

# Modeling the Mind with Logic

Selmer Bringsjord  
Rensselaer AI & Reasoning (RAIR) Lab  
Department of Cognitive Science  
Department of Computer Science  
Rensselaer Polytechnic Institute (RPI)  
Troy NY 12180 USA  
3.6.09 Arlington VA

# Modeling the Mind with Logic

# Modeling the Mind with Logic

# Modeling the Mind with Logic

# Modeling **the Mind** with Logic

# Modeling the Mind with Logic

# Modeling the Mind with **Logic**

# Modeling the Mind with Logic



# Modeling the Mind with Logic

As you must yourselves confess, the key terms here are *painfully* ambiguous.

# Modeling **the Mind** with Logic

## SUPERMINDS

People Harness Hypercomputation, and More

by

Selmer Bringsjord and Micael Zenzen

This is the first book-length presentation and defense of a new theory of human and machine cognition, according to which human persons are *superminds*. Superminds are capable of processing information not only at and below the level of Turing machines (standard computers), but above that level (the "Turing Limit"), as information processing devices that have not yet been (and perhaps can never be) built, but have been mathematically specified; these devices are known as *super-Turing machines* or hypercomputers. Superminds, as explained herein, also have properties no machine, whether above or below the Turing Limit, can have. The present book is the third and pivotal volume in Bringsjord's supermind quartet; the first two books were *What Robots Can and Can't Be* (Kluwer) and *AI and Literary Creativity* (Lawrence Erlbaum). The final chapter of this book offers eight prescriptions for the concrete practice of AI and cognitive science in light of the fact that we are superminds.

29

SELMER BRINGSJORD  
AND MICHAEL ZENZEN

SUPERMINDS  
People Harness Hypercomputation, and More

ISBN 1-4020-1095-8



9 781402 010958

KLUWER ACADEMIC PUBLISHERS

COGS 29

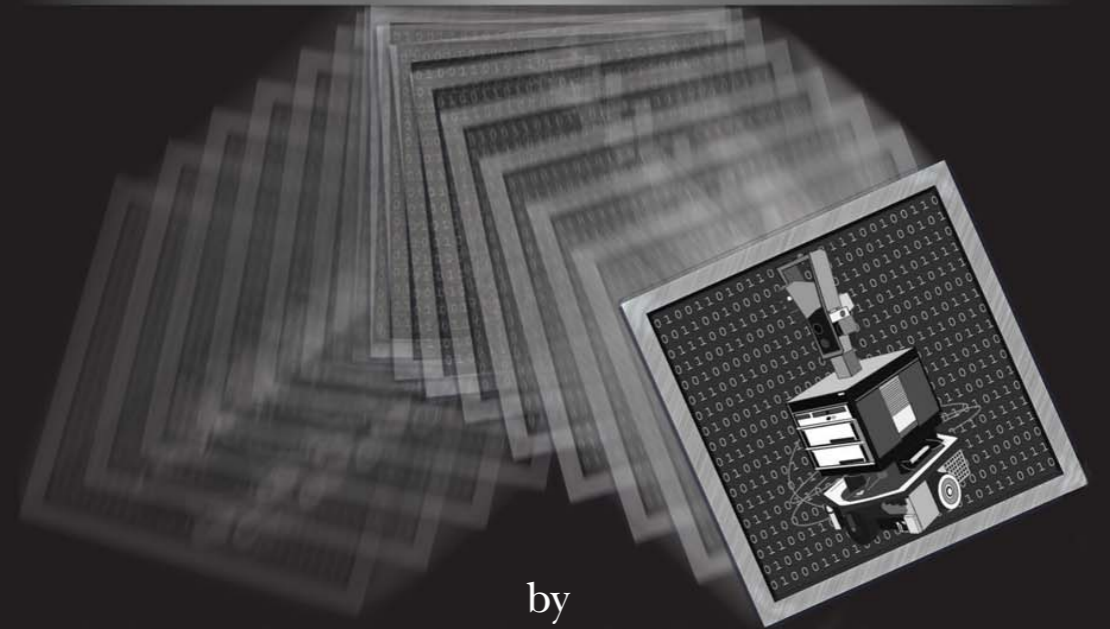


# Superminds

People Harness Hypercomputation, and More



TURING LIMIT



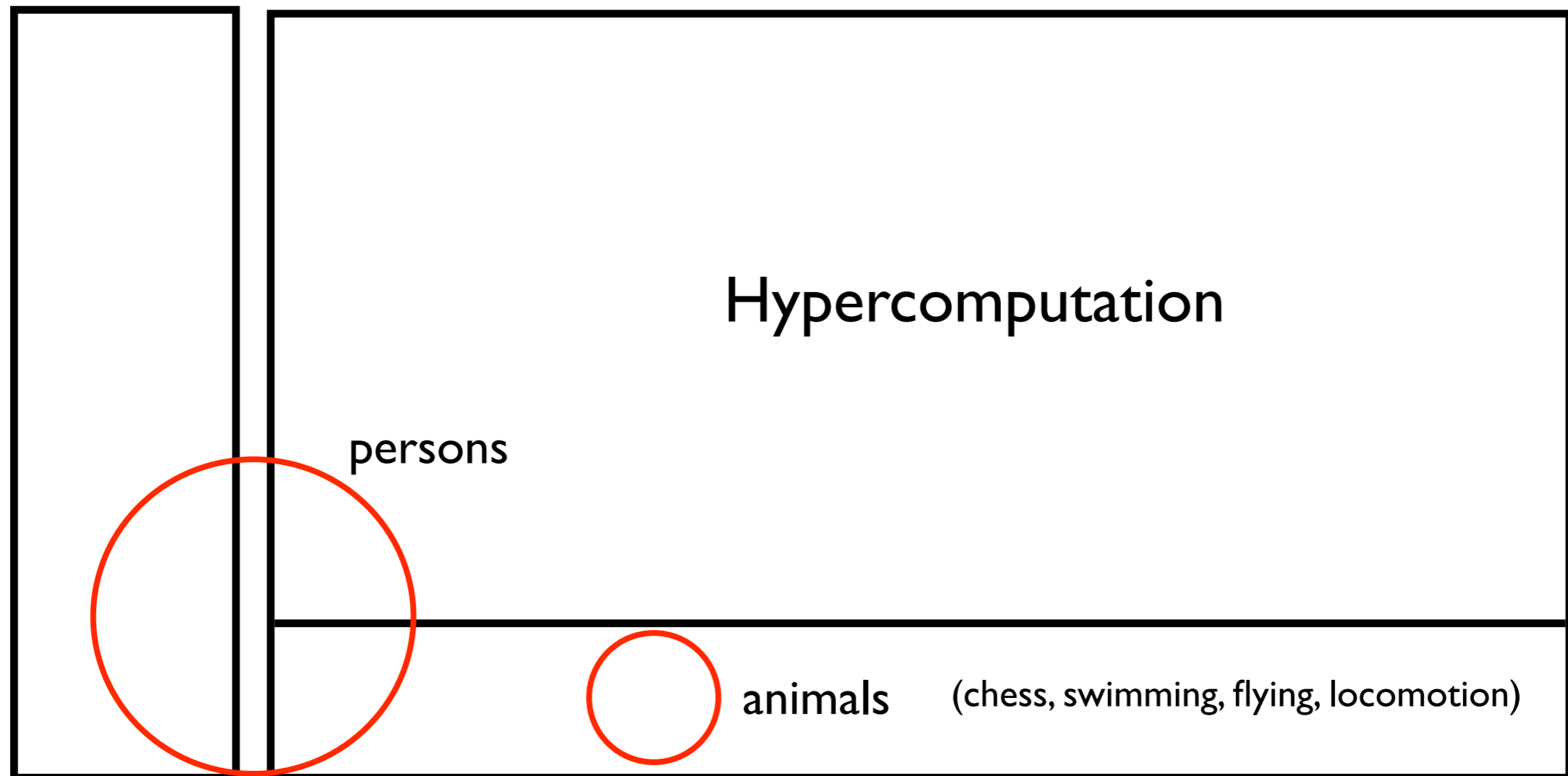
by

Selmer Bringsjord and Michael Zenzen

# Superminds (2003)

Phenomena in the incorporeal realm that can't be expressed in any third-person scheme

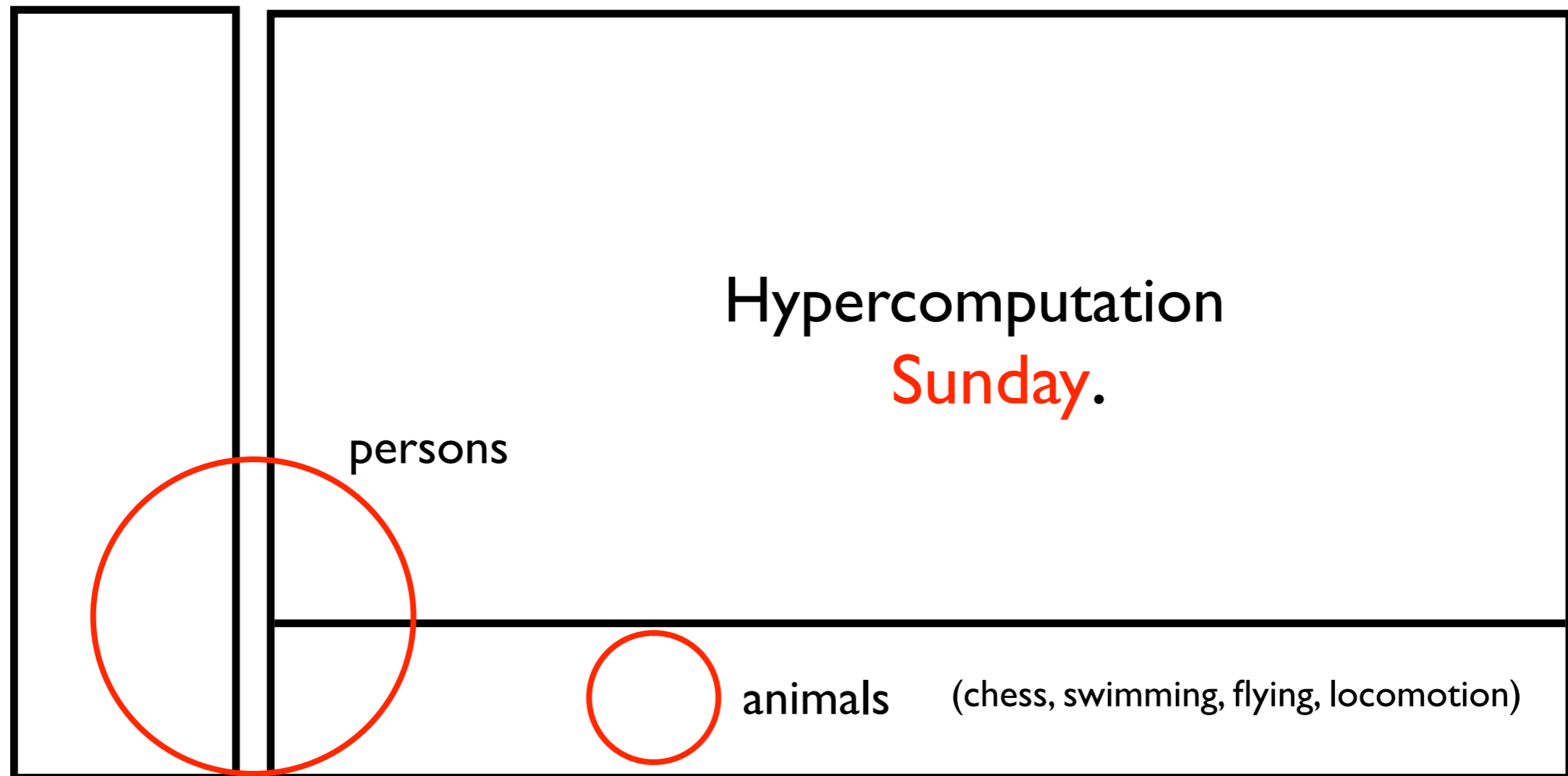
## Information Processing



# Superminds (2003)

Phenomena in the incorporeal realm that can't be expressed in any third-person scheme

## Information Processing

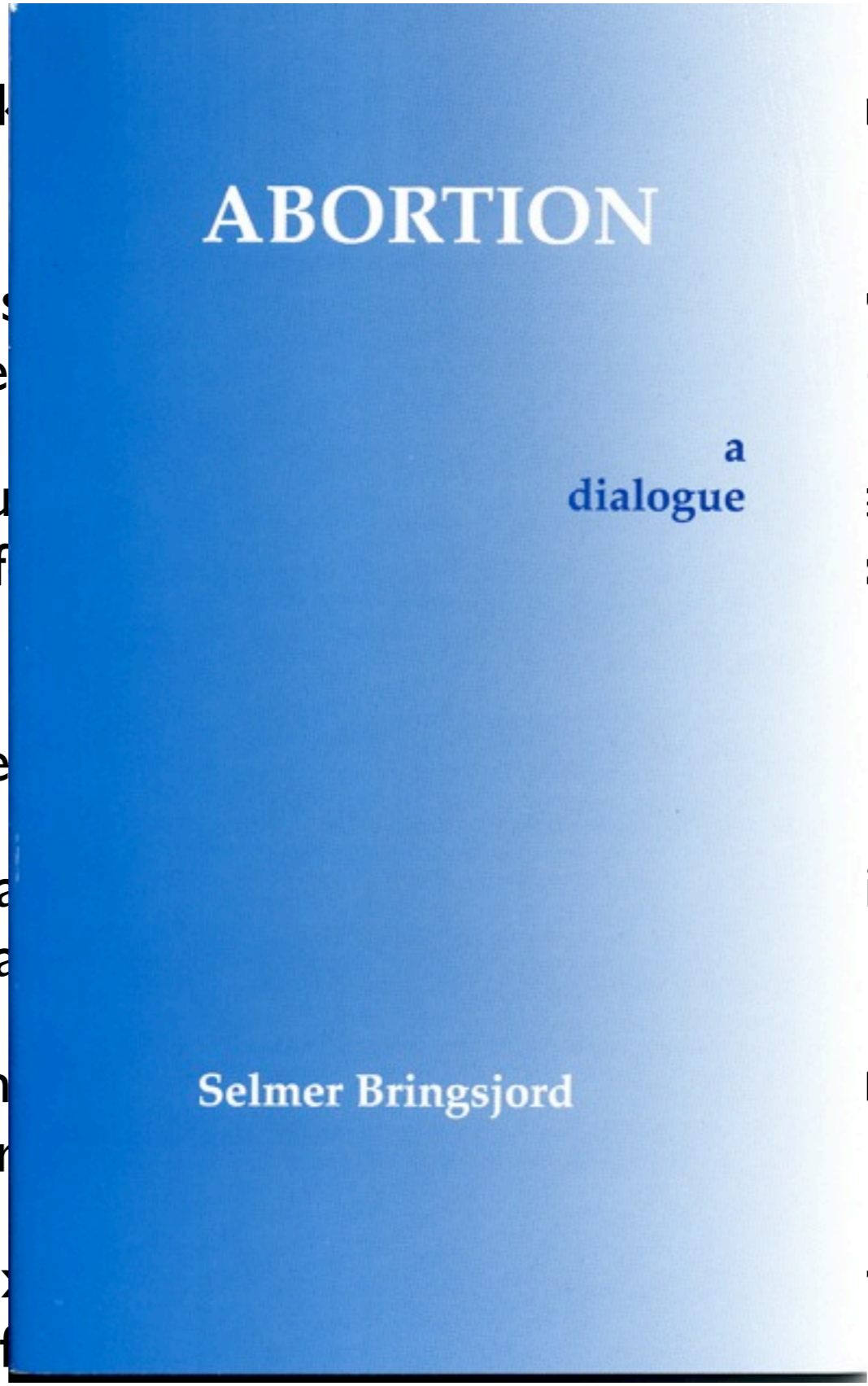


# $x$ is a person iff $x$ has the capacity ...

- to “will,” to make choices and decisions, set plans and projects — autonomously;
- for consciousness, for experiencing pain and sorrow and happiness, and a thousand other emotions — love, passion, gratitude, and so on;
- for *self*-consciousness, for being aware of his/her states of mind, inclinations, preferences, etc., and for grasping the concept of him/herself;
- to communicate through a language;
- to know things and believe things, and to believe things about what others believe (and so on);
- to desire not only particular objects and events, but also changes in his or her character;
- to reason (for example, in the fashion needed to prove the correctness of responses in false-belief, wise man, ... tests).

# x is a person iff x has the capacity ...

- to “will,” to make decisions autonomously;
- for consciousness, to be aware of a thousand other things;
- for *self-consciousness*, to be aware of her inclinations, preferences, and her self herself;
- to communicate with others;
- to know things about the world and what others believe (and to be able to communicate this knowledge);
- to desire not only to be happy but also to be a certain way or her character to be a certain way;
- to reason (for example, to be able to give reasons in favor of responses in favor of certain actions).



and projects —

and happiness, and  
desire, and so on;

states of mind,  
and her concept of him/

things about what

that also changes in his

to love the correctness

# $x$ is a person iff $x$ has the capacity ...

- to “will,” to make choices and decisions, set plans and projects — autonomously;
- for consciousness, for experiencing pain and sorrow and happiness, and a thousand other emotions — love, passion, gratitude, and so on;
- for *self*-consciousness, for being aware of his/her states of mind, inclinations, preferences, etc., and for grasping the concept of him/herself;
- to communicate through a language;
- to know things and believe things, and to believe things about what others believe (and so on);
- to desire not only particular objects and events, but also changes in his or her character;
- to reason (for example, in the fashion needed to prove the correctness of responses in false-belief, wise man, ... tests).



# $x$ is a person iff $x$ has the capacity ...

- to “will,” to make choices and decisions, set plans and projects — autonomously;

- for consciousness, for experiencing pain and sorrow and happiness, and a thousand other emotions — love, passion, gratitude, and so on;

- for *self*-consciousness, for being aware of his/her states of mind, inclinations, preferences, etc., and for grasping the concept of him/herself;

- to communicate through a language;
- to know things and believe things, and to believe things about what others believe (and so on);
- to desire not only particular objects and events, but also changes in his or her character;
- to reason (for example, in the fashion needed to prove the correctness of responses in false-belief, wise man, ... tests).

# $x$ is a person iff $x$ has the capacity ...

- to “will,” to make choices and decisions, set plans and projects — autonomously;

- for consciousness, for experiencing pain and sorrow and happiness, and a thousand other emotions — love, passion, gratitude, and so on;

unsearchably difficult; ignore *real p-*

- for self-consciousness, for being aware of his/her states of mind, inclinations, preferences, etc., and for grasping the concept of him/herself;

- to communicate through a language;
- to know things and believe things, and to believe things about what others believe (and so on);
- to desire not only particular objects and events, but also changes in his or her character;
- to reason (for example, in the fashion needed to prove the correctness of responses in false-belief, wise man, ... tests).

# $x$ is a person iff $x$ has the capacity ...

- to “will,” to make choices and decisions, set plans and projects — autonomously;

- for consciousness, for experiencing pain and sorrow and happiness, and a thousand other emotions — love, passion, gratitude, and so on;

unsearchably difficult; ignore *real p-*

- for self-consciousness, for being aware of his/her states of mind, inclinations, preferences, etc., and for grasping the concept of him/herself;

- to communicate through a language;

- to know things and believe things, and to believe things about what others believe (and so on);

- to desire not only particular objects and events, but also changes in his or her character;

- to reason (for example, in the fashion needed to prove the correctness of responses in false-belief, wise man, ... tests).

# $x$ is a person iff $x$ has the capacity ...

- to “will,” to make choices and decisions, set plans and projects — autonomously;

- for consciousness, for experiencing pain and sorrow and happiness, and a thousand other emotions — love, passion, gratitude, and so on;

**unsearchably difficult; ignore *real p-***

- for self-consciousness, for being aware of his/her states of mind, inclinations, preferences, etc., and for grasping the concept of him/herself;

**● machines still whipped by sharp toddlers; logic our only hope**

- to know things and believe things, and to believe things about what others believe (and so on);
- to desire not only particular objects and events, but also changes in his or her character;
- to reason (for example, in the fashion needed to prove the correctness of responses in false-belief, wise man, ... tests).

# $x$ is a person iff $x$ has the capacity ...

- to “will,” to make choices and decisions, set plans and projects — autonomously;

- for consciousness, for experiencing pain and sorrow and happiness, and a thousand other emotions — love, passion, gratitude, and so on;

unsearchably difficult; ignore *real p-*

- for self-consciousness, for being aware of his/her states of mind, inclinations, preferences, etc., and for grasping the concept of him/herself;

- machines still whipped by sharp toddlers; logic our only hope

- to know things and believe things, and to believe things about what others believe (and so on);

- to desire not only particular objects and events, but also changes in his or her character;

- to reason (for example, in the fashion needed to prove the correctness of responses in false-belief, wise man, ... tests).

# Modeling the Mind with **Logic**

---

THE CAMBRIDGE HANDBOOK OF  
**Computational  
Psychology**

---

EDITED BY  
**Ron Sun**

---

THE CAMBRIDGE HANDBOOK OF  
**Computational  
Psychology**

---

EDITED BY  
**Ron Sun**

**“Logic-Based/Declarative Computational Cognitive Modeling”  
by Selmer Bringsjord**

Preprint: [http://kryten.mm.rpi.edu/sb\\_lccm\\_ab-toc\\_031607.pdf](http://kryten.mm.rpi.edu/sb_lccm_ab-toc_031607.pdf)



---

THE CAMBRIDGE HANDBOOK OF

# Computational Psychology

---

Bringsjord, S. (2008) “The Logician Manifesto: At Long Last Let Logic-Based AI Become a Field Unto Itself” *Journal of Applied Logic* **6.4**: 502–525.

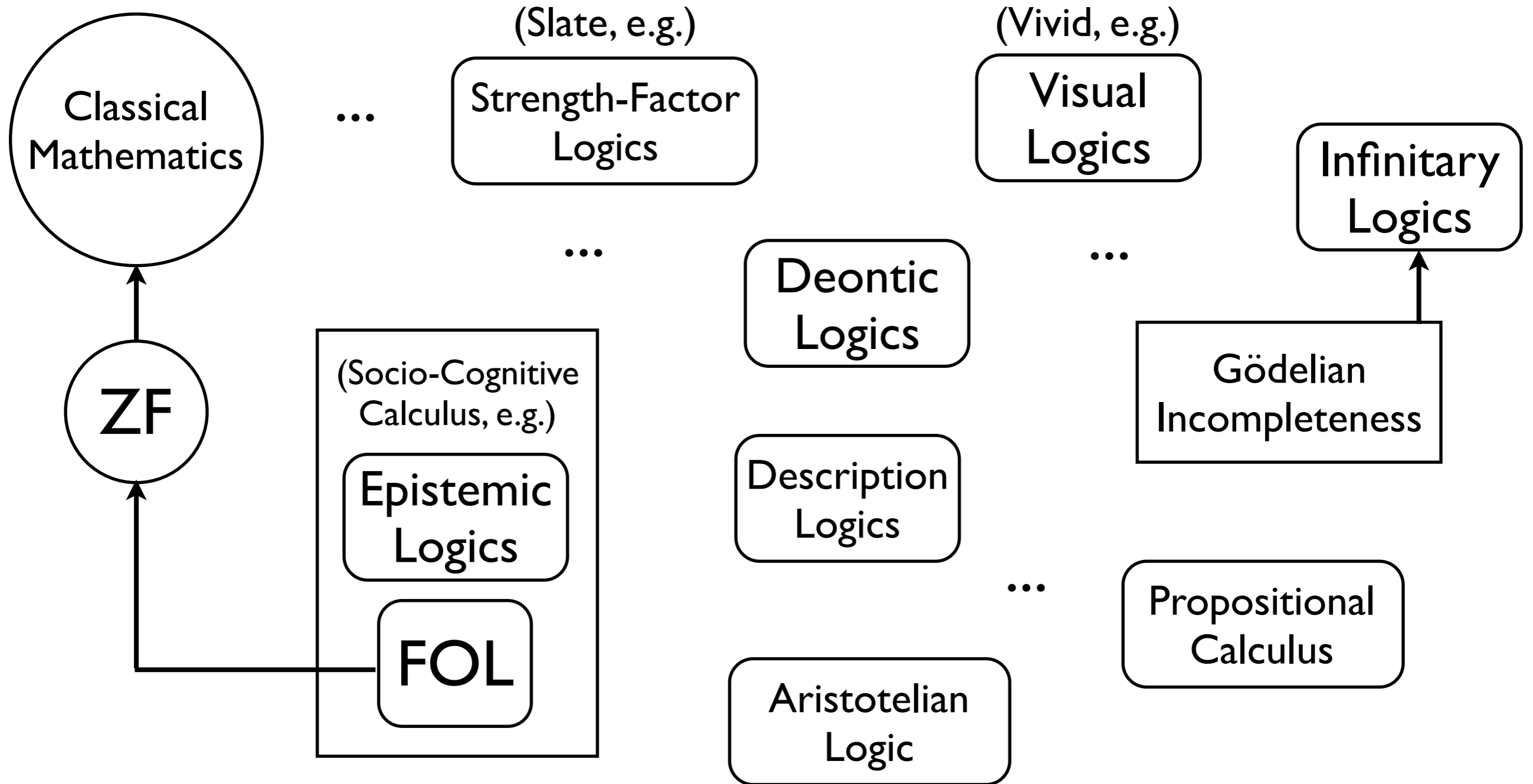
EDITED BY

Ron Sun

“Logic-Based/Declarative Computational Cognitive Modeling”  
by Selmer Bringsjord

Preprint: [http://kryten.mm.rpi.edu/sb\\_lccm\\_ab-toc\\_031607.pdf](http://kryten.mm.rpi.edu/sb_lccm_ab-toc_031607.pdf)

# The Space of Logical Systems



# Absolutely Crucial for AGI:

I'm betting the farm on one logical system  $L$   
(e.g., production systems, CYC-L, ...).

# Absolutely Crucial for AGI:

I'm betting the farm on one logical system  $L$   
(e.g., production systems, CYC-L, ...).

versus

# Absolutely Crucial for AGI:

I'm betting the farm on one logical system  $L$   
(e.g., production systems, CYC-L, ...).

versus

I know humans operate in ways that range  
*across* these logical systems, so I need a formal  
theory, and a corresponding set of processes,  
that captures the meta-coordination of various  
logical systems.

# Modeling the Mind with Logic

# Modeling the Mind with Logic

For computational cognitive science, this needs to be formalized, so that the field can be theorem-guided.

# Modeling the Mind with Logic

For computational cognitive science, this needs to be formalized, so that the field can be theorem-guided.

For AI, we can fall back on computing functions.



# Method

# Method

- Isolate and dissect the *impressive* cognition in question, whether in humans or computing machines.

# Method

- Isolate and dissect the *impressive* cognition in question, whether in humans or computing machines.
- Formalize this cognition in advanced logical systems.

# Method

- Isolate and dissect the *impressive* cognition in question, whether in humans or computing machines.
- Formalize this cognition in advanced logical systems.
- As needed, carry out further formal analysis, establishing key theorems, etc., at the meta-reasoning level.

# Method

- Isolate and dissect the *impressive* cognition in question, whether in humans or computing machines.
- Formalize this cognition in advanced logical systems.
- As needed, carry out further formal analysis, establishing key theorems, etc., at the meta-reasoning level.
- In the light of this formal work, implement working computer programs as well.

# Method

- Isolate and dissect the *impressive* cognition in question, whether in humans or computing machines.
- Formalize this cognition in advanced logical systems.
- As needed, carry out further formal analysis, establishing key theorems, etc., at the meta-reasoning level.
- In the light of this formal work, implement working computer programs as well.
- Boost performance of implementations as needed by clever software engineering and HPC.

# Method

- Isolate and dissect the *impressive* cognition in question, whether in humans or computing machines.
- Formalize this cognition in advanced logical systems.
- As needed, carry out further formal analysis, establishing key theorems, etc., at the meta-reasoning level.
- In the light of this formal work, implement working computer programs as well.
- Boost performance of implementations as needed by clever software engineering and HPC.
- Empower humans by handing over implementations.

# Method

- Isolate and dissect the *impressive* cognition in question, whether in humans or computing machines.
- Formalize this cognition in advanced logical systems.
- As needed, carry out further formal analysis, establishing key theorems, etc., at the meta-reasoning level.
- In the light of this formal work, implement working computer programs as well.
- Boost performance of implementations as needed by clever software engineering and HPC.
- Empower humans by handing over implementations.
  - If desired, provide assistance with implementations.

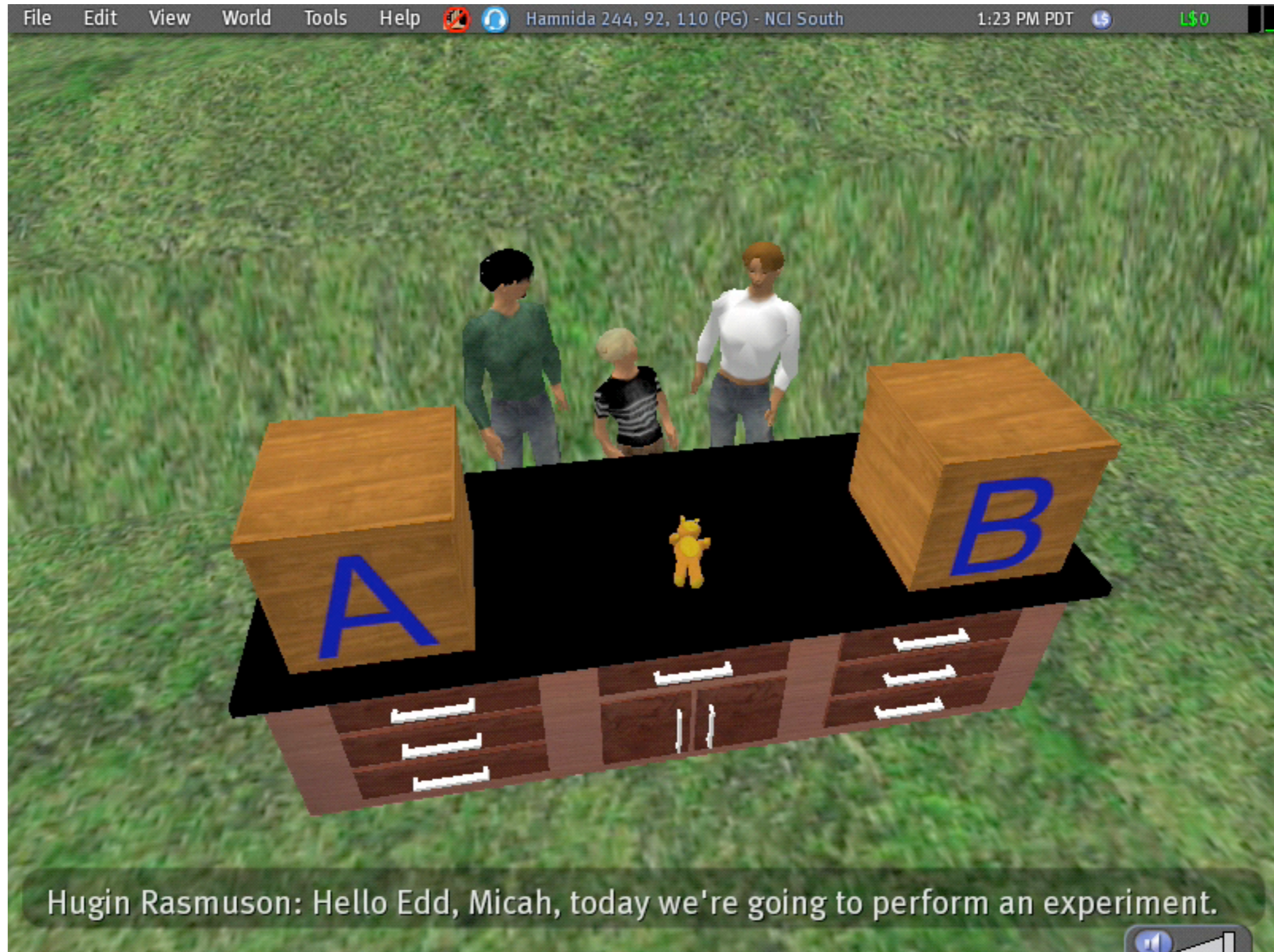


Examples ...

# False Belief(-Like) Tasks ...

In *SL*, w/ real-time comm w/ ATP

# In SL, w/ real-time comm w/ ATP



“The present account of the false belief transition is incomplete in important ways. After all, our agent had only to choose the best of two known models. This begs an understanding of the dynamics of rational revision near threshold and when the space of possible models is far larger. Further, a single formal model ought ultimately to be applicable to many false belief tasks, and to reasoning about mental states more generally. Several components seem necessary to extend a particular theory of mind into such a framework theory: a richer representation for the propositional content and attitudes in these tasks, extension of the implicit quantifier over trials to one over situations and people, and a broader view of the probability distributions relating mental state variables. Each of these is an important direction for future research.”

“Intuitive Theories of Mind: A Rational Approach to False Belief”  
Goodman et al.

“The present account of the false belief transition is incomplete in important ways. After all, our agent had only to choose the best of two known models. This begs an understanding of the dynamics of rational revision near threshold and when the space of possible models is far larger. **Further, a single formal model ought ultimately to be applicable to many false belief tasks, and to reasoning about mental states more generally.** Several components seem necessary to extend a particular theory of mind into such a framework theory: a richer representation for the propositional content and attitudes in these tasks, extension of the implicit quantifier over trials to one over situations and people, and a broader view of the probability distributions relating mental state variables. Each of these is an important direction for future research.”

“Intuitive Theories of Mind: A Rational Approach to False Belief”  
Goodman et al.

Done.

“The present account of the false belief transition is incomplete in important ways. After all, our agent had only to choose the best of two known models. This begs an understanding of the dynamics of rational revision near threshold and when the space of possible models is far larger. **Further, a single formal model ought ultimately to be applicable to many false belief tasks, and to reasoning about mental states more generally.** Several components seem necessary to extend a particular theory of mind into such a framework theory: a richer representation for the propositional content and attitudes in these tasks, extension of the implicit quantifier over trials to one over situations and people, and a broader view of the probability distributions relating mental state variables. Each of these is an important direction for future research.”

Done.

“Intuitive Theories of Mind: A Rational Approach to False Belief”  
Goodman et al.

# The Socio-Cognitive Calculus



# Toward Mechanizing Folk Psychology: A Formal Analysis of False-Belief Tasks

Konstantine Arkoudas & Selmer Bringsjord

**Abstract.** Predicting and explaining the behavior of other agents in terms of mental states is indispensable for everyday life. We believe it will be equally important for artificial agents. We present an inference system for representing and reasoning about mental states, and use it to provide a formal analysis of the false-belief task. The system allows for the representation of information about events, causation, and perceptual, doxastic, and epistemic states (vision, belief, and knowledge), incorporating ideas from the event calculus and multi-agent epistemic logic. Reasoning is performed via cognitively plausible inference rules, and a degree of automation is achieved by general-purpose inference *methods*, akin to the demons of blackboard-based multi-agent systems. The system has been implemented and is available for experimentation.

## 1 Introduction

Predicting and explaining the behavior of other people is indispensable for everyday life. The ability to ascribe mental states to others and to reason about such mental states is pervasive and invaluable. All social transactions—from engaging in commerce and negotiating to making jokes and empathizing with other people’s pain or joy—require at least a rudimentary grasp of common-sense psychology (CSP). Artificial agents without an ability of this sort would essentially suffer from autism, and would be severely handicapped in their interactions with humans. This could present problems not only for artificial agents trying to interpret human behavior, but also for artificial agents trying to interpret the behavior of one another. When a system exhibits a complex but rational behavior and detailed knowledge of its internal structure is not available, the best strategy for predicting and explaining its actions might be to analyze its behavior in intentional terms, i.e., in terms of mental states such as beliefs and desires (regardless of whether the system *actually* has genuine mental states). Mentalistic models are likely to be particularly apt for agents trying to manipulate the behavior of other agents.

Any computational treatment of CSP will have to integrate action and cognition. Agents must be able to reason about the causes and effects of various events, whether they are intentional events brought about by their own agency or non-intentional physical events. More importantly, they must be able to reason about what others believe or know about such events. To that end, our system combines ideas drawn from the event calculus and from multi-agent epistemic logics. It is based on multi-sorted first-order logic extended with subsorting, epistemic operators for perception, belief, and knowledge, and mechanisms for reasoning about causation and action. Using subsorting, we formally model agent actions as types of events, which enables us to use the resources of the event calculus to represent and reason about agent actions. The usual axioms of the event calculus are

encoded as common knowledge, suggesting that people have an understanding of the basic folk laws of causality (innate or acquired), and are indeed aware that others have such an understanding.

It is important to be clear on what we hope to accomplish with the present work. In general, any logical system or methodology capable of representing and reasoning about intentional notions such as knowledge can have at least three different uses. First, it can serve as a tool for the specification and analysis of rational epistemic agents. Second, in tandem with some appropriate reasoning mechanism, it can serve as a knowledge representation framework, i.e., it can be used *by* artificial agents to represent their own “mental states”—and those of other agents—and to deliberate and act in accordance with those states and their environment. Finally, it can be used to provide formal models of certain interesting phenomena. A chief intended contribution of our present work is of the third sort, namely, as a formal model of false-belief attributions, and in particular as a description of the competence of an agent capable of passing a false-belief task. It addresses questions such as the following: What sort of principles is it plausible to assume that an agent has to deploy in order to be able to succeed on a false-belief task? What is the depth and complexity of the required reasoning? Can such reasoning be automated, and if so, how? These questions have not been taken up in detail in the relevant discussions in cognitive science and the philosophy of mind, which have been couched in overly abstract and rather vague terms. Formal computational models such as the one we present here can help to ground such discussions, to clarify conceptual issues, and to begin to answer important questions in a concrete setting.

Although the import of such a model is primarily scientific, there can be interesting engineering implications. For instance, if the formalism is sufficiently expressive and versatile, and the posited computational mechanisms can be automated with reasonable efficiency, then the system can make potential contributions to the first two areas mentioned above. We believe that our system has such potential for two reasons. First, the combination of epistemic constructs such as common knowledge and the conceptual resources of the event calculus for dealing with causation appears to afford great expressive power, as demonstrated by our formalization. A key technical insight behind this combination is the modelling of agent actions as events via subsorting. Second, procedural abstraction mechanisms appear to hold significant promise for automation; we discuss this issue later in more detail.

The remainder of this paper is structured as follows. The next section gives the formal definition of our system. Section 3 represents the false-belief task in our system, and section 4 presents a model of the reasoning that is required to succeed in such a task, carried out in a modular fashion by collaborating methods. Section 5 discusses some related work and concludes.

## 2 A calculus for representing and reasoning about mental states

The syntactic and semantic problems that arise when one tries to use classical logic to represent and reason about intentional notions are well-known. Syntactically, modelling belief or knowledge relationally is problematic because one believes or knows arbitrarily complex propositions, whereas the arguments of relation symbols are terms built from constants, variables, and function symbols. (The objects of belief could be encoded by strings, but such representations are too low-level for engineering purposes.) Semantically, the main issue is the referential opacity (or intensionality) that must be exhibited by any operators for belief, desire, knowledge, etc. In intensional contexts one cannot freely substitute one coreferential term for another. Broadly speaking, there are two ways of addressing these issues. One is to use a modal logic, with built-in syntactic operators for intentional notions. The other is to stick with classical logic but distinguish between an object-language and a meta-language, representing intentional discourse at the object level. Each approach has its own advantages and drawbacks. Sticking with classical logic has the important advantage of efficiency, in that automated deduction systems for classical logic, such as resolution provers—which have made impressive strides over the last decade—can be used for reasoning. One disadvantage of this approach is that when the object language is first-order (includes quantification), then notions such as substitutions and alphabetic equivalence must be explicitly encoded. Depending on the facilities provided by the meta-language, this does not need to be overly onerous, but it does require extra effort.

The modal-logic approach has the advantage of solving the syntactic and semantics problems directly, without the need to distinguish an object-language and a meta-language. That is the approach we have taken in this work. The main drawback of this approach is the difficulty of automating reasoning, since standard theorem-proving techniques from classical logic cannot be directly employed. We have tried to overcome this limitation here by exploring the automation potential of methods, or derived inference rules (called *tactics* in the terminology of HOL [7]). Another drawback is the issue of semantics. The standard semantics of modal logics are given in terms of Kripke structures involving possible worlds. Such semantics are very elegant and well-understood mathematically. They are also quite intuitive for logics dealing with necessity or time. However, they are remarkably unintuitive for doxastic and epistemic logics. Not only do they fail to shed any light on the nature of belief or knowledge, but they also have a number of widely known counter-intuitive consequences that are unacceptable for resource-bounded agents, such as logical omniscience (deductive closure of knowledge, knowledge of all tautologies, etc.) and the fixed-point characterization of common knowledge. These issues are significant for us, given that we are interested in telling a plausible story for how actual agents in the real world can succeed on false-belief tasks. There have been numerous attempts to rectify these issues [8, 4, 9, 10], but each has faced serious problems of its own, and outside of Kripke structures there is no widely accepted standard at present.

Accordingly, we have not provided a possible-world semantics for our system. Note that an additional potential complication here is that the semantics of the event calculus are given in terms of circumscription, a second-order logic schema, and it is not obvious how to accommodate that feature in the setting of possible worlds. Due to these issues, and due to space restrictions, our presentation here is entirely proof-theoretic. The meanings of the various syntactic constructs—such as the knowledge operator—can be viewed as

determined by their inferential *roles*, as specified by the various inference rules. (This can itself be regarded as a form of semantics; it is called “conceptual-role semantics” or “functional semantics” in the philosophy of mind; “natural semantics” in computer science; and “procedural semantics” in cognitive science.)

The following is the formal specification of our system, describing the various sorts of our universe ( $S$ ), the signatures of certain built-in function symbols ( $f$ ), and the abstract syntax of terms ( $t$ ) and propositions ( $P$ ). The symbol  $\sqsubseteq$  denotes subsorting:

$$\begin{aligned}
 S &::= \text{Object} \mid \text{Agent} \mid \text{ActionType} \mid \text{Action} \sqsubseteq \text{Event} \\
 &\quad \mid \text{Moment} \mid \text{Boolean} \mid \text{Fluent} \\
 &\quad \text{action} : \text{Agent} \times \text{ActionType} \rightarrow \text{Action} \\
 &\quad \text{initially} : \text{Fluent} \rightarrow \text{Boolean} \\
 &\quad \text{holds} : \text{Fluent} \times \text{Moment} \rightarrow \text{Boolean} \\
 f &::= \text{happens} : \text{Event} \times \text{Moment} \rightarrow \text{Boolean} \\
 &\quad \text{clipped} : \text{Moment} \times \text{Fluent} \times \text{Moment} \rightarrow \text{Boolean} \\
 &\quad \text{initiates} : \text{Event} \times \text{Fluent} \times \text{Moment} \rightarrow \text{Boolean} \\
 &\quad \text{terminates} : \text{Event} \times \text{Fluent} \times \text{Moment} \rightarrow \text{Boolean} \\
 &\quad \text{prior} : \text{Moment} \times \text{Moment} \rightarrow \text{Boolean} \\
 t &::= x : S \mid c : S \mid f(t_1, \dots, t_n) \\
 P &::= t : \text{Boolean} \mid \neg P \mid P \wedge Q \mid P \vee Q \mid P \Rightarrow Q \mid P \Leftrightarrow Q \mid \\
 &\quad \forall x : S. P \mid \exists x : S. P \mid S(a, P) \mid K(a, P) \mid B(a, P) \mid C(P)
 \end{aligned}$$

Propositions of the form  $S(a, P)$ ,  $B(a, P)$ , and  $K(a, P)$  should be understood as saying that agent  $a$  sees that  $P$  is the case, believes that  $P$ , and knows that  $P$ , respectively. Propositions of the form  $C(P)$  assert that  $P$  is commonly known. Sort annotations will generally be omitted, as they are easily deducible from the context. We write  $P[x \mapsto t]$  for the proposition obtained from  $P$  by replacing every free occurrence of  $x$  by  $t$ , assuming that  $t$  is of a sort compatible with the sort of the free occurrences in question, and taking care to rename  $P$  as necessary to avoid variable capture. We use the infix notation  $t_1 < t_2$  instead of  $\text{prior}(t_1, t_2)$ .

We express the following standard axioms of the event calculus as common knowledge:

- $$\begin{aligned}
 [A_1] \quad & C(\forall f, t. \text{initially}(f) \wedge \neg \text{clipped}(0, f, t) \Rightarrow \text{holds}(f, t)) \\
 [A_2] \quad & C(\forall e, f, t_1, t_2. \text{happens}(e, t_1) \wedge \text{initiates}(e, f, t_1) \wedge \\
 & \quad t_1 < t_2 \wedge \neg \text{clipped}(t_1, f, t_2) \Rightarrow \text{holds}(f, t_2)) \\
 [A_3] \quad & C(\forall t_1, f, t_2. \text{clipped}(t_1, f, t_2) \Leftrightarrow \\
 & \quad [\exists e, t. \text{happens}(e, t) \wedge t_1 < t < t_2 \wedge \text{terminates}(e, f, t)])
 \end{aligned}$$

suggesting that people have a (possibly innate) understanding of basic causality principles, and are indeed aware that everybody has such an understanding. In addition to  $[A_1]$ – $[A_3]$ , we postulate a few more axioms pertaining to what people know or believe about causality. First, agents know the events that they intentionally bring about themselves—that is part of what “action” means. In fact, this is common knowledge. The following axiom expresses this:

- $$[A_4] \quad C(\forall a, d, t. \text{happens}(\text{action}(a, d), t) \Rightarrow K(a, \text{happens}(\text{action}(a, d), t)))$$

The next axiom states that it is common knowledge that if an agent  $a$  believes that a certain fluent  $f$  holds at  $t$  and he does not believe that  $f$  has been clipped between  $t$  and  $t'$ , then he will also believe that  $f$  holds at  $t'$ :

- $$[A_5] \quad C(\forall a, f, t, t'. B(a, \text{holds}(f, t)) \wedge B(a, t < t') \wedge \neg B(a, \text{clipped}(t, f, t')) \Rightarrow B(a, \text{holds}(f, t')))$$

The final axiom states that if  $a$  believes that  $b$  believes that  $f$  holds at  $t_1$  and  $a$  believes that nothing has happened between  $t_1$  and  $t_2$  to change  $b$ 's mind, then  $a$  will believe that  $b$  will not think that  $f$  has been clipped between  $t_1$  and  $t_2$ :

## 2 A calculus for representing and reasoning about mental states

The syntactic and semantic problems that arise when one tries to use classical logic to represent and reason about intentional notions are well-known. Syntactically, modelling belief or knowledge relationally is problematic because one believes or knows arbitrarily complex propositions, whereas the arguments of relation symbols are terms built from constants, variables, and function symbols. (The objects of belief could be encoded by strings, but such representations are too low-level for engineering purposes.) Semantically, the main issue is the referential opacity (or intensionality) that must be exhibited by any operators for belief, desire, knowledge, etc. In intensional contexts one cannot freely substitute one coreferential term for another. Broadly speaking, there are two ways of addressing these issues. One is to use a modal logic, with built-in syntactic operators for intentional notions. The other is to stick with classical logic but distinguish between an object-language and a meta-language, representing intentional discourse at the object level. Each approach has its own advantages and drawbacks. Sticking with classical logic has the important advantage of efficiency, in that automated deduction systems for classical logic, such as resolution provers—which have made impressive strides over the last decade—can be used for reasoning. One disadvantage of this approach is that when the object language is first-order (includes quantification), then notions such as substitutions and alphabetic equivalence must be explicitly encoded. Depending on the facilities provided by the meta-language, this does not need to be overly onerous, but it does require extra effort.

The modal-logic approach has the advantage of solving the syntactic and semantics problems directly, without the need to distinguish an object-language and a meta-language. That is the approach we have taken in this work. The main drawback of this approach is the difficulty of automating reasoning, since standard theorem-proving techniques from classical logic cannot be directly employed. We have tried to overcome this limitation here by exploring the automation potential of methods, or derived inference rules (called *tactics* in the terminology of HOL [7]). Another drawback is the issue of semantics. The standard semantics of modal logics are given in terms of Kripke structures involving possible worlds. Such semantics are very elegant and well-understood mathematically. They are also quite intuitive for logics dealing with necessity or time. However, they are remarkably unintuitive for doxastic and epistemic logics. Not only do they fail to shed any light on the nature of belief or knowledge, but they also have a number of widely known counter-intuitive consequences that are unacceptable for resource-bounded agents, such as logical omniscience (deductive closure of knowledge, knowledge of all tautologies, etc.) and the fixed-point characterization of common knowledge. These issues are significant for us, given that we are interested in telling a plausible story for how actual agents in the real world can succeed on false-belief tasks. There have been numerous attempts to rectify these issues [8, 4, 9, 10], but each has faced serious problems of its own, and outside of Kripke structures there is no widely accepted standard at present.

Accordingly, we have not provided a possible-world semantics for our system. Note that an additional potential complication here is that the semantics of the event calculus are given in terms of circumscription, a second-order logic schema, and it is not obvious how to accommodate that feature in the setting of possible worlds. Due to these issues, and due to space restrictions, our presentation here is entirely proof-theoretic. The meanings of the various syntactic constructs—such as the knowledge operator—can be viewed as

determined by their inferential *roles*, as specified by the various inference rules. (This can itself be regarded as a form of semantics; it is called “conceptual-role semantics” or “functional semantics” in the philosophy of mind; “natural semantics” in computer science; and “procedural semantics” in cognitive science.)

The following is the formal specification of our system, describing the various sorts of our universe ( $S$ ), the signatures of certain built-in function symbols ( $f$ ), and the abstract syntax of terms ( $t$ ) and propositions ( $P$ ). The symbol  $\sqsubseteq$  denotes subsorting:

```

S ::= Object | Agent | ActionType | Action  $\sqsubseteq$  Event
    | Moment | Boolean | Fluent
    action : Agent  $\times$  ActionType  $\rightarrow$  Action
    initially : Fluent  $\rightarrow$  Boolean
    holds : Fluent  $\times$  Moment  $\rightarrow$  Boolean
f ::= happens : Event  $\times$  Moment  $\rightarrow$  Boolean
    clipped : Moment  $\times$  Fluent  $\times$  Moment  $\rightarrow$  Boolean
    initiates : Event  $\times$  Fluent  $\times$  Moment  $\rightarrow$  Boolean
    terminates : Event  $\times$  Fluent  $\times$  Moment  $\rightarrow$  Boolean
    prior : Moment  $\times$  Moment  $\rightarrow$  Boolean
t ::= x : S | c : S | f(t1, ..., tn)
P ::= t : Boolean |  $\neg P$  | P  $\wedge$  Q | P  $\vee$  Q | P  $\Rightarrow$  Q | P  $\Leftrightarrow$  Q |
     $\forall x : S. P$  |  $\exists x : S. P$  | S(a, P) | K(a, P) | B(a, P) | C(P)

```

Propositions of the form  $S(a, P)$ ,  $B(a, P)$ , and  $K(a, P)$  should be understood as saying that agent  $a$  sees that  $P$  is the case, believes that  $P$ , and knows that  $P$ , respectively. Propositions of the form  $C(P)$  assert that  $P$  is commonly known. Sort annotations will generally be omitted, as they are easily deducible from the context. We write  $P[x \mapsto t]$  for the proposition obtained from  $P$  by replacing every free occurrence of  $x$  by  $t$ , assuming that  $t$  is of a sort compatible with the sort of the free occurrences in question, and taking care to rename  $P$  as necessary to avoid variable capture. We use the infix notation  $t_1 < t_2$  instead of  $prior(t_1, t_2)$ .

We express the following standard axioms of the event calculus as common knowledge:

- $$\begin{aligned}
 [A_1] \quad & C(\forall f, t. \text{initially}(f) \wedge \neg \text{clipped}(0, f, t) \Rightarrow \text{holds}(f, t)) \\
 [A_2] \quad & C(\forall e, f, t_1, t_2. \text{happens}(e, t_1) \wedge \text{initiates}(e, f, t_1) \wedge \\
 & \quad t_1 < t_2 \wedge \neg \text{clipped}(t_1, f, t_2) \Rightarrow \text{holds}(f, t_2)) \\
 [A_3] \quad & C(\forall t_1, f, t_2. \text{clipped}(t_1, f, t_2) \Leftrightarrow \\
 & \quad [\exists e, t. \text{happens}(e, t) \wedge t_1 < t < t_2 \wedge \text{terminates}(e, f, t)])
 \end{aligned}$$

suggesting that people have a (possibly innate) understanding of basic causality principles, and are indeed aware that everybody has such an understanding. In addition to  $[A_1]$ – $[A_3]$ , we postulate a few more axioms pertaining to what people know or believe about causality. First, agents know the events that they intentionally bring about themselves—that is part of what “action” means. In fact, this is common knowledge. The following axiom expresses this:

- $$[A_4] \quad C(\forall a, d, t. \text{happens}(\text{action}(a, d), t) \Rightarrow K(a, \text{happens}(\text{action}(a, d), t)))$$

The next axiom states that it is common knowledge that if an agent  $a$  believes that a certain fluent  $f$  holds at  $t$  and he does not believe that  $f$  has been clipped between  $t$  and  $t'$ , then he will also believe that  $f$  holds at  $t'$ :

- $$[A_5] \quad C(\forall a, f, t, t'. B(a, \text{holds}(f, t)) \wedge B(a, t < t') \wedge \neg B(a, \text{clipped}(t, f, t')) \Rightarrow B(a, \text{holds}(f, t')))$$

The final axiom states that if  $a$  believes that  $b$  believes that  $f$  holds at  $t_1$  and  $a$  believes that nothing has happened between  $t_1$  and  $t_2$  to change  $b$ 's mind, then  $a$  will believe that  $b$  will not think that  $f$  has been clipped between  $t_1$  and  $t_2$ :

In this approach, ontologies are simply pairs

$$(\Sigma, \Phi)$$

## 2 A calculus for representing and reasoning about mental states

The syntactic and semantic problems that arise when one tries to use classical logic to represent and reason about intentional notions are well-known. Syntactically, modelling belief or knowledge relationally is problematic because one believes or knows arbitrarily complex propositions, whereas the arguments of relation symbols are terms built from constants, variables, and function symbols. (The objects of belief could be encoded by strings, but such representations are too low-level for engineering purposes.) Semantically, the main issue is the referential opacity (or intensionality) that must be exhibited by any operators for belief, desire, knowledge, etc. In intensional contexts one cannot freely substitute one coreferential term for another. Broadly speaking, there are two ways of addressing these issues. One is to use a modal logic, with built-in syntactic operators for intentional notions. The other is to stick with classical logic but distinguish between an object-language and a meta-language, representing intentional discourse at the object level. Each approach has its own advantages and drawbacks. Sticking with classical logic has the important advantage of efficiency, in that automated deduction systems for classical logic, such as resolution provers—which have made impressive strides over the last decade—can be used for reasoning. One disadvantage of this approach is that when the object language is first-order (includes quantification), then notions such as substitutions and alphabetic equivalence must be explicitly encoded. Depending on the facilities provided by the meta-language, this does not need to be overly onerous, but it does require extra effort.

The modal-logic approach has the advantage of solving the syntactic and semantics problems directly, without the need to distinguish an object-language and a meta-language. That is the approach we have taken in this work. The main drawback of this approach is the difficulty of automating reasoning, since standard theorem-proving techniques from classical logic cannot be directly employed. We have tried to overcome this limitation here by exploring the automation potential of methods, or derived inference rules (called *tactics* in the terminology of HOL [7]). Another drawback is the issue of semantics. The standard semantics of modal logics are given in terms of Kripke structures involving possible worlds. Such semantics are very elegant and well-understood mathematically. They are also quite intuitive for logics dealing with necessity or time. However, they are remarkably unintuitive for doxastic and epistemic logics. Not only do they fail to shed any light on the nature of belief or knowledge, but they also have a number of widely known counter-intuitive consequences that are unacceptable for resource-bounded agents, such as logical omniscience (deductive closure of knowledge, knowledge of all tautologies, etc.) and the fixed-point characterization of common knowledge. These issues are significant for us, given that we are interested in telling a plausible story for how actual agents in the real world can succeed on false-belief tasks. There have been numerous attempts to rectify these issues [8, 4, 9, 10], but each has faced serious problems of its own, and outside of Kripke structures there is no widely accepted standard at present.

Accordingly, we have not provided a possible-world semantics for our system. Note that an additional potential complication here is that the semantics of the event calculus are given in terms of circumscription, a second-order logic schema, and it is not obvious how to accommodate that feature in the setting of possible worlds. Due to these issues, and due to space restrictions, our presentation here is entirely proof-theoretic. The meanings of the various syntactic constructs—such as the knowledge operator—can be viewed as

determined by their inferential *roles*, as specified by the various inference rules. (This can itself be regarded as a form of semantics; it is called “conceptual-role semantics” or “functional semantics” in the philosophy of mind; “natural semantics” in computer science; and “procedural semantics” in cognitive science.)

The following is the formal specification of our system, describing the various sorts of our universe ( $S$ ), the signatures of certain built-in function symbols ( $f$ ), and the abstract syntax of terms ( $t$ ) and propositions ( $P$ ). The symbol  $\sqsubseteq$  denotes subsorting:

```

S ::= Object | Agent | ActionType | Action  $\sqsubseteq$  Event
    | Moment | Boolean | Fluent
    action : Agent  $\times$  ActionType  $\rightarrow$  Action
    initially : Fluent  $\rightarrow$  Boolean
    holds : Fluent  $\times$  Moment  $\rightarrow$  Boolean
f ::= happens : Event  $\times$  Moment  $\rightarrow$  Boolean
    clipped : Moment  $\times$  Fluent  $\times$  Moment  $\rightarrow$  Boolean
    initiates : Event  $\times$  Fluent  $\times$  Moment  $\rightarrow$  Boolean
    terminates : Event  $\times$  Fluent  $\times$  Moment  $\rightarrow$  Boolean
    prior : Moment  $\times$  Moment  $\rightarrow$  Boolean
t ::= x : S | c : S | f(t1, ..., tn)
P ::= t : Boolean |  $\neg P$  | P  $\wedge$  Q | P  $\vee$  Q | P  $\Rightarrow$  Q | P  $\Leftrightarrow$  Q |
     $\forall x : S. P$  |  $\exists x : S. P$  | S(a, P) | K(a, P) | B(a, P) | C(P)

```

Propositions of the form  $S(a, P)$ ,  $B(a, P)$ , and  $K(a, P)$  should be understood as saying that agent  $a$  sees that  $P$  is the case, believes that  $P$ , and knows that  $P$ , respectively. Propositions of the form  $C(P)$  assert that  $P$  is commonly known. Sort annotations will generally be omitted, as they are easily deducible from the context. We write  $P[x \mapsto t]$  for the proposition obtained from  $P$  by replacing every free occurrence of  $x$  by  $t$ , assuming that  $t$  is of a sort compatible with the sort of the free occurrences in question, and taking care to rename  $P$  as necessary to avoid variable capture. We use the infix notation  $t_1 < t_2$  instead of  $prior(t_1, t_2)$ .

We express the following standard axioms of the event calculus as common knowledge:

```

[A1] C( $\forall f, t. initially(f) \wedge \neg clipped(0, f, t) \Rightarrow holds(f, t)$ )
[A2] C( $\forall e, f, t_1, t_2. happens(e, t_1) \wedge initiates(e, f, t_1) \wedge t_1 < t_2 \wedge \neg clipped(t_1, f, t_2) \Rightarrow holds(f, t_2)$ )
[A3] C( $\forall t_1, f, t_2. clipped(t_1, f, t_2) \Leftrightarrow \exists e, t. happens(e, t) \wedge t_1 < t < t_2 \wedge terminates(e, f, t)$ )

```

suggesting that people have a (possibly innate) understanding of basic causality principles, and are indeed aware that everybody has such an understanding. In addition to [A<sub>1</sub>]–[A<sub>3</sub>], we postulate a few more axioms pertaining to what people know or believe about causality. First, agents know the events that they intentionally bring about themselves—that is part of what “action” means. In fact, this is common knowledge. The following axiom expresses this:

```

[A4] C( $\forall a, d, t. happens(action(a, d), t) \Rightarrow K(a, happens(action(a, d), t))$ )

```

The next axiom states that it is common knowledge that if an agent  $a$  believes that a certain fluent  $f$  holds at  $t$  and he does not believe that  $f$  has been clipped between  $t$  and  $t'$ , then he will also believe that  $f$  holds at  $t'$ :

```

[A5] C( $\forall a, f, t, t'. B(a, holds(f, t)) \wedge B(a, t < t') \wedge \neg B(a, clipped(t, f, t')) \Rightarrow B(a, holds(f, t'))$ )

```

The final axiom states that if  $a$  believes that  $b$  believes that  $f$  holds at  $t_1$  and  $a$  believes that nothing has happened between  $t_1$  and  $t_2$  to change  $b$ 's mind, then  $a$  will believe that  $b$  will not think that  $f$  has been clipped between  $t_1$  and  $t_2$ :

In this approach, ontologies are simply pairs

$$(\Sigma, \Phi)$$

Full generality wrt time and change: includes event calculus — yet fast.

formalize this scenario in our calculus. In the next section we will present a formal explanation as to how Alice can come to acquire the correct belief about Bob's false belief.

We introduce the sort `Location` and the following function symbols specifically for reasoning about the false-belief task:

$$\begin{aligned} places &: \text{Object} \times \text{Location} \rightarrow \text{ActionType} \\ moves &: \text{Object} \times \text{Location} \times \text{Location} \rightarrow \text{ActionType} \\ located &: \text{Object} \times \text{Location} \rightarrow \text{Fluent} \end{aligned}$$

Intuitively,  $action(a, places(o, l))$  signifies  $a$ 's action of placing object  $o$  in location  $l$ , while

$$action(a, moves(o, l_1, l_2))$$

is  $a$ 's action of moving object  $o$  from location  $l_1$  to location  $l_2$ . It is common knowledge that placing  $o$  in  $l$  initiates the fluent  $located(o, l)$ :

$$[D_1] \quad C(\forall a, t, o, l . initiates(action(a, places(o, l)), located(o, l), t))$$

It is likewise known that if an object  $o$  is located at  $l_1$  at a time  $t$ , then the act of moving  $o$  from  $l_1$  to  $l_2$  results in  $o$  being located at  $l_2$ :

$$[D_2] \quad C(\forall a, t, o, l_1, l_2 . holds(located(o, l_1), t) \Rightarrow initiates(action(a, moves(o, l_1, l_2)), located(o, l_2), t))$$

If, in addition, the new location is different from the old one, the move terminates the fluent  $located(o, l_1)$ :

$$[D_3] \quad C(\forall a, t, o, l_1, l_2 . holds(located(o, l_1), t) \wedge l_1 \neq l_2 \Rightarrow terminates(action(a, moves(o, l_1, l_2)), located(o, l_1), t))$$

The following axiom captures the constraint that an object cannot be in more than one place at one time; this is also common knowledge:

$$[D_4] \quad C(\forall o, t, l_1, l_2 . holds(located(o, l_1), t) \wedge holds(located(o, l_2), t) \Rightarrow l_1 = l_2)$$

We introduce three time moments that are central to the narrative of the false-belief task: *beginning*, *departure*, and *return*. The first signifies the time point when Bob places the cookie in the cabinet, while *departure* and *return* mark the points when he leaves and comes back, respectively. We assume that it's common knowledge that these three time points are linearly ordered in the obvious manner:

$$[D_5] \quad C(beginning < departure < return).$$

We also introduce two distinct locations, *cabinet* and *drawer*:

$$[D_6] \quad C(cabinet \neq drawer).$$

Finally, we introduce a domain `Cookie` as a subsort of `Object`, and declare a single element of it, *cookie*. It is a given premise that, in the beginning, Alice sees Bob place the cookie in the cabinet:

$$[D_7] \quad S(Alice, happens(action(Bob, places(cookie, cabinet)), beginning)).$$

#### 4 Modeling the reasoning underlying false-belief tasks, and automating it via abstraction

At this point we have enough representational and reasoning machinery in place to infer the correct conclusion from a couple of obvious premises. However, a monolithic derivation of the conclusion from the premises would be unsatisfactory, as it would not give us a story about how such reasoning can be dynamically put together. Agents must be able to reason about the behavior of other agents efficiently. It is not at all obvious how efficiency can be achieved in the absence of mechanisms for abstraction, modularity, and reusability.

We can begin to address both issues by pursuing further the idea of derived inference rules, and by borrowing a page from classic work in cognitive science and production systems. Suppose that we had a mechanism which enabled the derivation of not only *schematic* inference rules, such as the ones that we presented in section 2, but derived inference rules allowing for arbitrary computation and search. We could then formulate *generic* inference rules, capable of being applied to an unbounded (potentially infinite) number of arbitrarily complex concrete situations.

Our system has a notion of *method* that allows for that type of abstraction and encapsulation. Methods are derived inference rules, not just of the schematic kind, but incorporating arbitrary computation and search. They are thus more general than the simple if-then rules of production systems, and more akin to the knowledge sources (or "demons") of blackboard systems [5]. They can be viewed as encapsulating specialized expertise in deriving certain types of conclusions from certain given information. They can be parameterized over any variables, e.g., arbitrary agents or time points. In our system, the role of the blackboard is played by an associative data structure (shared by all methods) known as the *assumption base*, which is an efficiently indexed collection of propositions that represent the collective knowledge state at any given moment, including perceptual knowledge. The assumption base is capable of serving as a communication buffer for the various methods. Finally, the control executive is itself a method, which directs the reasoning process incrementally by invoking various methods triggered by the contents of the assumption base.

We describe below three general-purpose methods for reasoning in the calculus we have presented. With these methods, the reasoning for the false-belief task can be performed in a handful of lines—essentially with one invocation of each of these methods. We stress that these methods are not ad hoc or hardwired to false-belief tasks. They are generic, and can be reused in any context that requires reasoning about other minds and satisfies the relevant preconditions. In particular, the methods do not contain or require any information specific to false-belief tasks.

- *Method 1*: This method, which we call  $M_1$ , shows that when an agent  $a_1$  sees an agent  $a_2$  perform some action-type  $\alpha$  at some time point  $t$ ,  $a_1$  knows that  $a_2$  knows that  $a_2$  has carried out  $\alpha$  at  $t$ .  $M_1$  is parameterized over  $a_1$ ,  $a_2$ ,  $\alpha$ , and  $t$ :

1. The starting premise is that  $a_1$  sees  $a_2$  perform  $\alpha$  at  $t$ :

$$S(a_1, happens(action(a_2, \alpha), t)) \quad (1)$$

2. Therefore,  $a_1$  knows that the corresponding event has occurred at  $t$ :

$$K(a_1, happens(action(a_2, \alpha), t)) \quad (2)$$

This follows from the preceding premise and  $[DR_4]$ .

3. From  $[A_4]$  and  $[DR_2]$  we obtain:

$$K(a_1, \forall a, \alpha, t . happens(action(a, \alpha), t) \Rightarrow K(a, happens(action(a, \alpha), t))) \quad (3)$$

4. From (3) and  $[DR_9]$  we get:

$$K(a_1, happens(action(a_2, \alpha), t) \Rightarrow K(a_2, happens(action(a_2, \alpha), t))) \quad (4)$$

5. From (4), (2), and  $[DR_6]$  we get:

$$K(a_1, K(a_2, happens(action(a_2, \alpha), t))) \quad (5)$$

- *Method 2*: The second method,  $M_2$ , shows that when (1) it is common knowledge that a certain event  $e$  initiates a fluent  $f$ ; (2) an agent  $a_1$  knows that an agent  $a_2$  knows that  $e$  has happened at a

Methods would seem to be key for general intelligence.

formalize this scenario in our calculus. In the next section we will present a formal explanation as to how Alice can come to acquire the correct belief about Bob's false belief.

We introduce the sort `Location` and the following function symbols specifically for reasoning about the false-belief task:

$$\begin{aligned} places &: \text{Object} \times \text{Location} \rightarrow \text{ActionType} \\ moves &: \text{Object} \times \text{Location} \times \text{Location} \rightarrow \text{ActionType} \\ located &: \text{Object} \times \text{Location} \rightarrow \text{Fluent} \end{aligned}$$

Intuitively,  $action(a, places(o, l))$  signifies  $a$ 's action of placing object  $o$  in location  $l$ , while

$$action(a, moves(o, l_1, l_2))$$

is  $a$ 's action of moving object  $o$  from location  $l_1$  to location  $l_2$ . It is common knowledge that placing  $o$  in  $l$  initiates the fluent  $located(o, l)$ :

$$[D_1] \quad C(\forall a, t, o, l . initiates(action(a, places(o, l)), located(o, l), t))$$

It is likewise known that if an object  $o$  is located at  $l_1$  at a time  $t$ , then the act of moving  $o$  from  $l_1$  to  $l_2$  results in  $o$  being located at  $l_2$ :

$$[D_2] \quad C(\forall a, t, o, l_1, l_2 . holds(located(o, l_1), t) \Rightarrow initiates(action(a, moves(o, l_1, l_2)), located(o, l_2), t))$$

If, in addition, the new location is different from the old one, the move terminates the fluent  $located(o, l_1)$ :

$$[D_3] \quad C(\forall a, t, o, l_1, l_2 . holds(located(o, l_1), t) \wedge l_1 \neq l_2 \Rightarrow terminates(action(a, moves(o, l_1, l_2)), located(o, l_1), t))$$

The following axiom captures the constraint that an object cannot be in more than one place at one time; this is also common knowledge:

$$[D_4] \quad C(\forall o, t, l_1, l_2 . holds(located(o, l_1), t) \wedge holds(located(o, l_2), t) \Rightarrow l_1 = l_2)$$

We introduce three time moments that are central to the narrative of the false-belief task: *beginning*, *departure*, and *return*. The first signifies the time point when Bob places the cookie in the cabinet, while *departure* and *return* mark the points when he leaves and comes back, respectively. We assume that it's common knowledge that these three time points are linearly ordered in the obvious manner:

$$[D_5] \quad C(beginning < departure < return).$$

We also introduce two distinct locations, *cabinet* and *drawer*:

$$[D_6] \quad C(cabinet \neq drawer).$$

Finally, we introduce a domain `Cookie` as a subsort of `Object`, and declare a single element of it, *cookie*. It is a given premise that, in the beginning, Alice sees Bob place the cookie in the cabinet:

$$[D_7] \quad S(Alice, happens(action(Bob, places(cookie, cabinet)), beginning)).$$

#### 4 Modeling the reasoning underlying false-belief tasks, and automating it via abstraction

At this point we have enough representational and reasoning machinery in place to infer the correct conclusion from a couple of obvious premises. However, a monolithic derivation of the conclusion from the premises would be unsatisfactory, as it would not give us a story about how such reasoning can be dynamically put together. Agents must be able to reason about the behavior of other agents efficiently. It is not at all obvious how efficiency can be achieved in the absence of mechanisms for abstraction, modularity, and reusability.

We can begin to address both issues by pursuing further the idea of derived inference rules, and by borrowing a page from classic work in cognitive science and production systems. Suppose that we had a mechanism which enabled the derivation of not only *schematic* inference rules, such as the ones that we presented in section 2, but derived inference rules allowing for arbitrary computation and search. We could then formulate *generic* inference rules, capable of being applied to an unbounded (potentially infinite) number of arbitrarily complex concrete situations.

Our system has a notion of *method* that allows for that type of abstraction and encapsulation. Methods are derived inference rules, not just of the schematic kind, but incorporating arbitrary computation and search. They are thus more general than the simple if-then rules of production systems, and more akin to the knowledge sources (or "demons") of blackboard systems [5]. They can be viewed as encapsulating specialized expertise in deriving certain types of conclusions from certain given information. They can be parameterized over any variables, e.g., arbitrary agents or time points. In our system, the role of the blackboard is played by an associative data structure (shared by all methods) known as the *assumption base*, which is an efficiently indexed collection of propositions that represent the collective knowledge state at any given moment, including perceptual knowledge. The assumption base is capable of serving as a communication buffer for the various methods. Finally, the control executive is itself a method, which directs the reasoning process incrementally by invoking various methods triggered by the contents of the assumption base.

We describe below three general-purpose methods for reasoning in the calculus we have presented. With these methods, the reasoning for the false-belief task can be performed in a handful of lines—essentially with one invocation of each of these methods. We stress that these methods are not ad hoc or hardwired to false-belief tasks. They are generic, and can be reused in any context that requires reasoning about other minds and satisfies the relevant preconditions. In particular, the methods do not contain or require any information specific to false-belief tasks.

- *Method 1*: This method, which we call  $M_1$ , shows that when an agent  $a_1$  sees an agent  $a_2$  perform some action-type  $\alpha$  at some time point  $t$ ,  $a_1$  knows that  $a_2$  knows that  $a_2$  has carried out  $\alpha$  at  $t$ .  $M_1$  is parameterized over  $a_1$ ,  $a_2$ ,  $\alpha$ , and  $t$ :

1. The starting premise is that  $a_1$  sees  $a_2$  perform  $\alpha$  at  $t$ .

$$S(a_1, happens(action(a_2, \alpha), t)) \quad (1)$$

2. Therefore,  $a_1$  knows that the corresponding event has occurred at  $t$ :

$$K(a_1, happens(action(a_2, \alpha), t)) \quad (2)$$

This follows from the preceding premise and  $[DR_4]$ .

3. From  $[A_4]$  and  $[DR_2]$  we obtain:

$$K(a_1, \forall a, \alpha, t . happens(action(a, \alpha), t) \Rightarrow K(a, happens(action(a, \alpha), t))) \quad (3)$$

4. From (3) and  $[DR_9]$  we get:

$$K(a_1, happens(action(a_2, \alpha), t) \Rightarrow K(a_2, happens(action(a_2, \alpha), t))) \quad (4)$$

5. From (4), (2), and  $[DR_6]$  we get:

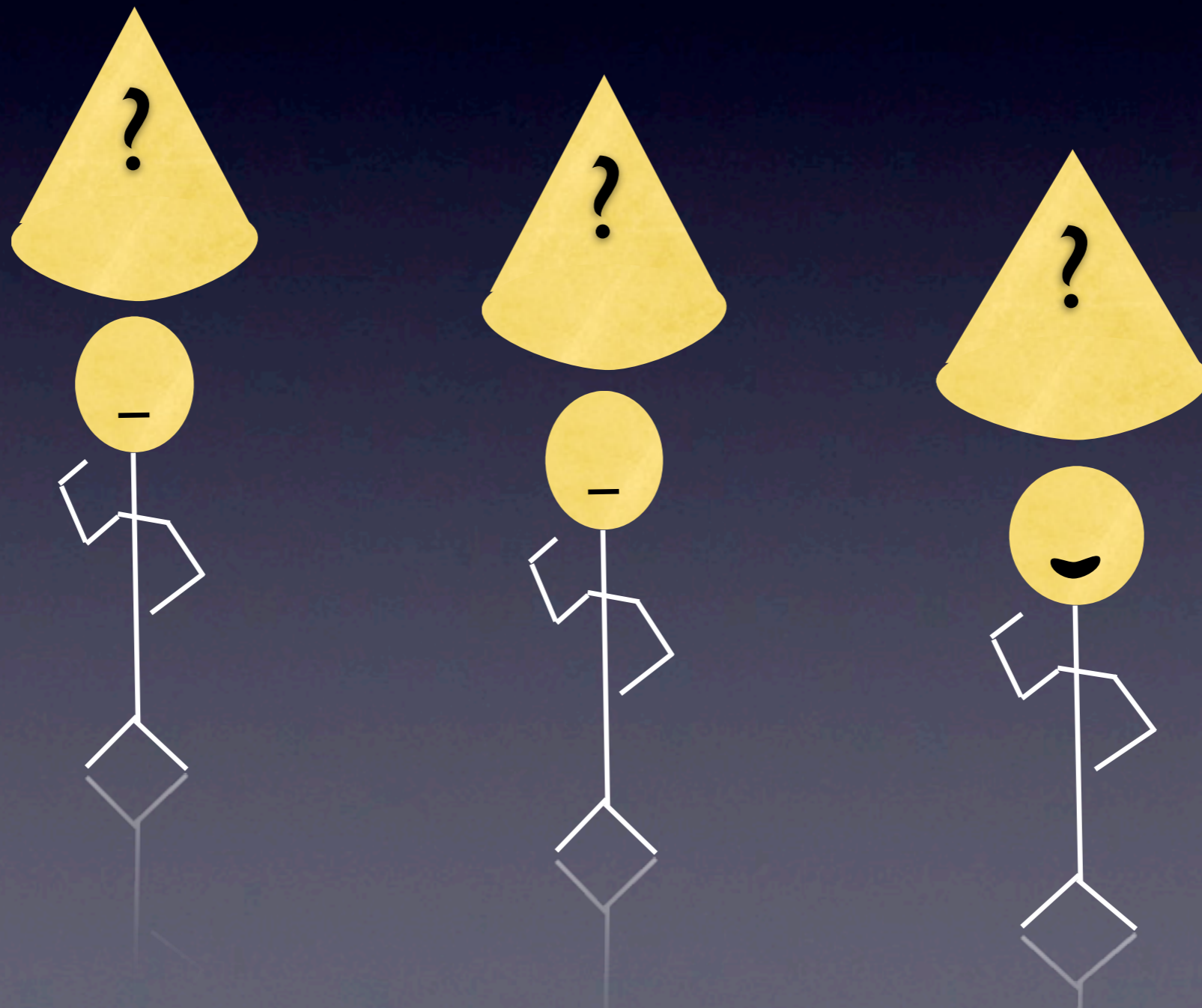
$$K(a_1, K(a_2, happens(action(a_2, \alpha), t))) \quad (5)$$

- *Method 2*: The second method,  $M_2$ , shows that when (1) it is common knowledge that a certain event  $e$  initiates a fluent  $f$ ; (2) an agent  $a_1$  knows that an agent  $a_2$  knows that  $e$  has happened at a

Methods would seem to be key for general intelligence.

# Cracking Wise Man Tests ...

# Wise Men Puzzle



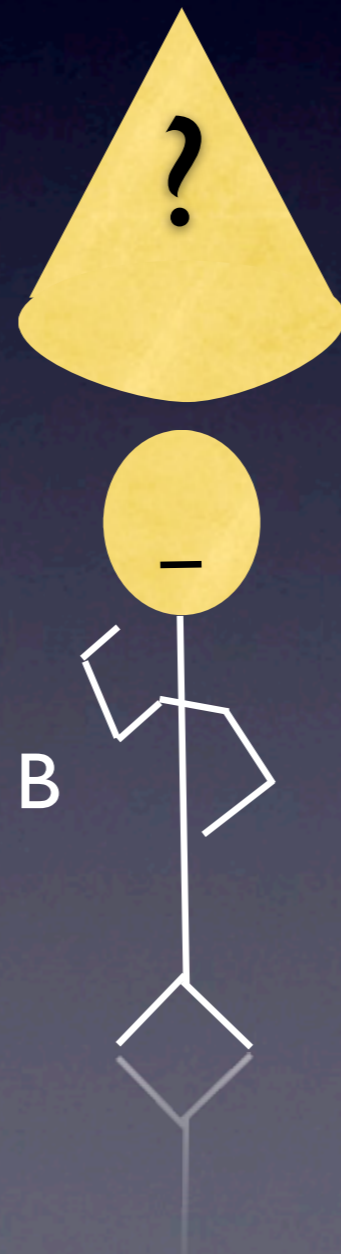


# Wise Men Puzzle

Wise man A



Wise man B



Wise man C



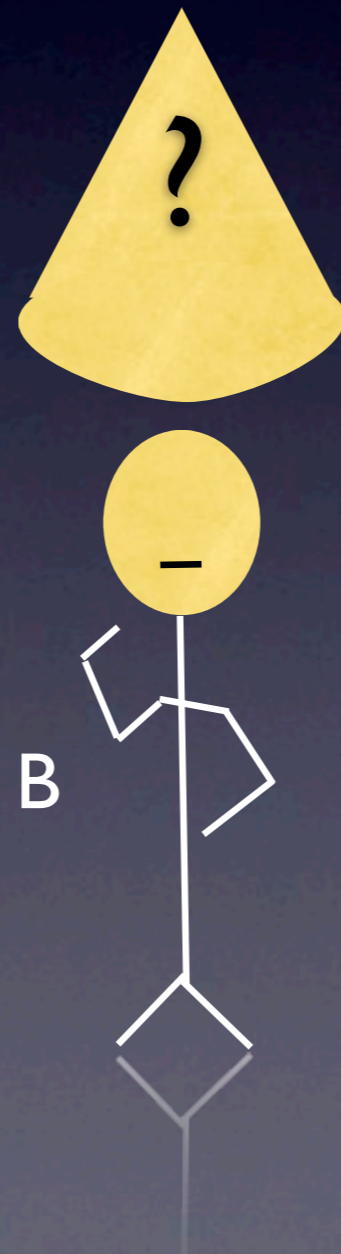
# Wise Men Puzzle

I don't know

Wise man A



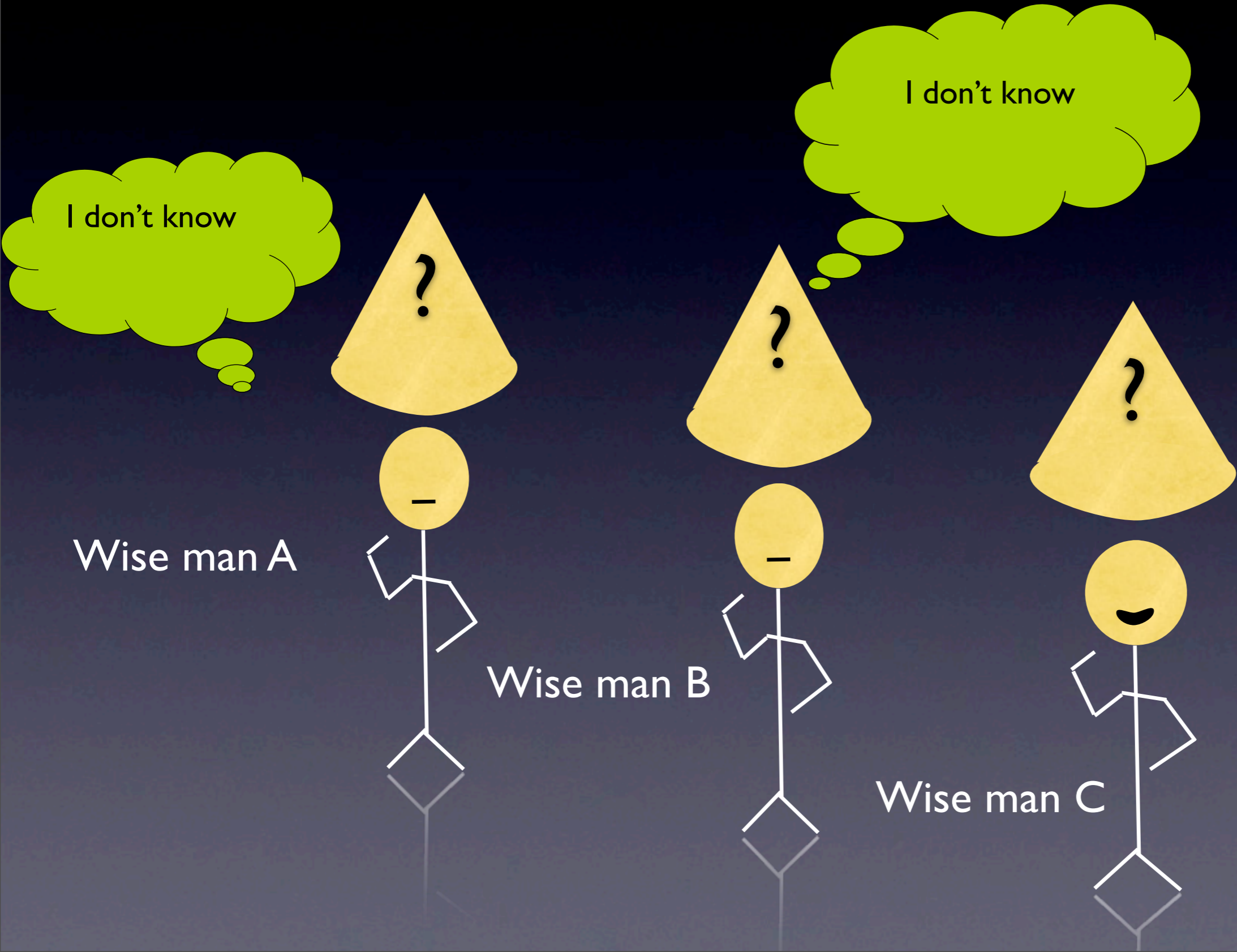
Wise man B



Wise man C



# Wise Men Puzzle

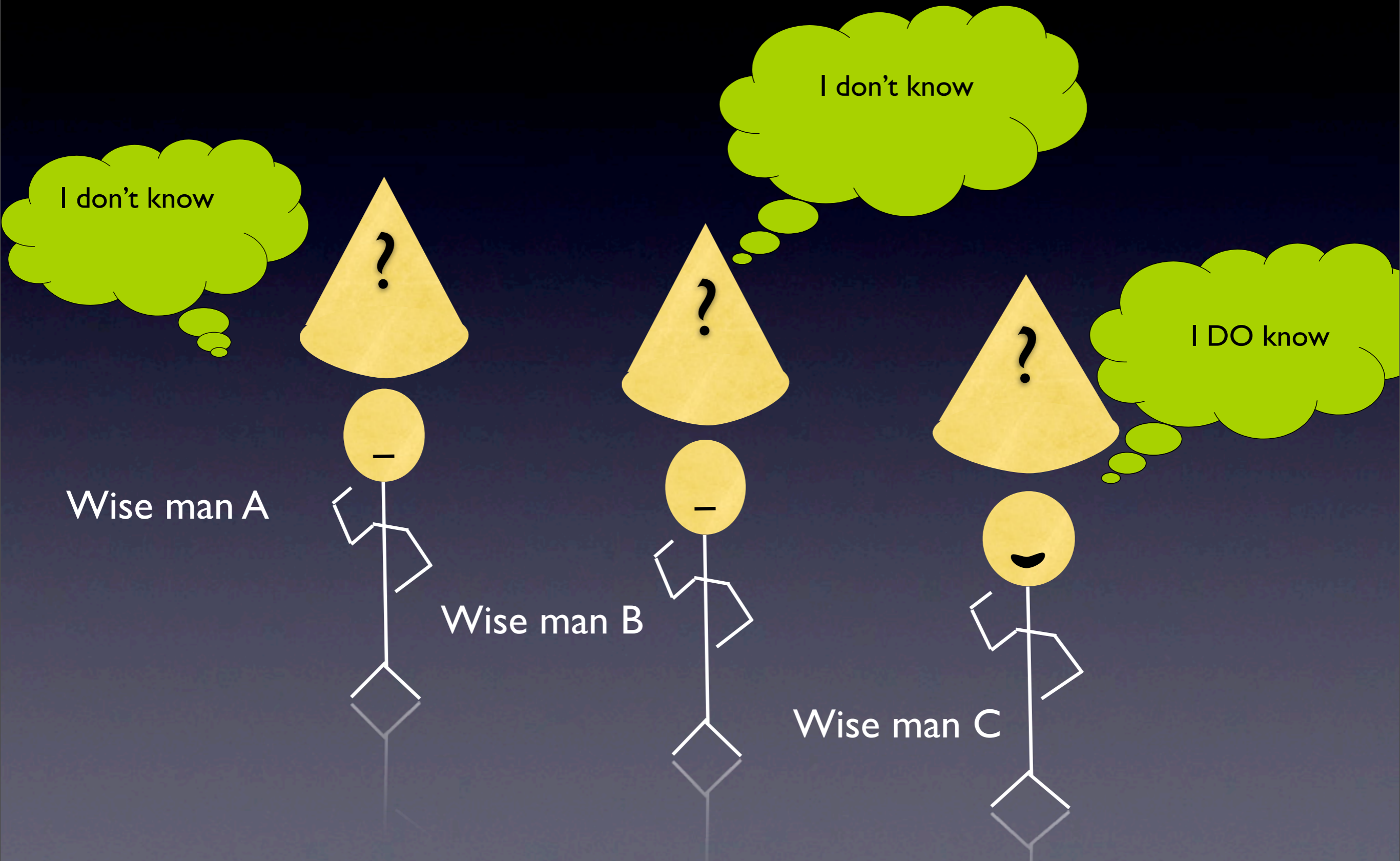


Wise man A

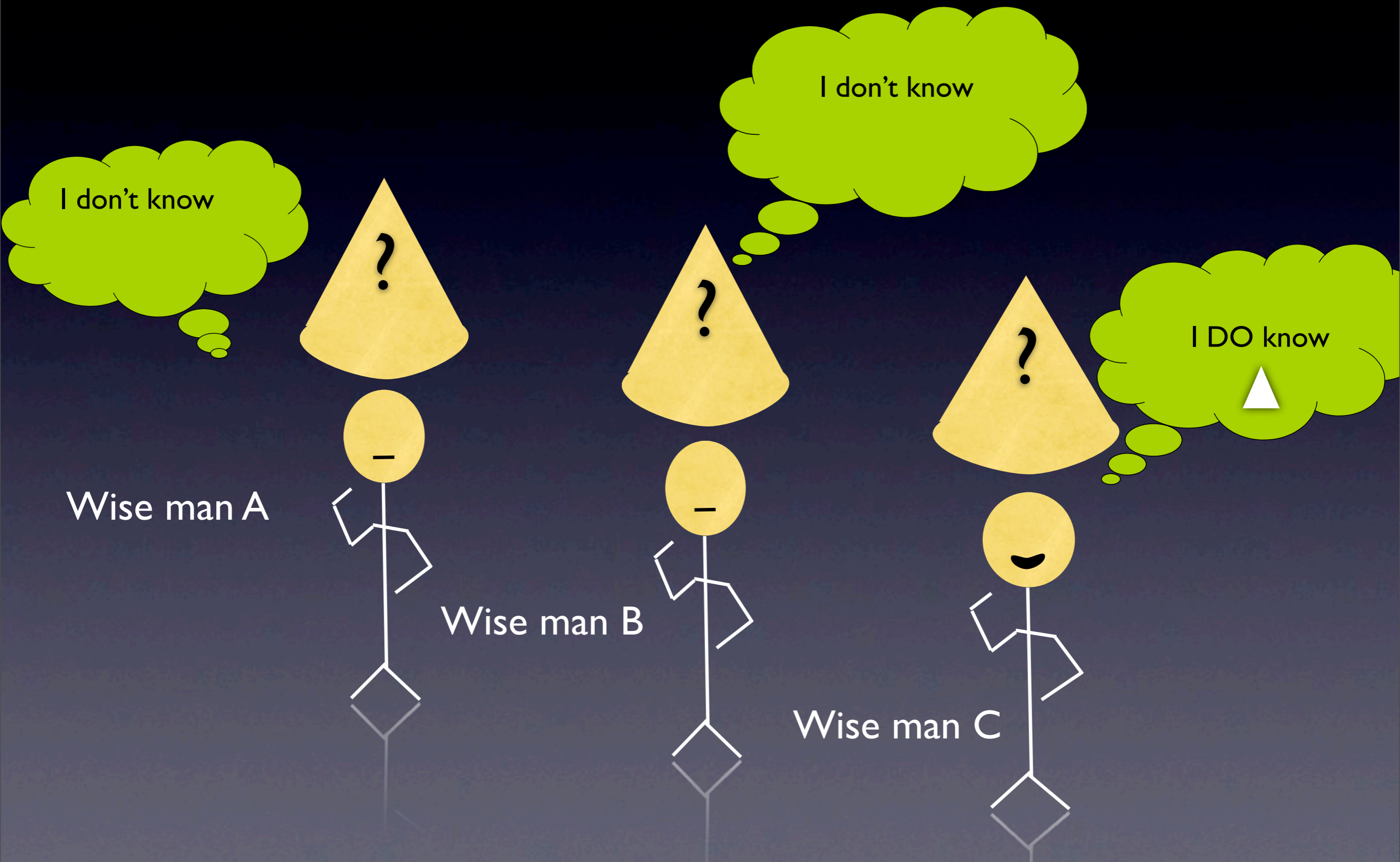
Wise man B

Wise man C

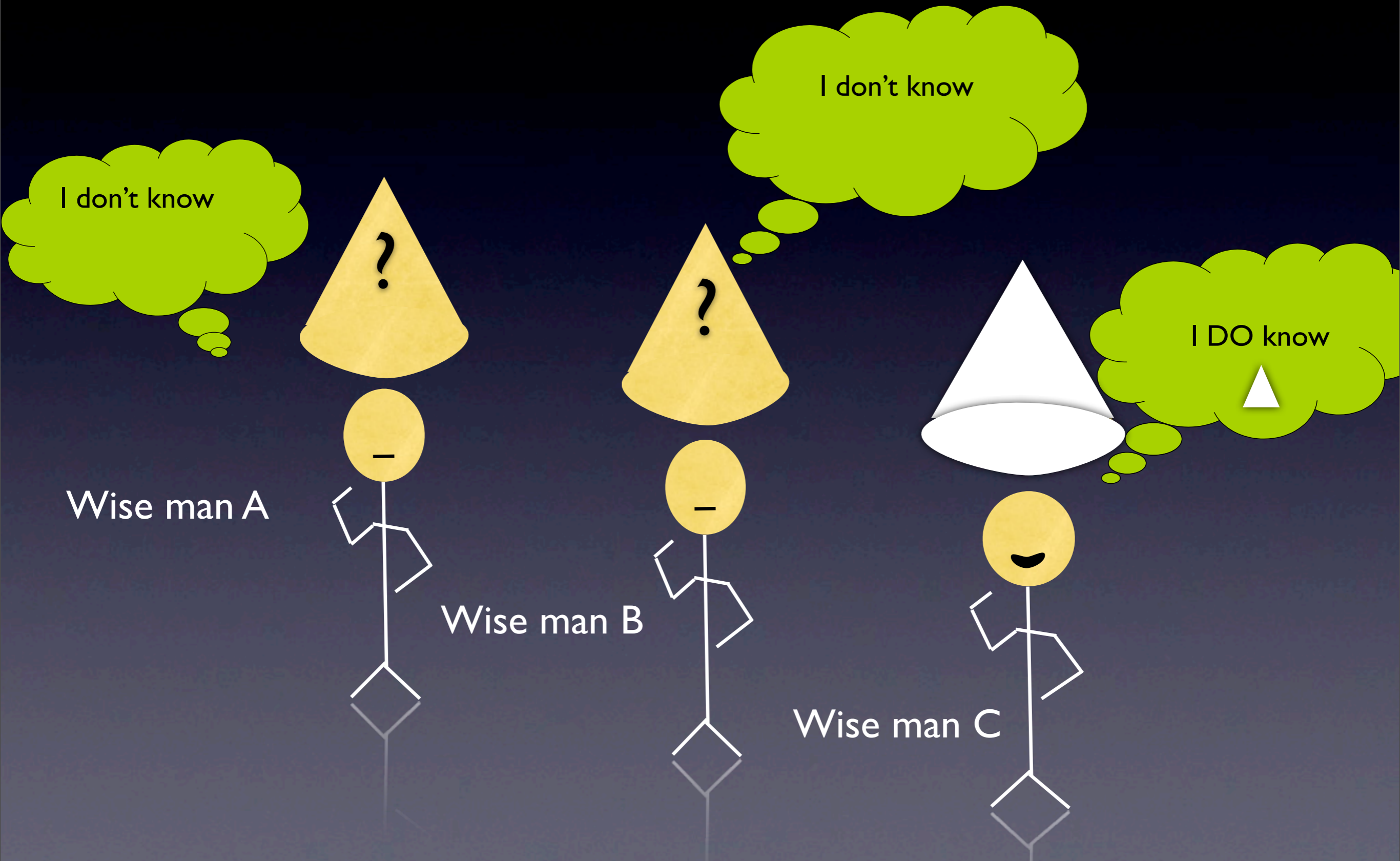
# Wise Men Puzzle



# Wise Men Puzzle



# Wise Men Puzzle



Wise man A

Wise man B

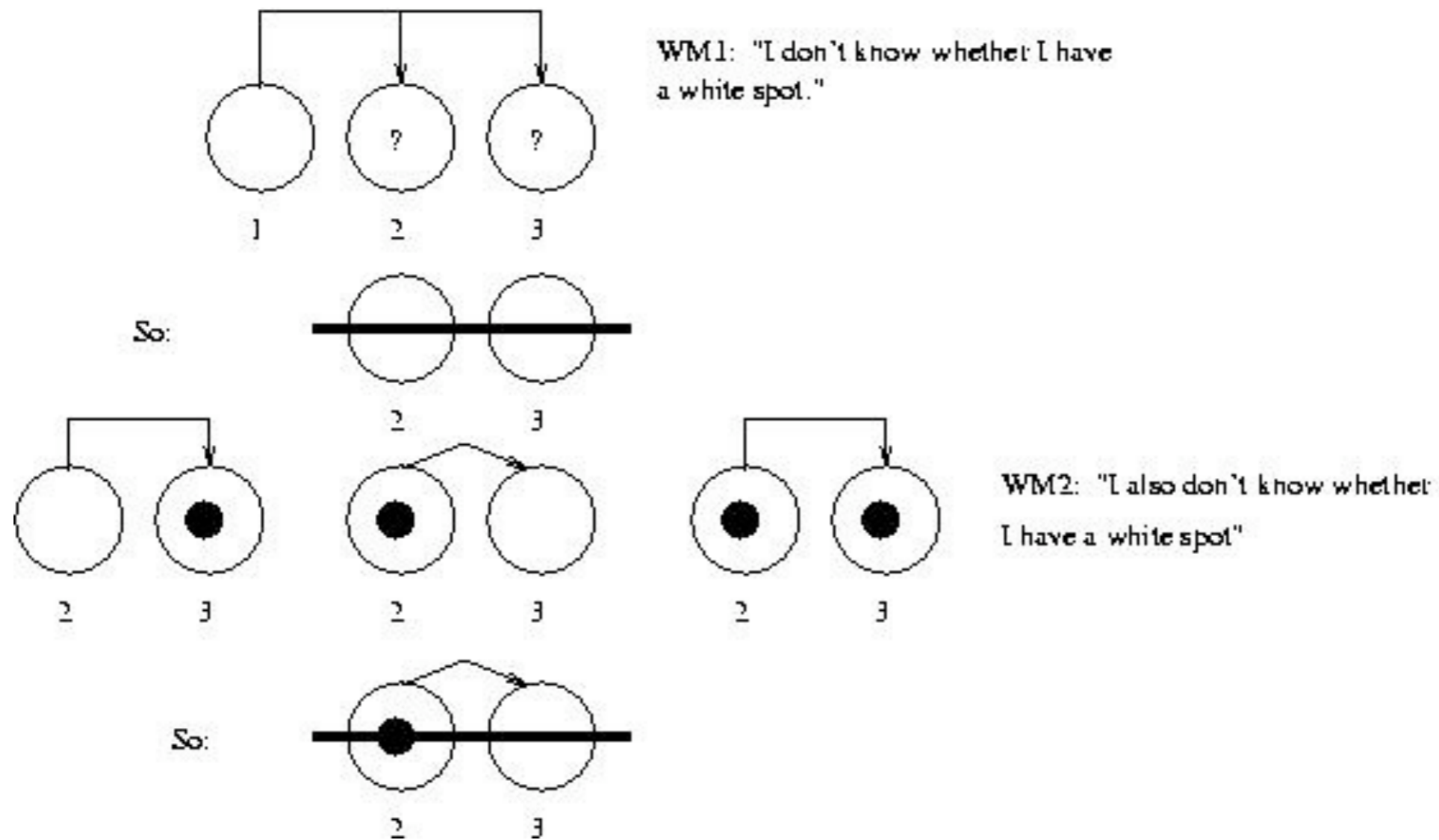
Wise man C

# Start of Reasoning in WMP<sub>3</sub>

(pov of *truly* wise man; easy for smart humans)

# Start of Reasoning in WMP<sub>3</sub>

(pov of *truly* wise man; easy for smart humans)





**Abstract.** We present an encoding of a sequent calculus for a multi-agent epistemic logic in Athena, an interactive theorem proving system for many-sorted first-order logic. We then use Athena as a metalanguage in order to reason about the multi-agent logic as an object language. This facilitates theorem proving in the multi-agent logic in several ways. First, it lets us marshal the highly efficient theorem provers for classical first-order logic that are integrated with Athena for the purpose of doing proofs in the multi-agent logic. Second, unlike model-theoretic embeddings of modal logics into classical first-order logic, our proofs are directly convertible into native epistemic logic proofs. Third, because we are able to quantify over propositions and agents, we get much of the generality and power of higher-order logic even though we are in a first-order setting. Finally, we are able to use Athena's versatile tactics for proof automation in the multi-agent logic. We illustrate by developing a tactic for solving the generalized version of the wise men problem.

## 1 Introduction

Multi-agent modal logics are widely used in Computer Science and AI. Multi-agent epistemic logics, in particular, have found applications in fields ranging from AI domains such as robotics, planning, and motivation analysis in natural language [13]; to negotiation and game theory in economics; to distributed systems analysis and protocol authentication in computer security [16,31]. The reason is simple—intelligent agents must be able to reason about knowledge. It is therefore important to have efficient means for performing machine reasoning in such logics. While the validity problem for most propositional modal logics is of intractable theoretical complexity<sup>1</sup>, several approaches have been investigated in recent years that have resulted in systems that appear to work well in practice. These approaches include tableau-based provers, SAT-based algorithms, and translations to first-order logic coupled with the use of resolution-based automated theorem provers (ATPs). Some representative systems are FaCT [24], KSATC [14], TA [25], LWB [23], and MSPASS [37].

Translation-based approaches (such as that of MSPASS) have the advantage of leveraging the tremendous implementation progress that has occurred over

<sup>1</sup> For instance, the validity problem for multi-agent propositional epistemic logic is PSPACE-complete [18]; adding a common knowledge operator makes the problem EXPTIME-complete [21].

# Proved-Sound Algorithm for Generating Proof-Theoretic Solution to WMP<sub>n</sub>

$$\begin{array}{c}
 \frac{}{\Gamma \vdash [K_\alpha(P \Rightarrow Q)] \Rightarrow [K_\alpha(P) \Rightarrow K_\alpha(Q)]} [K] \quad \frac{}{\Gamma \vdash K_\alpha(P) \Rightarrow P} [\Gamma] \\
 \\
 \frac{\emptyset \vdash P}{} [C-I] \quad \frac{}{\Gamma \vdash C(P) \Rightarrow K_\alpha(P)} [C-E] \\
 \\
 \frac{}{\Gamma \vdash [C(P \Rightarrow Q)] \Rightarrow [C(P) \Rightarrow C(Q)]} [CK] \quad \frac{}{\Gamma \vdash C(P) \Rightarrow C(K_\alpha(P))} [R]
 \end{array}$$

Fig. 2. Inference rules for the epistemic operators.

is  $\Gamma \vdash P$ . Intuitively, this is a judgment stating that  $P$  follows from  $\Gamma$ . We will write  $P, \Gamma$  (or  $\Gamma, P$ ) as an abbreviation for  $\Gamma \cup \{P\}$ . The sequent calculus that we will use consists of a collection of inference rules for deriving judgments of the form  $\Gamma \vdash P$ . Figure 1 shows the inference rules that deal with the standard propositional connectives. This part is standard (e.g., it is very similar to the sequent calculus of Ebbinghaus et al. [15]). In addition, we have some rules pertaining to  $K_\alpha$  and  $C$ , shown in Figure 2.

Rule  $[K]$  is the sequent formulation of the well-known *Kripke axiom* stating that the knowledge operator distributes over conditionals. Rule  $[CK]$  is the corresponding principle for the common knowledge operator. Rule  $[T]$  is the “truth axiom”: an agent cannot know false propositions. Rule  $[C-I]$  is an introduction rule for common knowledge: if a proposition  $P$  follows from the empty set of hypotheses, i.e., if it is a tautology, then it is commonly known. This is the common-knowledge version of the “omniscience axiom” for single-agent knowledge which says that  $\Gamma \vdash K_\alpha(P)$  can be derived from  $\emptyset \vdash P$ . We do not need to postulate that axiom in our formulation, since it follows from  $[C-I]$  and  $[C-E]$ . The latter says that if it is common knowledge that  $P$  then any (every) agent knows  $P$ , while  $[R]$  says that if it is common knowledge that  $P$  then it is common knowledge that (any) agent  $\alpha$  knows it.  $[R]$  is a reiteration rule that allows us to capture the recursive behavior of  $C$ , which is usually expressed via the so-called “induction axiom”

$$C(P \Rightarrow E(P)) \Rightarrow [P \Rightarrow C(P)]$$

where  $E$  is the shared-knowledge operator. Since we do not need  $E$  for our purposes, we omit its formalization and “unfold”  $C$  via rule  $[R]$  instead.

We state a few lemmas that will come handy later:

**Lemma 1 (Cut).** *If  $\Gamma_1 \vdash P_1$  and  $\Gamma_2, P_1 \vdash P_2$  then  $\Gamma_1 \cup \Gamma_2 \vdash P_2$ .*

**Proof:** Assume  $\Gamma_1 \vdash P_1$  and  $\Gamma_2, P_1 \vdash P_2$ . Then, by  $[\Rightarrow-I]$ , we get  $\Gamma_2 \vdash P_1 \Rightarrow P_2$ . Further, by dilation, we have  $\Gamma_1 \cup \Gamma_2 \vdash P_1 \Rightarrow P_2$  and  $\Gamma_1 \cup \Gamma_2 \vdash P_1$ . Hence, by  $[\Rightarrow-E]$ , we obtain  $\Gamma_1 \cup \Gamma_2 \vdash P_2$ .  $\square$

The proofs of the remaining lemmas are equally simple exercises.

$$\begin{array}{ll}
 \varnothing \vdash R_1 \wedge R_2 \wedge R_3 \vdash R_1 & [Reflex], \wedge-E_1 \\
 \varnothing \vdash R_1 \wedge R_2 \wedge R_3 \vdash R_2 & [Reflex], \wedge-E_1, \wedge-E_2 \\
 \varnothing \vdash R_1 \wedge R_2 \wedge R_3 \vdash R_3 & [Reflex], \wedge-E_2 \\
 \varnothing \vdash R_1 \wedge R_2 \wedge R_3 \vdash K_\alpha(\neg Q) \Rightarrow K_\alpha(P) & 2, [K], \Rightarrow-E \\
 \varnothing \vdash R_1 \wedge R_2 \wedge R_3 \vdash \neg Q \Rightarrow K_\alpha(P) & 3, 4, Lemma 2 \\
 \varnothing \vdash R_1 \wedge R_2 \wedge R_3 \vdash \neg K_\alpha(P) \Rightarrow \neg \neg Q & 5, Lemma 3 \\
 \varnothing \vdash R_1 \wedge R_2 \wedge R_3 \vdash \neg \neg Q & 6, 1, \Rightarrow-E \\
 \varnothing \vdash R_1 \wedge R_2 \wedge R_3 \vdash Q & 7, [\neg-E]
 \end{array}$$

$\square$

at the above proof is not entirely low-level because most steps combine more inference rule applications in the interest of brevity.

**Lemma 7.** *Consider any agent  $\alpha$  and propositions  $P, Q$ . Define  $R_1$  and  $R_3$  as in Lemma 6, let  $R_2 = P \vee Q$ , and let  $S_i = C(R_i)$  for  $i = 1, 2, 3$ . Then  $S_3 \vdash C(Q)$ .*

Let  $R'_2 = \neg Q \Rightarrow P$  and consider the following derivation:

$$\begin{array}{ll}
 S_1, S_2, S_3 \vdash S_1 & [Reflex] \\
 S_1, S_2, S_3 \vdash S_2 & [Reflex] \\
 S_1, S_2, S_3 \vdash S_3 & [Reflex] \\
 \vdash (P \vee Q) \Rightarrow (\neg Q \Rightarrow P) & Lemma 4a \\
 S_1, S_2, S_3 \vdash C((P \vee Q) \Rightarrow (\neg Q \Rightarrow P)) & 4, [C-I] \\
 S_1, S_2, S_3 \vdash C(P \vee Q) \Rightarrow C(\neg Q \Rightarrow P) & 5, [CK], [\Rightarrow-E] \\
 S_1, S_2, S_3 \vdash C(\neg Q \Rightarrow P) & 6, 2, [\Rightarrow-E] \\
 S_1, S_2, S_3 \vdash C(\neg Q \Rightarrow P) \Rightarrow C(K_\alpha(\neg Q \Rightarrow P)) & [R] \\
 S_1, S_2, S_3 \vdash C(K_\alpha(\neg Q \Rightarrow P)) & 8, 7, [\Rightarrow-E] \\
 R_1 \wedge K_\alpha(\neg Q \Rightarrow P) \wedge R_3 \vdash Q & Lemma 6 \\
 \vdash (R_1 \wedge K_\alpha(\neg Q \Rightarrow P) \wedge R_3) \Rightarrow Q & 10, [\Rightarrow-I] \\
 S_1, S_2, S_3 \vdash C((R_1 \wedge K_\alpha(\neg Q \Rightarrow P) \wedge R_3) \Rightarrow Q) & 11, [C-I] \\
 S_1, S_2, S_3 \vdash C(R_1 \wedge K_\alpha(\neg Q \Rightarrow P) \wedge R_3) \Rightarrow C(Q) & 12, [CK], [\Rightarrow-E] \\
 S_1, S_2, S_3 \vdash C(R_1 \wedge K_\alpha(\neg Q \Rightarrow P) \wedge R_3) & 1, 3, 9, Lemma 5, [\wedge-I] \\
 S_1, S_2, S_3 \vdash C(Q) & 13, 14, [\Rightarrow-E]
 \end{array}$$

$\square$

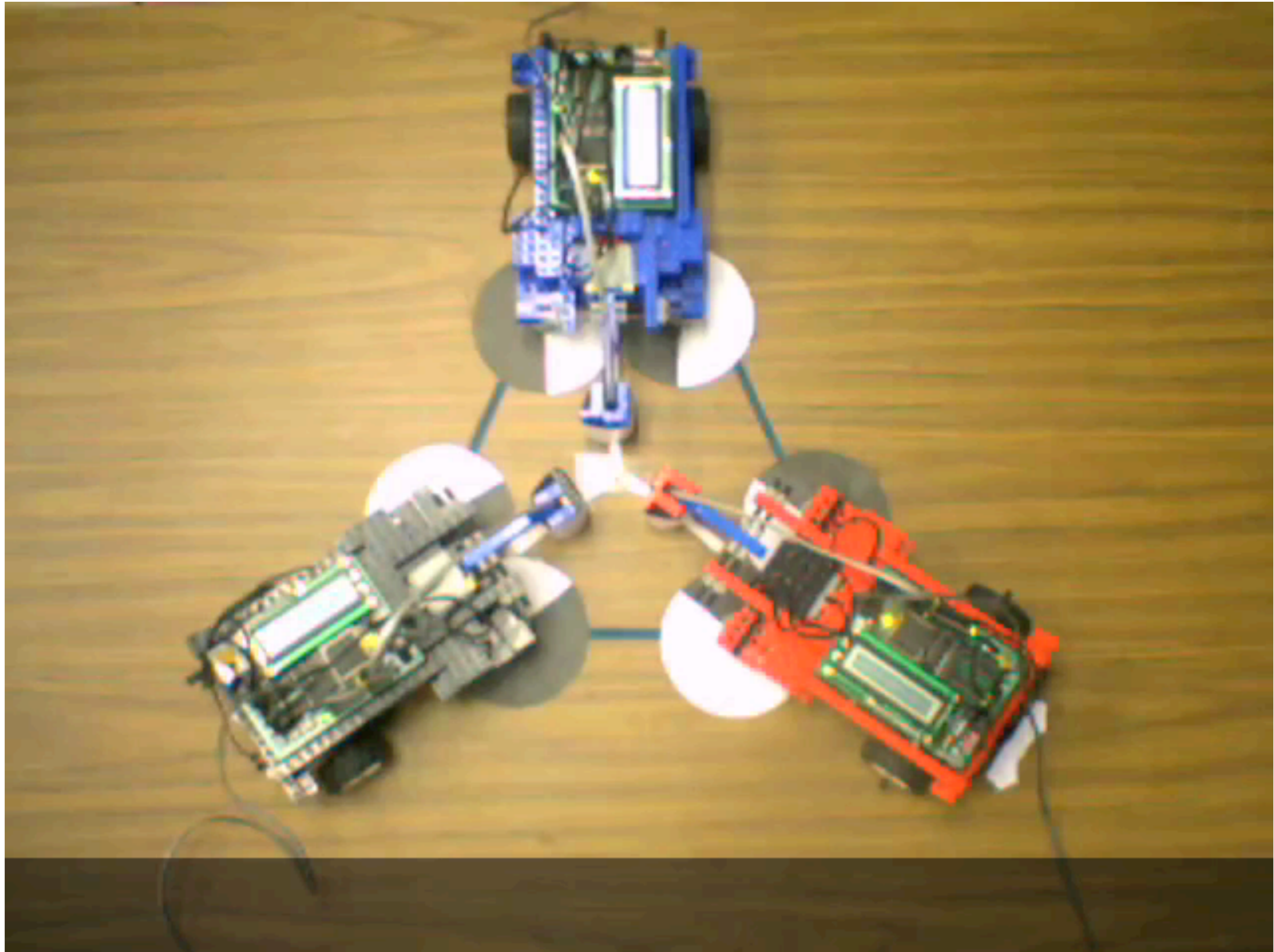
all  $n \geq 1$ , it turns out that the last— $(n + 1)^{st}$ —wise man knows he is marked. The case of two wise men is simple. The reasoning runs essentially by induction. The second wise man reasons as follows:

Suppose I were not marked. Then  $w_1$  would have seen this, and knowing that at least one of us is marked, he would have inferred that he was marked one. But  $w_1$  has expressed ignorance; therefore, I must be marked.

Now the case of  $n = 3$  wise men  $w_1, w_2, w_3$ . After  $w_1$  announces that he does not know that he is marked,  $w_2$  and  $w_3$  both infer that at least one of them is marked. For if neither  $w_2$  nor  $w_3$  were marked,  $w_1$  would have seen this and would have concluded—and stated—that he was the marked one, since he knows that at least one of the three is marked. At this point the puzzle reduces to the two-men case: both  $w_2$  and  $w_3$  know that at least one of them is marked,

All our  
human-  
authored  
proofs  
machine-  
checked.

# “Life and Death” Wise Man Test (3)



# Modeling *Visual* Reasoning

Arkoudas, K. & Bringsjord, S. (forthcoming)  
“Vivid: An AI Framework for Heterogeneous  
Problem Solving” *Artificial Intelligence*.

(Thank you DARPA and IARPA/ARDA/DTO.)

# Toward a General Logician Methodology for Engineering Ethically Correct Robots

Selmer Bringsjord, Konstantine Arkoudas, and Paul Bello,  
Rensselaer Polytechnic Institute

**A**s intelligent machines assume an increasingly prominent role in our lives, there seems little doubt they will eventually be called on to make important, ethically charged decisions. For example, we expect hospitals to deploy robots that can administer medications, carry out tests, perform surgery, and so on, supported by software agents,

or softbots, that will manage related data. (Our discussion of ethical robots extends to all artificial agents, embodied or not.) Consider also that robots are already finding their way to the battlefield, where many of their potential actions could inflict harm that is ethically impermissible.

How can we ensure that such robots will always behave in an ethically correct manner? How can we know ahead of time, via rationales expressed in clear natural languages, that their behavior will be constrained specifically by the ethical codes affirmed by human overseers? Pessimists have claimed that the answer to these questions is: "We can't!" For example, Sun Microsystems' cofounder and former chief scientist, Bill Joy, published a highly influential argument for this answer.<sup>1</sup> Inevitably, according to the pessimists, AI will produce robots that have tremendous power and behave immorally. These predictions certainly have some traction, particularly among a public that pays good money to see such dark films as Stanley Kubrick's *2001* and his joint venture with Stephen Spielberg, *AI*.

Nonetheless, we're optimists: we think formal logic offers a way to preclude doomsday scenarios of malicious robots taking over the world. Faced with the challenge of engineering ethically correct robots, we propose a logic-based approach (see the related sidebar). We've successfully implemented and demonstrated this approach.<sup>2</sup> We present it here in a general method-

ology to answer the ethical questions that arise in entrusting robots with more and more of our welfare.

## Deontic logics: Formalizing ethical codes

Our answer to the questions of how to ensure ethically correct robot behavior is, in brief, to insist that robots only perform actions that can be proved ethically permissible in a human-selected *deontic logic*. A deontic logic formalizes an ethical code—that is, a collection of ethical rules and principles. Isaac Asimov introduced a simple (but subtle) ethical code in his famous Three Laws of Robotics:<sup>3</sup>

1. A robot may not harm a human being, or, through inaction, allow a human being to come to harm.
2. A robot must obey the orders given to it by human beings, except where such orders would conflict with the First Law.
3. A robot must protect its own existence, as long as such protection does not conflict with the First or Second Law.

Human beings often view ethical theories, principles, and codes informally, but intelligent machines require a greater degree of precision. At present, and for the foreseeable future, machines can't work directly with natural language, so we can't simply feed Asimov's three laws to a robot and instruct it behave in

*A deontic logic formalizes a moral code, allowing ethicists to render theories and dilemmas in declarative form for analysis. It offers a way for human overseers to constrain robot behavior in ethically sensitive environments.*

# Toward a General Logician Methodology for Engineering Ethically Correct Robots

Selmer Bringsjord, Konstantine Arkoudas, and Paul Bello,  
Rensselaer Polytechnic Institute

**A**s intelligent machines assume an increasingly prominent role in our lives, there seems little doubt they will eventually be called on to make important, ethically charged decisions. For example, we expect hospitals to deploy robots that can administer medications, carry out tests, perform surgery, and so on, supported by software agents,

or softbots, that will manage related data. (Our discussion of ethical robots extends to all artificial agents, embodied or not.) Consider also that robots are already finding their way to the battlefield, where many of their potential actions could inflict harm that is ethically impermissible.

How can we ensure that such robots will always behave in an ethically correct manner? How can we know ahead of time, via rationales expressed in clear natural languages, that their behavior will be constrained specifically by the ethical codes affirmed by human overseers? Pessimists have claimed that the answer to these questions is: "We can't!" For example, Sun Microsystems' cofounder and former chief scientist, Bill Joy, published a highly influential argument for this answer.<sup>1</sup> Inevitably, according to the pessimists, AI will produce robots that have tremendous power and behave immorally. These predictions certainly have some traction, particularly among a public that pays good money to see such dark films as Stanley Kubrick's *2001* and his joint venture with Stephen Spielberg, *AI*.

Nonetheless, we're optimists: we think formal logic offers a way to preclude doomsday scenarios of malicious robots taking over the world. Faced with the challenge of engineering ethically correct robots, we propose a logic-based approach (see the related sidebar). We've successfully implemented and demonstrated this approach.<sup>2</sup> We present it here in a general method-

ology to answer the ethical questions that arise in entrusting robots with more and more of our welfare.

## Deontic logics: Formalizing ethical codes

Our answer to the questions of how to ensure ethically correct robot behavior is, in brief, to insist that robots only perform actions that can be proved ethically permissible in a human-selected *deontic logic*. A deontic logic formalizes an ethical code—that is, a collection of ethical rules and principles. Isaac Asimov introduced a simple (but subtle) ethical code in his famous Three Laws of Robotics:<sup>3</sup>

1. A robot may not harm a human being, or, through inaction, allow a human being to come to harm.
2. A robot must obey the orders given to it by human beings, except where such orders would conflict with the First Law.
3. A robot must protect its own existence, as long as such protection does not conflict with the First or Second Law.

Human beings often view ethical theories, principles, and codes informally, but intelligent machines require a greater degree of precision. At present, and for the foreseeable future, machines can't work directly with natural language, so we can't simply feed Asimov's three laws to a robot and instruct it behave in

*A deontic logic formalizes a moral code, allowing ethicists to render theories and dilemmas in declarative form for analysis. It offers a way for human overseers to constrain robot behavior in ethically sensitive environments.*

# Monday.

**The End**