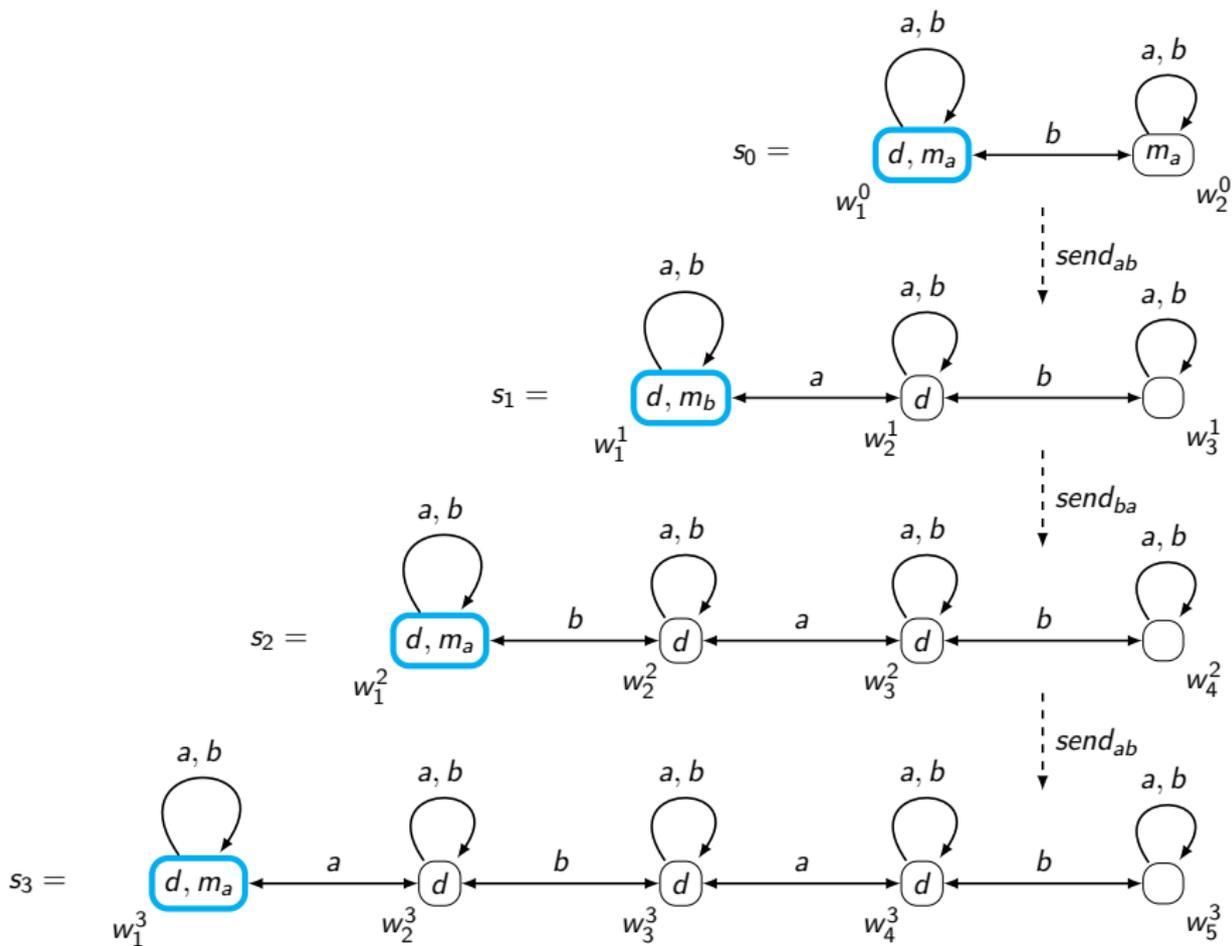
Nodes are **worlds**, edges are **indistinguishability edges** (as long as we're on S5).



# The product update in dynamic epistemic logic



$$s_0 \otimes send_{ab} \models K_a d \wedge K_b d \wedge \neg K_a K_b d$$



## Epistemic planning tasks

**Definition.** An **epistemic planning task** (or simply a **planning task**)  $T = (s_0, A, \varphi_g)$  consists of an epistemic state  $s_0$  called the **initial state**; a finite set of epistemic actions  $A$ ; and a **goal formula**  $\varphi_g$  of the epistemic language.

**Definition.** A **solution** to a planning task  $T = (s_0, A, \varphi_g)$  is a sequence of actions  $\alpha_1, \alpha_2, \dots, \alpha_n$  from  $A$  such that for all  $1 \leq i \leq n$ ,  $\alpha_i$  is applicable in  $s_0 \otimes \alpha_1 \otimes \dots \otimes \alpha_{i-1}$  and

$$s_0 \otimes \alpha_1 \otimes \alpha_2 \otimes \dots \otimes \alpha_n \models \varphi_g.$$

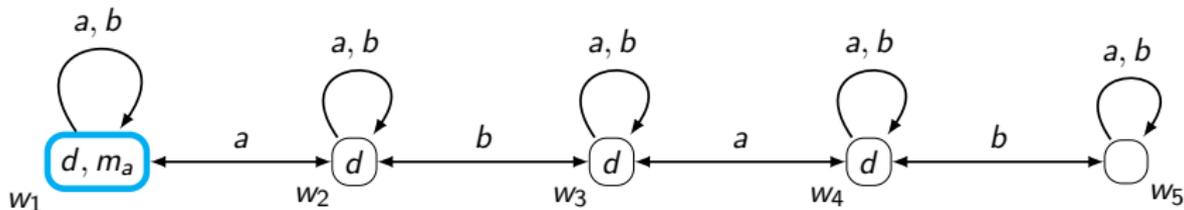
**Example.** Let  $s_0$  be the initial state of the coordinated attack problem. Let  $A = \{send_{ab}, send_{ba}\}$ . Then the following are planning tasks:

1.  $T = (s_0, A, Cd)$ , where  $C$  denotes common knowledge. It has no solution.
2.  $T = (s_0, A, E^n d)$ , where  $E$  denotes “everybody knows” and  $n \geq 1$ . It has a solution of length  $n$ .

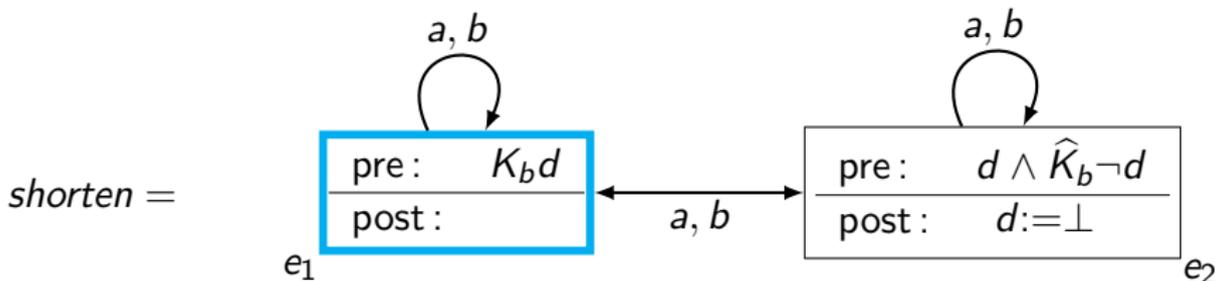
This formalisation of the coordinated attack problem is from [Bolander et al., 2019].

## Shortening the chain

Consider a chain of the form produced by the message-passing domain above:



Using preconditions of modal depth 1 we can also shorten the chain by 1:



Then it is only a short step to have multiple chains that can grow and shrink and then to encode two-counter machines  $\Rightarrow$  undecidability!

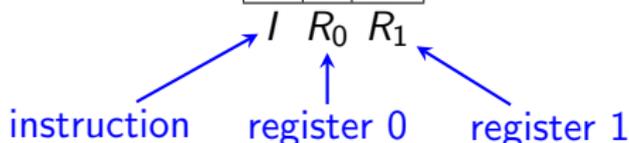
# Two-counter machines

Undecidability of the plan existence problem in epistemic planning (whether a solution exists) can be done by a reduction of the halting problem for **two-counter machines**:

**Configurations:**

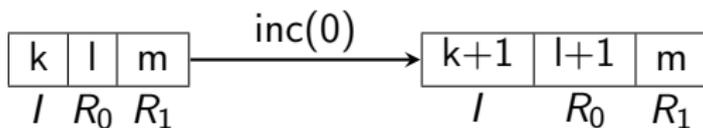
k	l	m
---	---	---

, where  $k, l, m \in \mathbb{N}$ .



**Instruction set:**  $\text{inc}(0), \text{inc}(1), \text{jump}(j), \text{jzdec}(0, j), \text{jzdec}(1, j), \text{halt}$ .

**Computation step example:**



*The halting problem for two-counter machines is undecidable*  
[Minsky, 1967].

## Proof idea for undecidability of epistemic planning

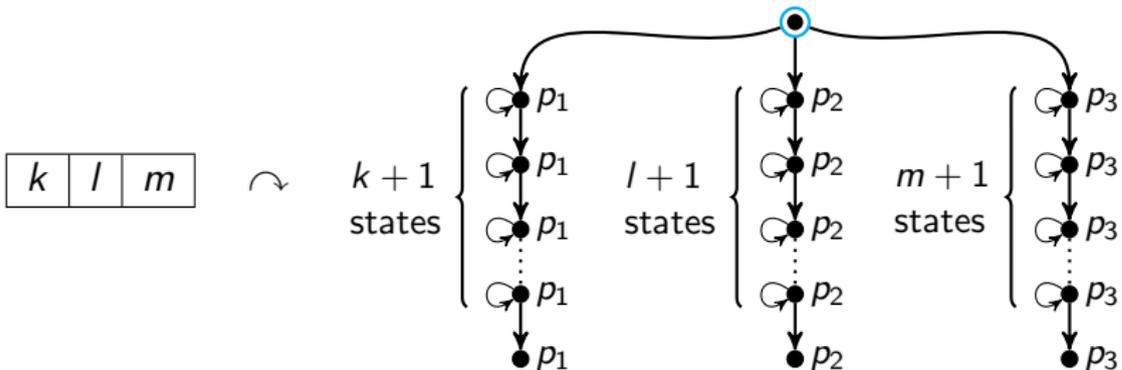
For each two-counter machine, construct a corresponding planning task where:

- The **initial state** encodes the initial configuration of the machine.
- The **epistemic actions** encode the instructions of the machine.
- The **goal formula** is true of all epistemic states representing halting configurations of the machine.

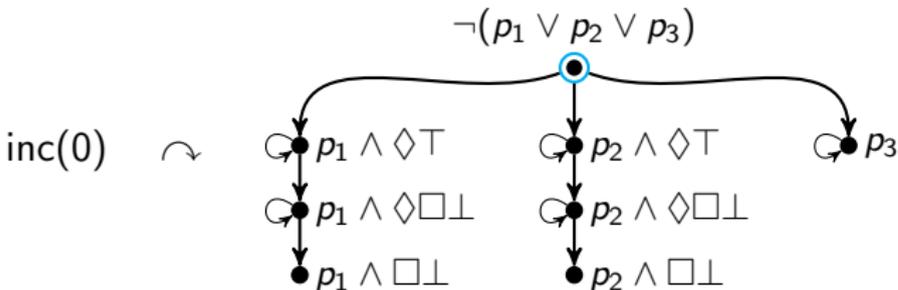
Then show that the two-counter machine halts iff the corresponding planning task has a solution. (Execution paths of the planning task encodes computations of the machine).

# Encodings

Encoding configurations as epistemic states:

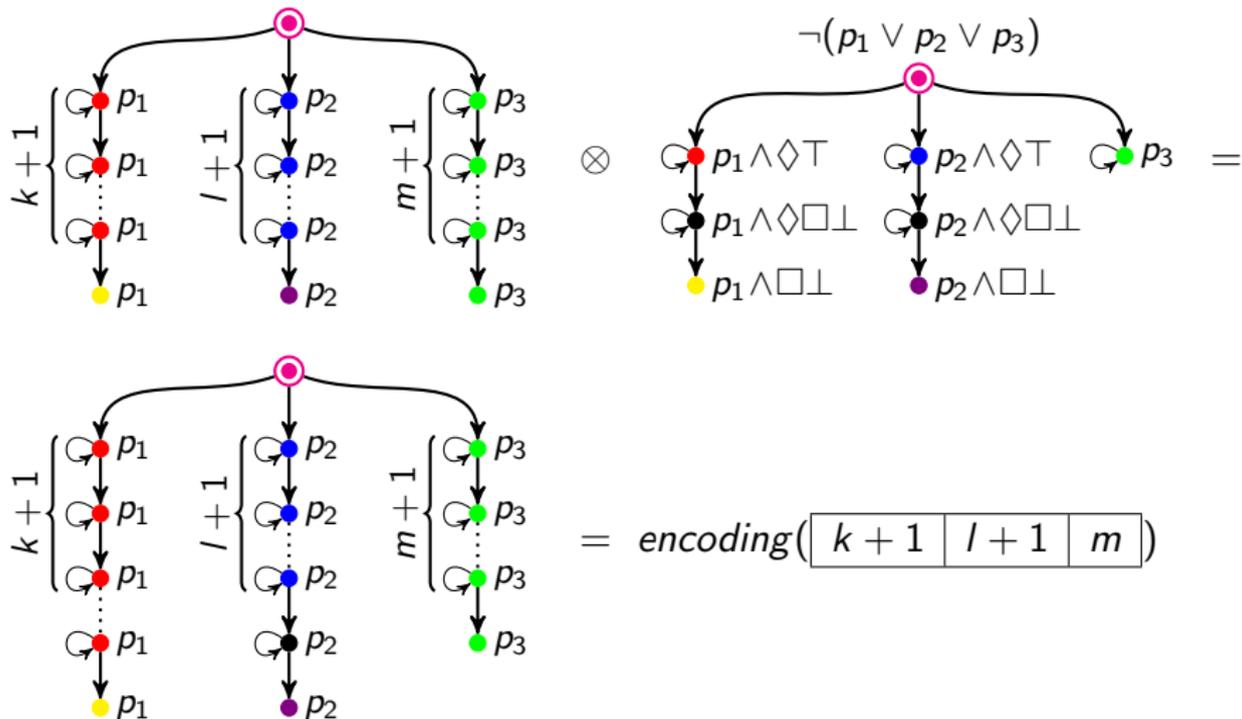


Encoding instructions as epistemic actions (note: only preconditions!):



The computation step  $\boxed{k \mid l \mid m} \xrightarrow{\text{inc}(0)} \boxed{k+1 \mid l+1 \mid m}$  is mimicked by:

$$\text{encoding}(\boxed{k \mid l \mid m}) \otimes \text{encoding}(\text{inc}(0)) =$$



## Plan existence and classes of planning tasks

**Definition.** Let  $\mathcal{T}$  be a class of planning tasks. By  $\text{PlanEx-}\mathcal{T}$  we denote the following decision problem, called the **plan existence problem** on  $\mathcal{T}$ : Given a planning task  $T \in \mathcal{T}$ , does  $T$  have a solution?

We here only consider the following classes:

- $\mathcal{T}(m, n)$  with  $m, n \in \mathbb{N} \cup \{\infty\}$ : Class of planning tasks where the preconditions are of modal depth  $\leq m$  and the postconditions are of modal depth  $\leq n$ .
- $\mathcal{T}(m, -1)$  with  $m \in \mathbb{N} \cup \{\infty\}$ : Class of planning tasks where the preconditions are of modal depth  $\leq m$  and there are no postconditions (purely epistemic).

**Example.** The coordinated attack problem is in  $\mathcal{T}(0, 0)$ . As we will later see,  $\text{PlanEx-}\mathcal{T}(0, 0)$  is decidable.

# Reductions between plan existence problems

[Bolander et al., 2019] (under submission) proves the following polynomial-time reductions for all  $m, n$ :

1.  $\text{PlanEx-}\mathcal{T}(m, n) \leq^P \text{PlanEx-}\mathcal{T}(m + k, n + l)$  for all  $k, l \geq 0$ .
2.  $\text{PlanEx-}\mathcal{T}(m, n) \leq^P \text{PlanEx-}\mathcal{T}(0, 1)$ .

- : proved decidable
- : decidable by reduction
- : proved undecidable
- : undecidable by reduction

## Decidability theorem.

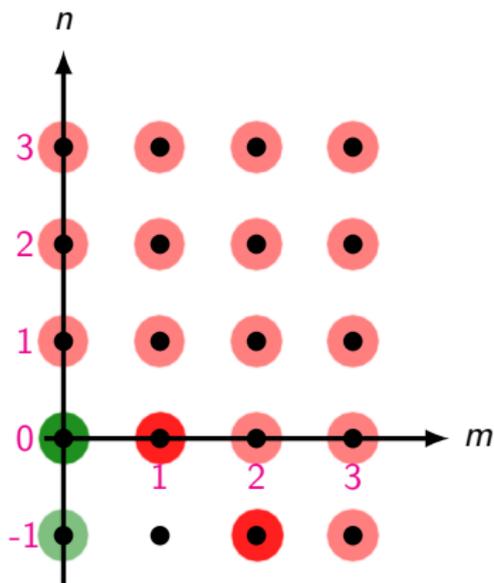
$\text{PlanEx-}\mathcal{T}(0, 0)$  is decidable.

## Undecidability theorem 1.

$\text{PlanEx-}\mathcal{T}(2, -1)$  is undecidable.

## Undecidability theorem 2.

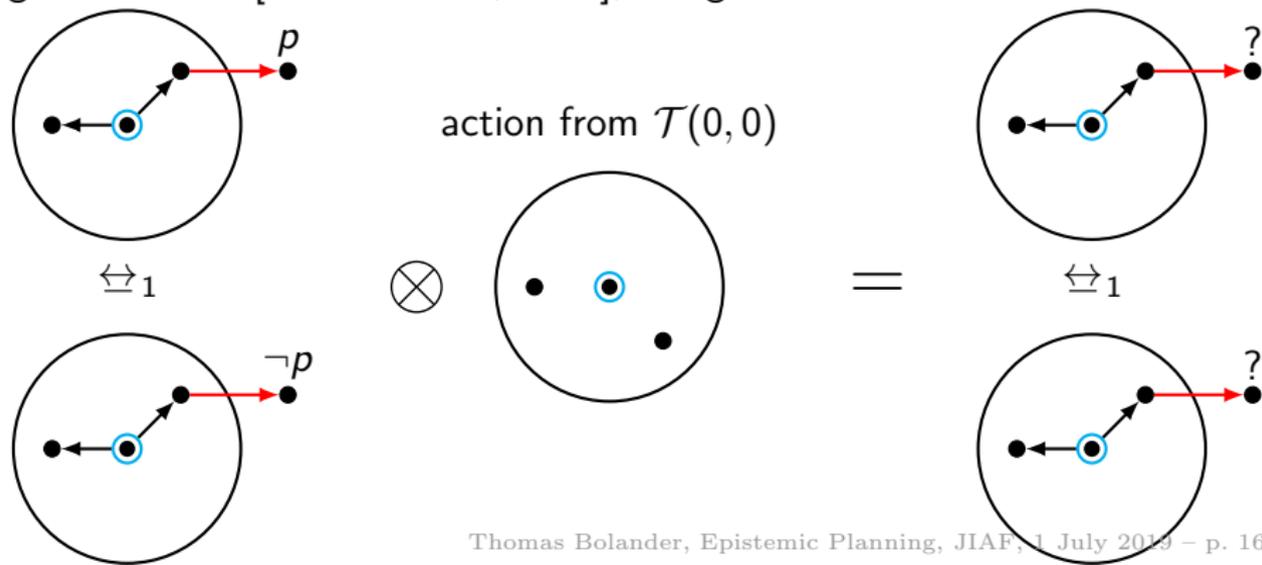
$\text{PlanEx-}\mathcal{T}(1, 0)$  is undecidable.



# Decidability theorem

**Theorem.**  $\text{PlanEx-}\mathcal{T}(0,0)$  is decidable.

**Proof idea:** Originally proved in [Yu et al., 2013], exploiting that  $k$ -bisimilarity is preserved when doing product update with epistemic actions having propositional pre- and post-conditions. Intuitively because the events of such actions can not look deeper into the model (they can only relate locally to the worlds in which they apply). The proof was generalised in [Aucher et al., 2014], using automatic structures.



# Undecidability theorem 1

Theorem ([Aucher and Bolander, 2013])

*PlanEx-T*( $\infty, -1$ ) is undecidable.

**Proof idea:** This is the two-counter machine reduction shown earlier. Preconditions of arbitrary modal depth was used to refer to—and modify—the value in the instruction counter, e.g. for jumping to another instruction.

Theorem ([Charrier et al., 2016])

*PlanEx-T*(2, -1) is undecidable.

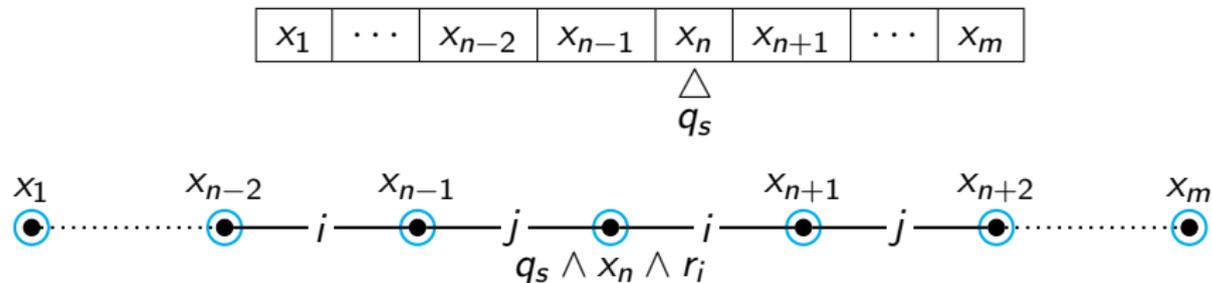
**Proof idea:** Strengthening of the proof above to avoid preconditions of arbitrary modal depth.

## Undecidability theorem 2

Theorem ([Bolander and Andersen, 2011])

$\text{PlanEx-T}(1,0)$  is undecidable.

**Proof idea:** Reduction of Halting problem for Turing machines. @fix formulation: States (epistemic models) encode IDs of the Turing machine, actions (event models) encode transitions of the Turing machine.



[Cong et al., 2018] strengthen the result, by showing that it still holds with only 2 agents and 6 propositions. The proof uses cellular automata instead of Turing machines, but otherwise uses a similar reduction.

# Epistemic planning for human-robot collaboration



Epistemic planning for knowing when when to interfere: Only provides information to the human when she has a wrong belief, and when the information is required in order for the human to be able solve the task.

- **Sub-symbolic AI** (mainly deep learning): face/object recognition, skeleton tracking, speech-to-text.
- **Symbolic AI** (epistemic logic, epistemic planning): logical reasoning, planning, perspective-taking.

# References I



**Aucher, G. and Bolander, T. (2013).**

**Undecidability in Epistemic Planning.**

In Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence (IJCAI) pp. 27–33,.



**Aucher, G., Maubert, B. and Pinchinat, S. (2014).**

**Automata Techniques for Epistemic Protocol Synthesis.**

In Proceedings 2nd International Workshop on Strategic Reasoning, (Mogavero, F., Murano, A. and Vardi, M. Y., eds), vol. 146, of Electronic Proceedings in Theoretical Computer Science pp. 97–103,.



**Bolander, T. and Andersen, M. B. (2011).**

**Epistemic Planning for Single- and Multi-Agent Systems.**

*Journal of Applied Non-Classical Logics* 21, 9–34.



**Bolander, T., Charrier, T., Pinchinat, S. and Schwarzentruher, F. (2019).**

**A roadmap of decidability and complexity results for DEL-based epistemic planning.**

Technical report Under submission.



**Charrier, T., Maubert, B. and Schwarzentruher, F. (2016).**

**On the Impact of Modal Depth in Epistemic Planning.**

In Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, USA, July 12-15, 2016.



**Cong, S. L., Pinchinat, S. and Schwarzentruher, F. (2018).**

**Small Undecidable Problems in Epistemic Planning.**

In Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden. pp. 4780–4786,.

# References II



Minsky, M. (1967).

Computation.

Prentice-Hall.



Yu, Q., Wen, X. and Liu, Y. (2013).

Multi-agent epistemic explanatory diagnosis via reasoning about actions.

In Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence (IJCAI) pp. 27–33,.