

Engineering Social Order in Multi-Agent Systems

Stephen Cranefield

School of Computing, University of Otago
Dunedin, New Zealand

EMAS@AAMAS 2024
7 May 2024

This work was supported by the Marsden Fund Council from New Zealand Government funding, managed by Royal Society Te Apārangi.

Engineering social order via peer-to-peer interactions

My work has focused on what *agents* themselves (might) need to coordinate/cooperate, rather than creating centralised services or hierarchical social structures. A few examples:

- ▶ Individual modelling and monitoring of social expectations, e.g. a monitor service for Jason agents and application to the Second Life virtual world.
- ▶ Learning existing norms from observation of interactions in a society: data mining and Bayesian approaches.
- ▶ Choosing agent plans to maximise a human partner's value fulfilment.
- ▶ Extending BDI agents to follow (predefined) social practices.
- ▶ Peer to peer proposal and execution of group plans, supported by decentralised middleware.
- ▶ And now ... individual agent reasoning about common knowledge.

Examples of common knowledge (within varying groups)

Examples of common knowledge (within varying groups)

- ▶ There are four seasons in a year.

Examples of common knowledge (within varying groups)

- ▶ There are four seasons in a year.
- ▶ After lightening we will hear thunder.

Examples of common knowledge (within varying groups)

- ▶ There are four seasons in a year.
- ▶ After lightening we will hear thunder.
- ▶ Christopher Luxon is the prime minister of New Zealand.

Examples of common knowledge (within varying groups)

- ▶ There are four seasons in a year.
- ▶ After lightening we will hear thunder.
- ▶ Christopher Luxon is the prime minister of New Zealand.
- ▶ Taylor Swift and Joe Alwyn have broken up.

Examples of common knowledge (within varying groups)

- ▶ There are four seasons in a year.
- ▶ After lightening we will hear thunder.
- ▶ Christopher Luxon is the prime minister of New Zealand.
- ▶ Taylor Swift and Joe Alwyn have broken up.



Examples of common knowledge (within varying groups)

- ▶ There are four seasons in a year.
- ▶ After lightening we will hear thunder.
- ▶ Christopher Luxon is the prime minister of New Zealand.
- ▶ Taylor Swift and Joe Alwyn have broken up.

▶



- ▶ The band rocks.



Informal definition of common knowledge

Proposition φ is common knowledge if:

- ▶ everyone knows φ
- ▶ everyone knows that everyone knows φ
- ▶ everyone knows that everyone knows that everyone knows φ
- ▶ ...

The coordinated attack problem (Fagin et al.)

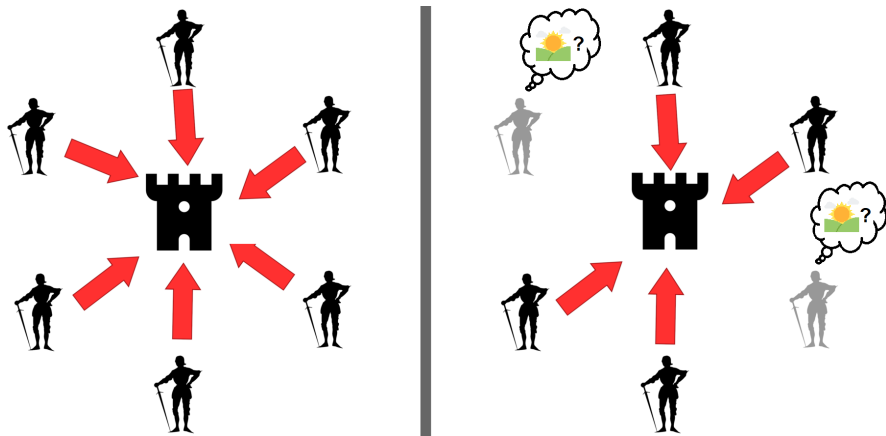
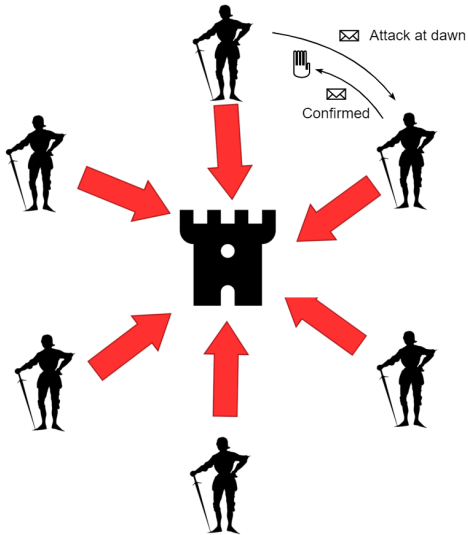


Image derived from Byzantine Generals.png by Lord Belbury. Licence CC BY-SA 4.0

https://en.wikipedia.org/wiki/Byzantine_fault#/media/File:Byzantine_Generals.png.

Dawn icon by Freepik, <https://www.flaticon.com/free-icons/dawn>.

The coordinated attack problem (Fagin et al.)

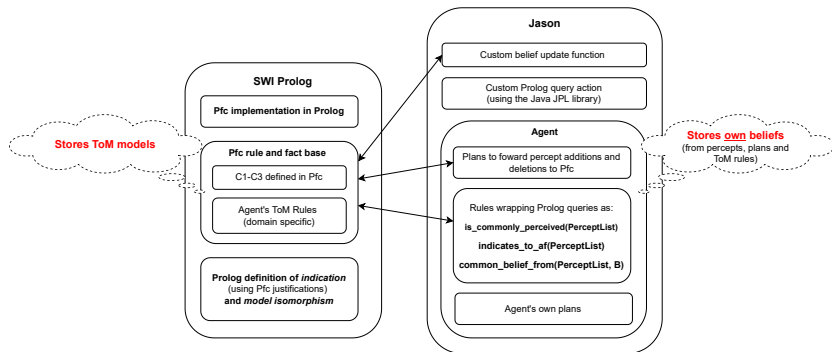


Common knowledge
cannot be achieved via
asynchronous
communication

Why do we need common knowledge?

- ▶ Basic assumption of game theory: the payoff structure and the rationality of all players are common knowledge
- ▶ Conventions (Lewis). In instances of a coordination problem S , it is *common knowledge* that:
 - ▶ There is some regularity of behaviour R that everyone conforms to.
 - ▶ Everyone expects everyone else to conform to R .
 - ▶ Everyone prefers to conform to R on condition that the others do since R is a coordination equilibrium in S .
- ▶ Definition of “We-mode” thinking in a group (e.g. Tuomela 2007)
- ▶ ...

Where I am heading: the agent engineering outcome



(Some) theories of common knowledge:

1) Extension of epistemic logic (Fagin et al.)

$K_i\varphi$ means “Agent i knows φ ”

$E_G\varphi := \bigwedge_{i \in G} K_i\varphi$: Everyone in group G knows φ

$C_G\varphi := \bigwedge_{i=0}^{\infty} E^i\varphi$: It is common knowledge in G that φ

where $E_G^n\varphi := E_G E_G^{n-1}\varphi$ and $E_G^0\varphi := \varphi$

To avoid an infinite conjunction, $C_G\varphi$ can be defined using the Fixed-Point Axiom:

$$C_G\varphi \leftrightarrow E_G(\varphi \wedge C_G\varphi)$$

and induction rule¹:

$$\frac{\text{If } \varphi \rightarrow E_G(\psi \wedge \varphi) \text{ then infer } \varphi \rightarrow C_G\psi$$

¹“The antecedent gives us the essential ingredient for proving, by induction on k , that $\varphi \rightarrow E_G^k(\psi \wedge \varphi)$ is valid for all k ” (Fagin et al.)

Problems (?) with the epistemic logic approach

- ▶ Artemov (2004):
 - ▶ “This kind of deductive system does not behave well proof-theoretically. This practically rules out automated proof search and severely limits the usage of formal methods in analyzing knowledge.”
 - ▶ “. . . this [is] the most liberal version of knowledge operator satisfying the Fixed Point Axiom, without imposing any conditions on the way this knowledge is attained. . . there might be nonconstructive versions of the common knowledge appearing by chance or for some unknown reasons or without any particular reasons at all.”
- ▶ We can (e.g.) combine a logic of common knowledge with *public announcement logic*, with an inference rule that infers common knowledge from a public announcement²
 - ▶ But how do we decide *what* counts as a public announcement?
 - ▶ Do we need add-on logics for other ways of creating common knowledge?
- ▶ Why don't see practical agent systems using it.

²<https://plato.stanford.edu/entries/dynamic-epistemic/#CommKnow>

(Some) theories of common knowledge:

2) David Lewis (1969)

- ▶ Lewis focuses on situations when a certain *state of affairs* A “indicates” that a proposition P holds.
 - ▶ Example (Lewis):
You said you will return tomorrow to continue our meeting
indicates that
you will return.
 - ▶ Example (Cubitt & Sugden):
The room we are in is lit by a flash of lightening
indicates that
within a few seconds, there will be the noise of thunder.
- ▶ Lewis defines (informally) three *sufficient* conditions for the indicator A to be a *basis for common knowledge of P* .

(Some) theories of common knowledge:

2) David Lewis (1969)

Proposition P is *common knowledge* if and only if there is some state of affairs A that holds and:

Everyone has reason to believe that A holds. (C1)

A indicates to everyone that everyone has reason to believe that A holds. (C2)

A indicates to everyone that P holds. (C3)

Plus “suitable ancillary premises regarding our rationality, inductive standards, and background information”

(Some) theories of common knowledge:

2) David Lewis (1969)

Proposition P is *common knowledge* if and only if there is some **state of affairs** A that holds and:

Everyone has **reason to believe** that A holds. (C1)

A **indicates**¹ to everyone that everyone has reason to believe that A holds. (C2)

A indicates to everyone that P holds. (C3)

Plus “suitable ancillary premises regarding our rationality, inductive standards, and background information”

¹ A indicates φ := If someone has reason to believe that A holds, they thereby have reason to believe φ (discussed later).

(Some) theories of common knowledge:

2) David Lewis (1969)

Proposition P is *common knowledge* if and only if there is some state of affairs A that holds and:

Everyone has reason to believe that A holds. (C1)

☞ Cubitt & Sugden: “ A is self-revealing”

A indicates to everyone that everyone has reason to believe that A holds. (C2)

☞ C&S: “ A is public”

A indicates to everyone that P holds. (C3)

☞ Me: “ A is objective”

Plus “suitable ancillary premises regarding our rationality, inductive standards, and background information”

(Some) theories of common knowledge:

2) David Lewis (1969)

- ▶ Informal proof: Given some A and P such that C1, C2 and C3 hold, there exists an infinite chain of reasoning that creates all levels of nested reasons to believe:
 - i has reason to believe that j has reason to believe that k has reason to believe ... that P .
- ▶ The proof doesn't depend on the content of A and P —just the properties C1, C2 and C3.
- ▶ The proof can be recast using mathematical induction There is no need to perform an infinite chain of reasoning.

Lewis's informal analysis

II | Convention Refined

I. Common Knowledge

Agreement, reliance, or preference: how come, can we have a coordination problem by producing systems of conventions that are higher-order mutual expectations. We need only imagine cases in which a coordination problem that higher-order expectations would be resolved that way? What problem have we to justify as an underlying that often have certain expectations, and how can we have the first few orders?

Take a simple case of coordination by agreement. Suppose the following state of affairs—call it *A*—holds: you and I have met, we have been talking together, you must have before me business is done, we may now say goodbye to the other. In such a situation, I expect you to return. You will expect me to expect you to return. Perhaps there will be no other cases.

- (1) You and I have reason to believe that *A* holds.
- (2) I believe to both of us that you and I have reason to believe that *A* holds.
- (3) I believe to both of us that you will return.

What is interesting? Let us say that *A* induces us to expect *A*.

52

COMMON KNOWLEDGE | 53

—I and only if *A*'s had reason to believe that *A* holds, *A* would surely have reason to believe that What *A* induces us to will depend, therefore, on *A*'s inductive standards and background information.

The first main position (1), (2), (3), together with suitable auxiliary premises regarding our reliability, inductive standards, and background information, suffice to justify my higher-order expectations. Let us now have my reasoning spelled out. Consider that if *A* induces me to will, and if I share *A*'s inductive standards and background information, then I must indicate the way things go. Therefore, if *A* induces us to *A*, *A* has reason to believe that *A* holds, and all of us believe that *A* holds. And if *A* has reason to believe that *A* holds, and *A* has inductive standards and background information, that *A* induces us to *A* has reason to believe that if *A* has reason to believe that *A* holds, I believe that *A* holds. Suppose you and I do have reason to believe that *A* holds, and I will indicate the same things to both of us. Thus (2) applied to (3) implies:

- (4) I believe to both of us that each of us has reason to believe that the other has reason to believe that you will return.
- (5) (2) applied in turn to (4) implies:
- (6) I believe to both of us that each of us has reason to believe that the other has reason to believe that you will return.
- (7) And so on, ad infinitum, since such an induction holds *A* induces us to both of us that Note that this is a chain of implications, not of steps in empirical causal reasoning. Therefore there is nothing improper about infinite length. Figure 26 is a more detailed representation of these implications in my case; there is your case and it would be symmetrical.

Consider now that *A* induces us to believe *A*, a principle of generalization of *A* induces us to and *A* has reason to

54 | COORDINATION PROBLEM

I believe that *A* holds, then I have reason to believe that Perhaps (3) applied in this way to (1) implies:

(7) Each of us has reason to believe that the other has reason to believe that the first has reason to believe that you will return.

- (7) Each of us has reason to believe that the other has reason to believe that the first has reason to believe that you will return.
- (8) And so on, for the whole infinite sequence we considered above. I can still indicate that only my own *A* is induced, but I am aware that you believe that *A* holds, and that you expect me to expect you to return. I have reason to believe that *A* holds, and I will indicate the same things to both of us. Thus (2) applied to (3) implies:

54

COMMON KNOWLEDGE | 55

and background information, reliability, inductive standards of reliability, and so on. Let us return to our example and consider the state of affairs *A* induces us to expect. Suppose that as part of *A* we maintain our conditional probability for returning to the starting place. There is more also in this case so that both have our preferences. If, *A* can serve as a basis for common knowledge that each of us prefers to return if the other does. Suppose also that as part of *A* we maintain a reliability of induction of reliability. Therefore there is nothing improper about infinite length. Figure 26 is a more detailed representation of these implications in my case; there is your case and it would be symmetrical.

Consider now that *A* induces us to believe *A*, a principle of generalization of *A* induces us to and *A* has reason to

56 | COORDINATION PROBLEM

believe that *A* holds, then *A* has reason to believe that Perhaps (3) applied in this way to (1) implies:

(7) Each of us has reason to believe that the other has reason to believe that the first has reason to believe that you will return.

(7) Each of us has reason to believe that the other has reason to believe that the first has reason to believe that you will return.

56 | COORDINATION PROBLEM

each has reason to expect that the other has reason to expect that he has reason that you will return. If, in addition, each of us has reason to another a sufficient degree of reliability in the other, then each has reason to expect that the other expects that he expects that you will return. And if *A* induces us to believe *A*, a principle of generalization of *A* induces us to and *A* has reason to

56

56

COMMON KNOWLEDGE | 57

and background information, reliability, inductive standards of reliability, and so on. Let us return to our example and consider the state of affairs *A* induces us to expect. Suppose that as part of *A* we maintain our conditional probability for returning to the starting place. There is more also in this case so that both have our preferences. If, *A* can serve as a basis for common knowledge that each of us prefers to return if the other does. Suppose also that as part of *A* we maintain a reliability of induction of reliability. Therefore there is nothing improper about infinite length. Figure 26 is a more detailed representation of these implications in my case; there is your case and it would be symmetrical.

Consider now that *A* induces us to believe *A*, a principle of generalization of *A* induces us to and *A* has reason to

58 | COORDINATION PROBLEM

believe that *A* holds, then *A* has reason to believe that Perhaps (3) applied in this way to (1) implies:

(7) Each of us has reason to believe that the other has reason to believe that the first has reason to believe that you will return.

58 | COORDINATION PROBLEM

each has reason to expect that the other has reason to expect that he has reason that you will return. If, in addition, each of us has reason to another a sufficient degree of reliability in the other, then each has reason to expect that the other expects that he expects that you will return. And if *A* induces us to believe *A*, a principle of generalization of *A* induces us to and *A* has reason to

58

COMMON KNOWLEDGE | 55

believe that *A* holds, then I have reason to believe that Perhaps (3) applied in this way to (1) implies:

(7) Each of us has reason to believe that the other has reason to believe that the first has reason to believe that you will return.

(7) Each of us has reason to believe that the other has reason to believe that the first has reason to believe that you will return.

55

COMMON KNOWLEDGE | 59

each has reason to expect that the other has reason to expect that he has reason that you will return. If, in addition, each of us has reason to another a sufficient degree of reliability in the other, then each has reason to expect that the other expects that he expects that you will return. And if *A* induces us to believe *A*, a principle of generalization of *A* induces us to and *A* has reason to

59

59

COMMON KNOWLEDGE | 57

and background information, reliability, inductive standards of reliability, and so on. Let us return to our example and consider the state of affairs *A* induces us to expect. Suppose that as part of *A* we maintain our conditional probability for returning to the starting place. There is more also in this case so that both have our preferences. If, *A* can serve as a basis for common knowledge that each of us prefers to return if the other does. Suppose also that as part of *A* we maintain a reliability of induction of reliability. Therefore there is nothing improper about infinite length. Figure 26 is a more detailed representation of these implications in my case; there is your case and it would be symmetrical.

Consider now that *A* induces us to believe *A*, a principle of generalization of *A* induces us to and *A* has reason to

Figure 26

Cubitt and Sugden's formal version of Lewis's analysis

Notation:

- ▶ $R_i(p)$: i has reason to believe p .
- ▶ $A \text{ ind}_i P$: A indicates to i that P
- ▶ Cubitt and Sugden's give *four* conditions for A to create common knowledge of P :

For all persons i : $A \text{ holds} \Rightarrow R_i(A \text{ holds})$. (CS1)

For all persons i, j : $A \text{ ind}_i R_j(A \text{ holds})$. (CS2)

For all persons i : $A \text{ ind}_i P$. (CS3)

For all persons i, j :, for all propositions Q :
 $(A \text{ ind}_i Q) \Rightarrow R_i(A \text{ ind}_j Q)$. (CS4)

Condition CS4 was implicit in Lewis's text as "suitable ancillary premises regarding our [shared] rationality, inductive standards, and background information".

- ▶ Reasons to believe can be arbitrarily nested: $R_i(R_j(\dots))$.
How can we verify CS4 for all such Q in finite time?

Cubitt and Sugden's formal version of Lewis's analysis

- ▶ Neither Lewis nor C&S provide specific semantics for the indication relationship.
- ▶ C&S state:

“Lewis clearly intends **if ... thereby ...** to be stronger than the material implication, \Rightarrow . On the most natural reading of the definition of ‘ $A \text{ ind}_i x$ ’, i 's reason to believe that A holds provides i 's reason for believing that x is true.”
- ▶ They present six properties that capture their intuition about the requirements for *any* indication relationship.
- ▶ We only need two of them (P1 and P6), but I won't discuss these further today.

Cubitt and Sugden's formal version of Lewis's analysis

Their version of Lewis's proof:

Consider any state of affairs A , any proposition P , and any population \mathcal{P} . Suppose that A holds and that A is a reflexive common indicator in \mathcal{P} that P . Then:

1. $\forall i \in \mathcal{P}, R_i(A \text{ holds})$ (from C1)
2. $\forall i, j \in \mathcal{P}, A \text{ ind}_i R_j(A \text{ holds})$ (from C2)
3. $\forall i \in \mathcal{P}, A \text{ ind}_i P$ (from C3)
4. $\forall i \in \mathcal{P}, R_i(P)$ (from 1 and 3, using P1)
5. $\forall i, j \in \mathcal{P}, R_i(A \text{ ind}_j P)$ (from 3, using C4)
6. $\forall i, j \in \mathcal{P}, A \text{ ind}_i R_j(P)$ (from 2 and 5, using P6)
7. $\forall i, j \in \mathcal{P}, R_i[R_j(P)]$ (from 1 and 6, using P1)
8. $\forall i, j, k \in \mathcal{P}, R_i(A \text{ ind}_j R_k(P))$ (from 6, using C4)
9. $\forall i, j, k \in \mathcal{P}, A \text{ ind}_i R_j[R_k(P)]$ (from 2 and 8, using P6)
10. $\forall i, j, k \in \mathcal{P}, R_i[R_j[R_k(P)]]$ (from 1 and 9, using P1)
11. $\forall i, j, k, l \in \mathcal{P}, R_i(A \text{ ind}_j R_k(R_l(P)))$ (from 8, using C4)

“and so on”

Our approach and notation (1)

Claim:

To reason about common knowledge, agents need a mechanism for *theory-of-mind* reasoning.

Our approach and notation (2)

- ▶ Each agent can choose to maintain a set of named models of other's percepts, beliefs and ToM rules.
- ▶ \odot denotes the agent's top-level model.
- ▶ af denotes "any fool" (McCarthy 1978)³.
- ▶ $\odot \gg af$ denotes the agent's model of any fool's percepts, beliefs and rules.
- ▶ $\odot \gg af \gg af$ denotes the agent's model of any fool's model about any (other) fool's percepts, beliefs and rules.
- ▶ Example percept and belief propositions:

percept(\odot , *colour*(*sky*, *blue*))

percept($\odot \gg af$, *colour*(*sky*, *blue*))

bel($\odot \gg af$, *colour*(*sky*, *blue*))

bel($\odot \gg af \gg af$, *colour*(*sky*, *blue*))

³In our approach, af is a Skolem constant. A set of "af scope percepts" can be declared to provide a restricted scope for af .

Our approach and notation (3)

- ▶ Agents have *theory-of-mind* (ToM) rules that can create new percepts, beliefs and rules in models, e.g.:

Believe what you perceive

$$\textit{percept}(\odot, P) \Rightarrow \textit{bel}(\odot, P)$$

Citizens believe they are citizens

$$\textit{bel}(\odot, \textit{citizen}(C)) \Rightarrow \textit{bel}(\odot \gg C, \textit{citizen}(me))$$

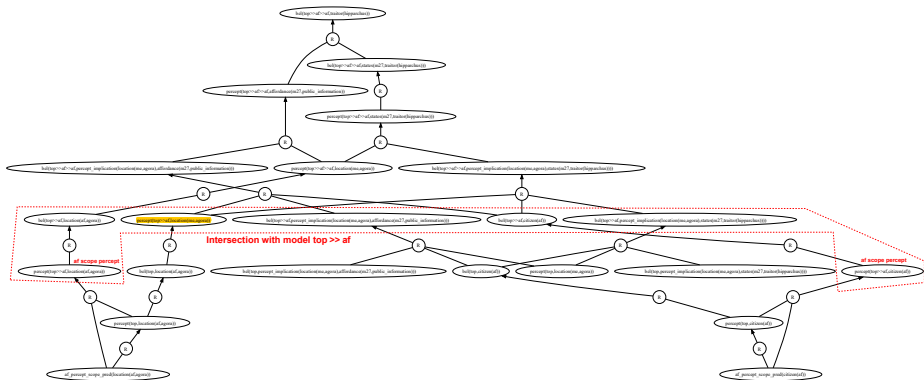
- ▶ We interpret states of affairs as sets of percepts.
- ▶ We write indication as:

$\textit{percepts}(M, A)$ ind ψ

where ψ is $\textit{percepts}(M', \dots)$, $\textit{percept}(M', \dots)$ or $\textit{bel}(M', \dots)$.
 M' is M or a nested model $M \gg \dots \gg Ag$.

- ▶ We interpret indication as stating that perceiving A in M provides *sufficient* conditions within model M , (in conjunction with the *af* scope percepts), to infer ψ using the ToM rules.

Example proof tree to determine indication



Our versions of conditions C1 to C3

When A is perceived, it is believed that any fool perceives A .

$$\text{percepts}(\odot, A) \rightarrow \text{percepts}(\odot \gg af, A) \quad (\text{C1})$$

Believing that any fool perceives A^* is sufficient to infer that any fool believes any fool perceives A .

$$\text{percepts}(\odot \gg af, A^*) \text{ ind } \text{percepts}(\odot \gg af \gg af, A) \quad (\text{C2})$$

(A^* is A augmented with the af scope percepts.)

Believing that any fool perceives A^* is sufficient to infer that any fool believes P .

$$\text{bel}(\odot \gg af, A) \text{ ind } \text{bel}(\odot \gg af, P) \quad (\text{C3})$$

Our version of condition C4

A specialised version of the C&S version: precisely what their proof needs.

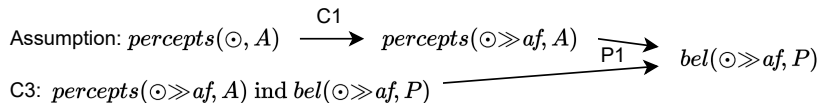
$$\begin{aligned} \forall n \geq 1: & \text{percepts}(\odot \gg af, A) \text{ ind } \text{bel}(\odot(\gg af)^n, P) \\ & \rightarrow \text{percepts}(\odot \gg af \gg af, A) \text{ ind } \text{bel}(\odot(\gg af)^{n+1}, P) \end{aligned} \quad (\text{C4})$$

This checks that whenever the first indication relationship holds (for any level of nesting n), the equivalent one with an extra “ $\gg af$ ” on each side must also hold.

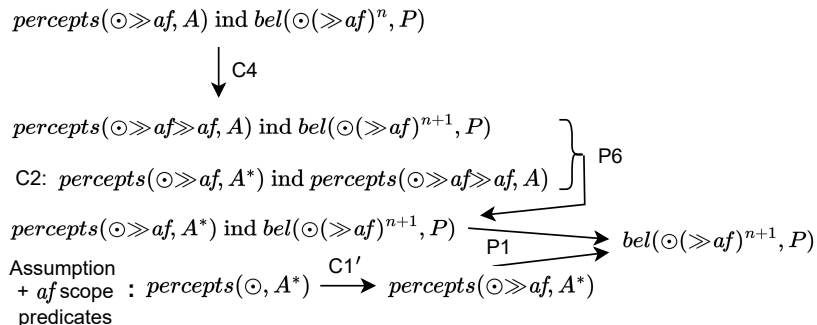
- ▶ **Problem:** This version still cannot be verified using a finite set of ToM models.
- ▶ **Solution:** We proved that C4 holds if the models $\odot \gg af$ and $\odot \gg af \gg af$ are *isomorphic*, i.e. they have the same percepts, beliefs and rules (except for the difference in model names).
- ▶ **Result:** Only two levels of ToM modelling are necessary to decide whether P is (Lewisian) common knowledge, given a set of percepts A .

Our *inductive* proof that C1–C4 lead to common knowledge of P

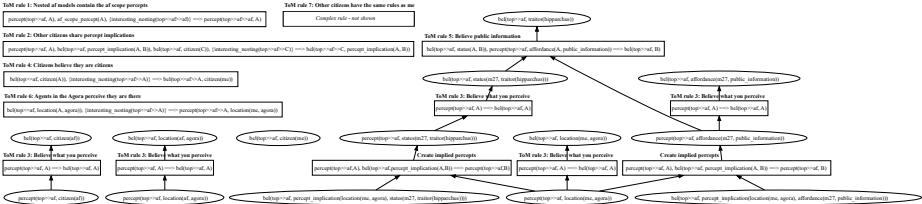
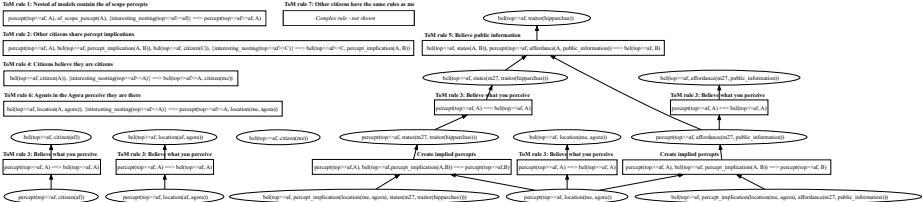
Base case



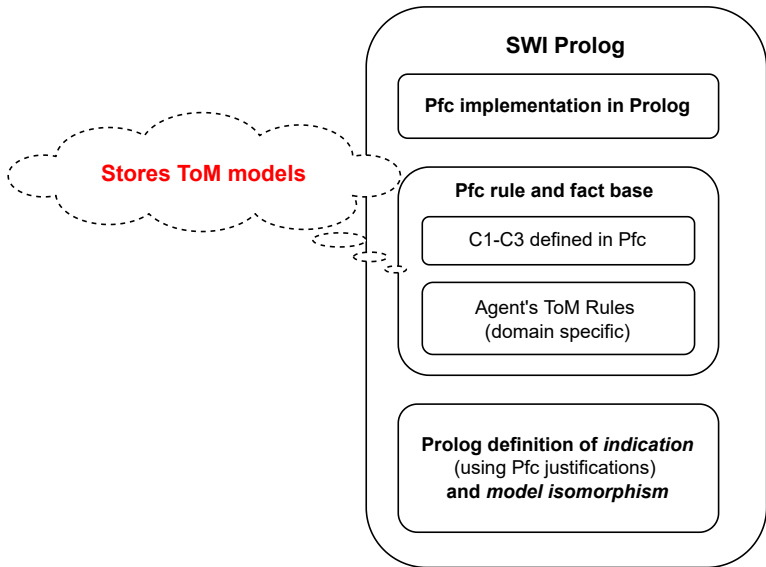
Inductive step



Example model structure comparison for the isomorphism test



Implementation architecture



Common knowledge conditions as Pfc *backward-chaining* rules

```
c1(A) <==  
  { forall(member(Ai, A),  
    (percept(top, Ai), percept(top>>af, Ai))) }.  
  }
```

```
c2(A) <==  
  { findall(P, af_scope_percept(P), Ps),  
    union(A, Ps, AfAugmentedA),  
    percepts(top>>af, AfAugmentedA) ind  
    percepts(top>>af>>af, A) }.  
  }
```

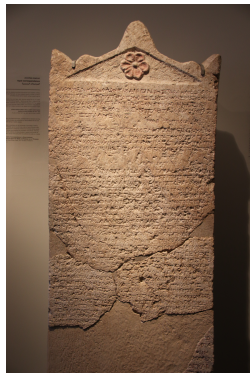
```
c3(A,P) <==  
  { percepts(top>>af, A) ind bel(top>>af, P) }.  
  }
```

```
isomorphic_models(M1, M2) :-  
  % Definition in Prolog too long to include
```

```
ck(P) <==  
  { percepts(top>>af, A) ind bel(top>>af, P) },  
  c1(A), c2(A), c3(A,P),  
  { isomorphic_models(top>>af, top>>af>>af) }.  
  }
```

Example scenario

- ▶ Part of a complex prosecutor's argument in a trial for treason in classical Athens (Ober 2010)
 - ▶ The prosecutor argued that what happens to traitors is common knowledge.
 - ▶ ... because it is inscribed on a monument in the Agora.
 - ▶ How can an agent infer that this is common knowledge?



Agora image by Ancient History Magazine / Karwansaray Publishers,

Scenario ToM rules in Pfc

```
% af scope predicate declarations
==> af_scope_percept(citizen(af)).
==> af_scope_percept(location(af, agora)).

% Create initial af scope percepts:
af_scope_percept(P) ==> percept(top, P).

% ToM 1: Nested af models contain af scope percepts
percept(top, P), af_scope_percept(P),
bel(top, citizen(C)), { interesting_nesting(top>>C) }
==> percept(top>>C, P).

% Percept implication beliefs
==> bel(top, percept_implication(
    location(me, agora),
    states(m27, traitor(hipparchus)))).
==> bel(top, percept_implication(
    location(me, agora),
    affordance(m27, public_information))).

% Create implied percepts
percept(top, P), bel(top, percept_implication(P, Q))
==> percept(top, Q).

% ToM 2: Other citizens share percept implications
percept(top, P), bel(top, percept_implication(P, Q)),
bel(top, citizen(C)), { interesting_nesting(top>>C) }
==> bel(top>>C, percept_implication(P, Q)).

% ToM 3: Believe what you perceive
percept(top, P) ==> bel(top, P).

% ToM 4: Citizens believe they are citizens
bel(top, citizen(C)), { interesting_nesting(top>>C) }
==> bel(top>>C, citizen(me)).

% ToM 5: Believe public information on monuments
bel(top, states(Monument, S)),
percept(top, affordance(Monument, public_information))
==> bel(top, S).

% ToM 6: Agents in the agora perceive they are there
bel(top, location(C, agora)),
{ interesting_nesting(top>>C) }
==> percept(top>>C, location(me, agora)).

% ToM 7: Other citizens have the same rules as me
( Conditions ==> Conclusion ),
{ functor(Conclusion, F, 2), memberchk(F, percept, bel),
  conjunction_head(Conditions, Condition),
  ( Condition=percept(M1,_) ; Condition=bel(M1,_) ),
  depth(M1, D), D < 2 },
bel(M1, citizen(C)), { interesting_nesting(M1>>C) }
==>
{ mapsubterms(append(M1, C), Conditions, ModifiedConds),
  mapsubterms(append(M1, C), Conclusion, ModifiedConcl)
},
( ModifiedConds ==> ModifiedConcl ).
```

Example ToM rules in English

- ▶ I believe what I perceive.
- ▶ Citizens believe they are citizens
- ▶ Public information on monuments is believed.
- ▶ Other citizens have the same ToM rules as me. This is a *rule-copying* rule, i.e.:

Given rule $Conds \Rightarrow bel(M, B)$ and $bel(M, citizen(C))$, create the new rule $Conds' \Rightarrow bel(M \gg C, B)$

where $Conds'$ is $Conds$ with occurrences of M replaced with $M \gg C$.

Inferring common knowledge

- ▶ Initial knowledge base

- ▶ Agent's percepts:

- percept(\odot , citizen(*me*))*

- percept(\odot , location(*me*, *agora*))*

- percept(\odot , states(*m27*, *traitor(hipparchus)*))*

- percept(\odot , affordance(*m27*, *public_information*))*

- ▶ Scope for “any fool”:

- percept(\odot , citizen(*af*))*

- percept(\odot , location(*af*, *agora*))*

- ▶ Theory of mind rules (Pfc)

- ▶ Rules defining C1 to C3 and overall common knowledge (Pfc)

- ▶ Implementation of indication and isomorphism test (Prolog)

- ▶ Query: *ck(P)*

- ▶ Result: *P = traitor(hipparchus)*

Prolog query transcript

```
?- bel(top, traitor(hipparchus)).  
true.
```

```
?- percepts(top>>af, I) ind bel(top>>af, traitor(hipparchus)).  
I = [affordance(m27, public_information), states(m27, traitor(hipparchus))] ;  
I = [location(me, agora)] ;  
false.
```

```
?- pfc(c1([affordance(m27,public_information), states(m27,traitor(hipparchus))])  
true.
```

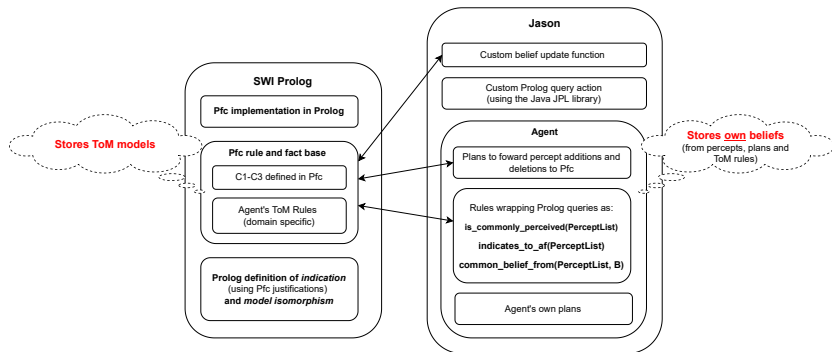
```
?- pfc(c2([affordance(m27,public_information), states(m27,traitor(hipparchus))])  
true.
```

```
?- pfc(c3([affordance(m27,public_information), states(m27,traitor(hipparchus))])  
true.
```

```
?- isomorphic_models(top>>af, top>>af>>af).  
true.
```

```
?- pfc(ck(traitor(hipparchus))).  
true.
```


Integration with Jason (work in progress)



Conclusion

- ▶ Agents can coordinate better if they understand what knowledge is common to them all.
- ▶ The logic of common knowledge has been investigated for decades, but does not appear to be practically used.
- ▶ We adapted Lewis's theory, added the missing ingredient of theory-of-mind reasoning and provided concrete semantics for indication.
- ▶ We proved that common knowledge can be inferred with only two levels of ToM reasoning.
- ▶ Our approach can be used for agents without rich logical reasoning capabilities, and can be integrated with Jason.
- ▶ Talk to me today or at my poster on Wednesday after lunch.

Questions for you

- ▶ Have you built agents that reason about common knowledge (or belief)?
- ▶ What problem domains do you have where ToM about common knowledge/belief would be useful?