

INFOMLSAI Logics for Safe AI

Coursework 2

Coursework released: 10 May 2021, on Blackboard
Coursework due: 23:59 21 May 2021, on Blackboard
Submission format: a folder containing a pdf and an .ispl file, one per group

Please do the coursework in groups of 2-3 people. Submit a single zipped folder on Blackboard for your group. The folder should contain the pdf file and the ispl file. Please name the folder group-X, where X is the number of your group, and state in the pdf file and comments in the ispl file the names of the members of the group.

Tasks that can be done in Week 3 (w/c 10 May)

The following tasks can be done after watching the lectures on Epistemic Logic and Epistemic Logic with Common and Distributed Knowledge.

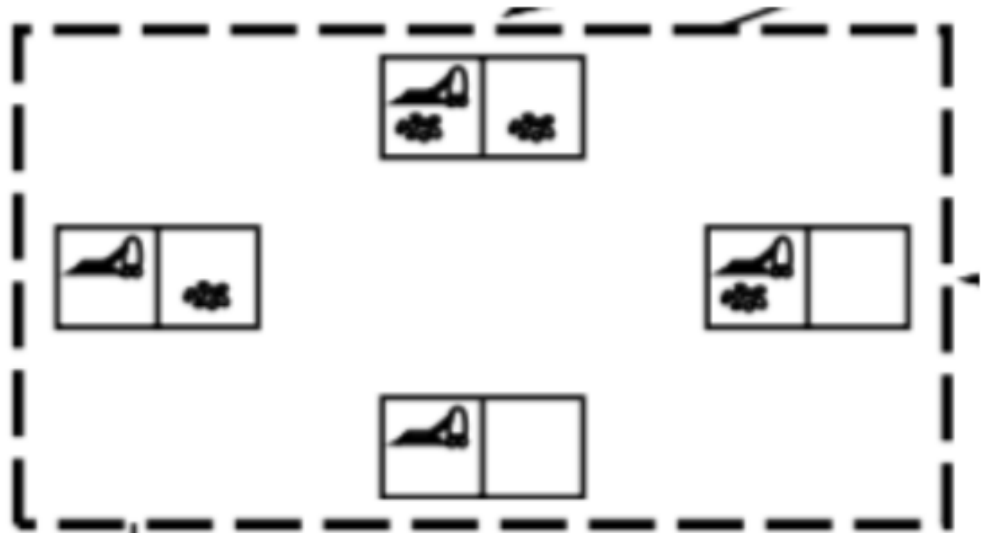
W3-1 Consider three agents, a , b and c who each were dealt one of the following cards: a card with one dot, a card with two dots, or a card with three dots. Each agent knows their card but not the cards of the other agents. This scenario was described in the lecture on epistemic logic, and an incomplete Kripke model corresponding to it is shown in the slides. Provide a completed Kripke model M_{abc} (specify the states, indistinguishability relation for all three agents, and a valuation). The propositions are $a1$ for agent a has the card with 1 dot, $a2$ for agent a has the card with 2 dots, etc. for $a3, b1, b2, b3, c1, c2, c3$.

W3-2 In the model above, in the state q_{123} where a has the card with 1 dot, b has the card with 2 dots, and c has the card with 3 dots, does it hold that:

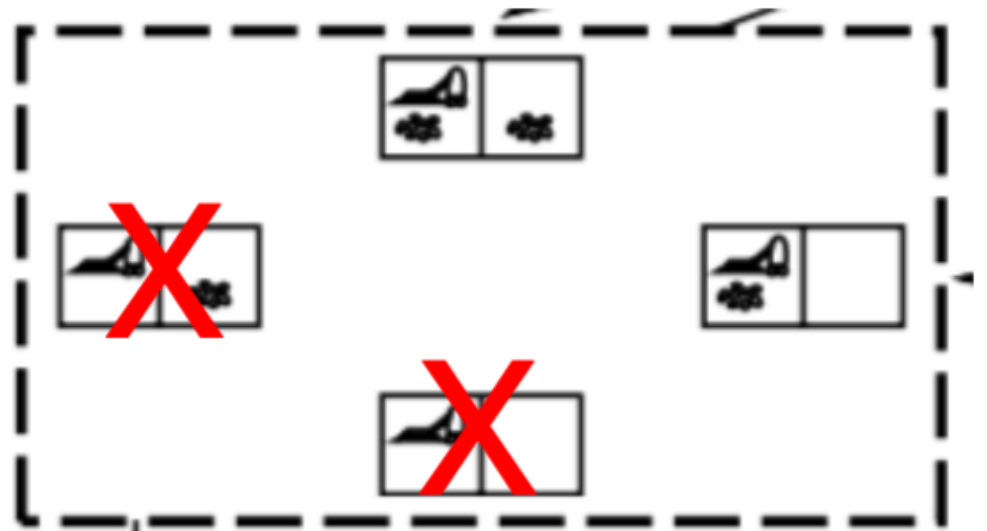
- it is distributed knowledge between a , b and c that $a1 \wedge b2 \wedge c3$. Express this in $ELCD$ and argue why it is true (or false).
- it is common knowledge between a , b and c that $a1 \vee b2 \vee c3$. Express this in $ELCD$ and argue why it is true (or false).
- it is common knowledge between a , b and c that $a1 \vee a2 \vee a3$. Express this in $ELCD$ and argue why it is true (or false).

W3-3 A basic action in epistemic planning is a *truthful public announcement*: some formula φ is announced, and as a result, each agent in the system considers

possible only the states where φ is true. For example, if the vacuum cleaner agent hears the truthful public announcement $\neg \text{cleanA}$, it will eliminate the states where cleanA holds, and will be left with just two possibilities:



↓ announcement of $\neg \text{cleanA}$



Formally, a truthful public announcement of φ is applied to a pair of a Kripke model $M = (St, \{\sim_i: i \in Agt\}, \mathcal{V})$ and a state $q \in St$ such that $M, q \models \varphi$ ¹ and produces an updated Kripke model $M^\varphi = (St^\varphi, \{\sim_i^\varphi: i \in Agt\}, \mathcal{V}^\varphi)$ where $St^\varphi = St \cap \{q' \mid M, q' \models \varphi\}$, and for each i , \sim_i^φ is \sim_i restricted to the remaining states, and \mathcal{V}^φ is \mathcal{V} restricted to the remaining states. For more background, see the entry on Public Announcement Logic in the Stanford Encyclopaedia of Philosophy.

Question: In the Kripke model M_{abc} above, and the state q_{123} , is there a truthful public announcement φ such that after this announcement b and c know which card each agent has but a does not know which cards b and c have?

Formally, is there a formula φ such that:

- $M_{abc}, q_{123} \models \varphi$
- $M_{abc}^\varphi, q_{123} \models E_{\{b,c\}}(a1 \wedge b2 \wedge c3)$
- $M_{abc}^\varphi, q_{123} \models \neg K_a(b2 \wedge c3)$

If yes, give the formula and describe M_{abc}^φ . Does φ also guarantee common knowledge $(C_{\{b,c\}}(a1 \wedge b2 \wedge c3))$?

Tasks that can be done during Week 4 (w/c 17 May)

The following tasks can be done after watching the lectures on interpreted systems, CTLK, and model checking CTLK.

W4-1 Represent the scenario of two robots and a carriage where the first robot can detect the colour of the surface and the second robot the texture, as an interpreted system in ISPL. The scenario is like the one in the reader on p.44, but both robots can push the carriage. Robot 1 pushes clockwise, and robot 2 pushes counter-clockwise. If both robots push at the same time, the carriage does not move. Define group g as $\{\text{Robot1}, \text{Robot2}\}$. The following formulas should be true in the initial state q_0 :

Formula	ISPL Property
$\text{pos2} \rightarrow \neg K_1 \text{pos2}$	$\text{pos2} \rightarrow \neg !K(\text{Robot1}, \text{pos2}) ;$
$\text{pos2} \rightarrow K_1 \neg \text{pos1}$	$\text{pos2} \rightarrow K(\text{Robot1}, \neg \text{pos1}) ;$
$\text{pos2} \rightarrow K_2 K_1 \neg \text{pos1}$	$\text{pos2} \rightarrow K(\text{Robot2}, K(\text{Robot1}, \neg \text{pos1})) ;$
$\text{pos2} \rightarrow D_{\{1,2\}} \text{pos2}$	$\text{pos2} \rightarrow DK(g, \text{pos2}) ;$
$\text{pos2} \rightarrow \neg E_{\{1,2\}} \text{pos2}$	$\text{pos2} \rightarrow \neg !GK(g, \text{pos2}) ;$
$\text{pos2} \rightarrow E_{\{1,2\}} \neg \text{pos1}$	$\text{pos2} \rightarrow GK(g, \neg \text{pos1}) ;$
$\text{pos2} \rightarrow \neg C_{\{1,2\}} \text{pos2}$	$\text{pos2} \rightarrow \neg !GCK(g, \text{pos2}) ;$
$\text{pos2} \rightarrow \neg C_{\{1,2\}} \neg \text{pos1}$	$\text{pos2} \rightarrow \neg !GCK(g, \neg \text{pos1}) ;$
$AGC_{\{1,2\}}(\text{pos0} \vee \text{pos1} \vee \text{pos2})$	$AG \text{ GCK}(g, (\text{pos0} \text{ or } \text{pos1} \text{ or } \text{pos2}))$
$EFE_{\{1,2\}} \neg \text{pos1}$	$EF \text{ GK}(g, \neg \text{pos1})$

¹ (M, q) is called a pointed Kripke model

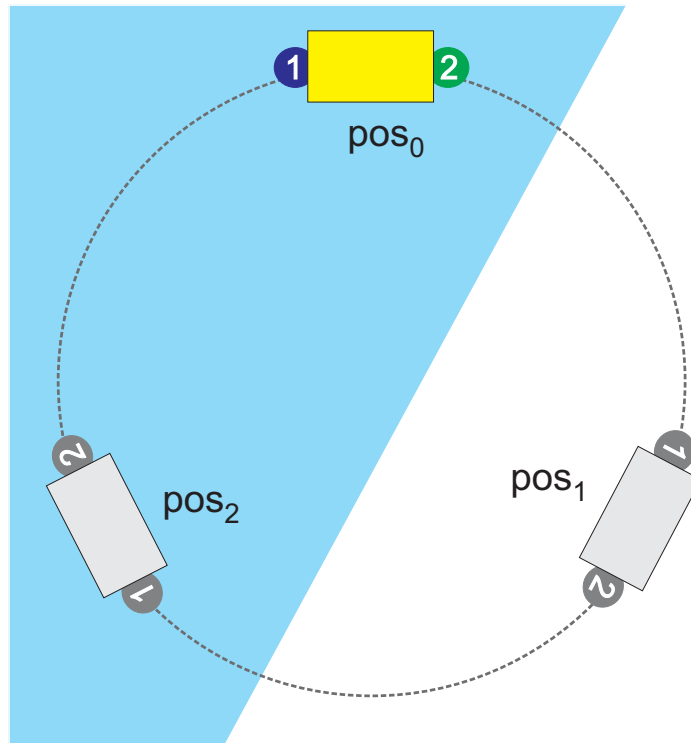


Figure 1: Two robots and a carriage

W4-2 Give a model checking algorithm for CTLK that uses only the existential preimage function (i.e., universal properties are checked without using the universal preimage).