

DYNAMIC AGENT SAFETY LOGIC: THEORY AND APPLICATIONS

A Thesis presented to
the Faculty of the Graduate School
at the University of Missouri

In Partial Fulfillment
of the Requirements for the Degree
Doctor of Philosophy

by
Seth Ahrenbach
Dr. Rohit Chadha, Thesis Supervisor
MAY 2019

The undersigned, appointed by the Dean of the Graduate School, have examined the dissertation entitled:

FORMAL METHODS WITH
DYNAMIC AGENT SAFETY LOGIC

presented by Seth Ahrenbach,
a candidate for the degree of Doctor of Philosophy and hereby certify that, in their opinion, it is worthy of acceptance.

Dr. Rohit Chadha

Dr. Alwyn Goodloe

Dr. William Harrison

Dr. Paul Weirich

ACKNOWLEDGMENTS

This page is where you would acknowledge all those who helped you with your academic research.

TABLE OF CONTENTS

ACKNOWLEDGMENTS	ii
LIST OF TABLES	vi
LIST OF FIGURES	vii
ABSTRACT	viii
CHAPTER	
1 Introduction	1
2 Background	11
2.1 Game Theory	11
2.2 Modal Logic	15
2.2.1 Modal Logic Syntax and Semantics	17
2.2.2 Epistemic Logic	21
2.2.3 Why S5 Is Appealing	23
2.2.4 Doxastic Logic	24
2.2.5 Knowledge and Belief Combined	25
2.3 Dynamic Epistemic Logic	28
2.4 Formal Methods Tools	30
3 Dynamic Agent Safety Logic	33
3.1 Syntax and Semantics	34
3.1.1 Syntax	34
3.1.2 Metalanguage	35
3.1.3 Semantics	40
3.1.4 Hilbert System	40

3.1.5	Static Base	42
3.1.6	Desiderata of a Static Base	44
3.2	Interpreting \mathbf{K}_i and \mathbf{B}_i	55
3.2.1	Justification	58
3.3	Indistinguishability Satisfied	61
3.4	Avoiding The Collapse	65
3.5	Soundness and Completeness	66
3.5.1	Soundness	66
3.5.2	Completeness	73
3.6	Discussion	88
3.7	Dynamic Extension	90
3.7.1	Inevitability and Minimum Rationality	93
4	Case Studies	95
4.1	Case Study and Mechanization	95
4.1.1	Air France 447	96
4.1.2	Mechanization in Coq	99
4.2	Additional Case Studies	109
4.2.1	Copa 201 and Asiana 214	109
4.3	Safety Properties	112
4.4	Decision Problem	114
4.4.1	Summary	115
5	Löb’s Obstacle	117
5.1	Löb’s Theorem	119
5.2	Discussion on Löb’s Obstacle	122
5.3	Epistemic Logics that Crash	126

5.3.1	$\mathcal{S5}$ Epistemic Logic	126
5.3.2	Hintikka's $\mathcal{S4}$ Epistemic Logic	127
5.3.3	Kraus and Lehmann System	128
5.3.4	\mathcal{DASL}	129
5.4	Avoiding Löb	132
5.4.1	Philosophical Considerations	136
5.5	Summary	139
6	Summary and concluding remarks	140
	BIBLIOGRAPHY	143
	VITA	152

LIST OF TABLES

Table		Page
2.1	Hintikka's Combined Logic of Knowledge and Belief	26
2.2	Kraus-Lehmann Combined Logic of Knowledge and Belief	27
3.1	\mathcal{DASL} semantics	40
3.2	\mathcal{DASL} 's dynamic model relations.	40
3.3	Hilbert System of \mathcal{DASL}	41
3.4	Static Base of \mathcal{DASL}	43
3.5	The reduction axioms of \mathcal{DASL}	91
3.6	\mathcal{DASL} axiom schemas of dynamic non-reductive character.	93
5.1	Logic of Kraus and Lehmann	128
5.2	Logic of \mathcal{DASL}	129
5.3	Static Logic of \mathcal{DASL}^-	133

LIST OF FIGURES

Figure		Page
2.1	A game between players A and B	12
2.2	\mathcal{M} : A simple counterexample using possible world semantics.	17
2.3	$M1$: The model before Alice announces “ φ ”.	29
2.4	$M2$: The model after Alice announces “ φ ”.	30
3.1	Action Model for “Alice flips a coin”	36
3.2	Action Model for “Alice peeks”	36
3.3	Action Model for “Alice secretly peeks”	37
3.4	A counterexample to $\mathbf{B}_i \varphi \Rightarrow \mathbf{K}_i \varphi$	65
5.1	A counterexample to $\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi$	131
5.2	A counterexample to $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{B}_i \varphi$	134
5.3	A counterexample to $\mathbf{B}_i (\mathbf{B}_i \varphi \Rightarrow \varphi) \Rightarrow \mathbf{B}_i \varphi$	135

ABSTRACT

This is the abstract of your dissertation project. It should not exceed one page.

Chapter 1

Introduction

In this doctoral thesis, we present a logic for reasoning about the connection among four aspects of human agency: knowledge, belief, mere action, and safe action. A mere action is an action that has either been performed or not; its evaluation as a good or bad action, relative to some normative domain, is not considered. A safe action, on the other hand, is a mere action that is evaluated in a safety domain. Other normative domains could be considered without much change to the logic being required, but in this dissertation I confine myself to the domain of aviation safety.

Knowledge and belief are aspects of human agency that have been well-studied by philosophers, logicians, and computer scientists. Epistemology is a core field of philosophy, studying what knowledge is, and logicians have formalized these concepts. Computer scientists have taken these formalized systems and applied them to computational domains ranging from artificial intelligence to information security. The intersection of these fields at the knowledge nexus provides a rich variety of topics to study.

Similarly, the study of action has a history of richness at the intersection of philosophy, logic, and computer science. The relationship between actions and knowledge in particular finds rich application in the fields of economics and game theory. The

logics for reasoning about these things have been actively developed over the past several decades, with a species of dynamic modal logic coming to the forefront as particularly well-suited to the task.

Dynamic logic itself originated in computer science as a way to formally study the executions of programs. Philosophers and logicians then borrowed it back and married it with epistemic logic in order to generate logics of knowledge and action. The exchange of ideas between fields often occurs through the logics they use to study their subjects.

This dissertation picks up the thread from philosophical logic, specifically the efforts to use dynamic epistemic logic to model behavior in the rarefied world of games and toy examples. There researchers have made progress formalizing simple scenarios where one agent announces to the other something, and this causes the other agent’s knowledge to increase. We extend that logic here in two ways. First, we experiment with its application to real humans piloting aircrafts. The environment is more complex, and challenges the logic to become a more realistic representation of human knowledge. To address this, we weaken the axioms of knowledge, and include axioms of belief and the relationship between knowledge and belief. The second extension concerns the normative aspect of actions. Dynamic epistemic logics impose preconditions on actions that determine whether they can be performed. We extend this idea with another precondition governing whether the action can be safely performed. It is a simple idea that opens up a wide new world of expressive power for the logics, as the underlying normative precondition function can be almost any normative aspect one cares to consider. One could swap out our safety precondition function with an ethical precondition function, or a legal precondition function, and develop similar dynamic logics of ethical or legal agency. Those applications are beyond the scope of this dissertation, however.

In order to validate the above extensions, we rely on rigorous tools to check our

math, so to speak. Many other researchers have made great strides in mechanizing the proof and model theories of modal logics in the interactive theorem proving tool Coq. We follow their lead and similarly extend their efforts to bear out the logic we develop in the target application domain. We modify a deep embedding of dynamic epistemic logic in Cog to create a deep embedding of our logic. We verify the theorems of the object language in the target domain in Coq, as well as the soundness and completeness proofs in the metalanguage in Coq. While this is not a novel approach to logic research, it represents a more rigorous and computational approach to the methods of logic research that help move the dissertation from the merely philosophical domain and into the computer science domain. We do not develop any new theorem proving tools or techniques in the thesis, but our use of the interactive theorem prover itself is a feature that distinguishes this dissertation as one of computer science.

We have chosen the target domain of aviation safety, and it has a nice property that makes it a good candidate domain for extending the logic. Rather than modelling a wide array of behaviors and information that apply to human behavior generally, we need only model the behavior and information that relates to flying an airplane. We need to model the instrument readings and the pilot’s inputs to the flight controls. This significantly cuts down the expressive space we need to handle. However, it is still a real human in a real world situation, requiring the logic to relax certain assumptions about the reasoning power of the agents and the information they have about the state of the world. Instrument readings can be conflicting, and pilots can be unaware of this conflict. This is not typically a situation logicians have tried to model in the past. Additionally, aviation safety is already a domain familiar with the use of formal methods to increase assurance, so the approach is not so novel as to be completely foreign to the target domain’s researchers. This gives us related work to compare and contrast the present efforts to.

In particular, various logics exist for modelling safety properties of avionics sys-

tems, and tools for verifying that systems satisfy those properties. Aviation historically has been one of the few domains where the cost, perceived or actual, of formal verification is justified by the risks associated with system failure. We have engaged with this research community to some extent through collaborative research, conference presentations, and journal articles, and have received very helpful feedback. In some ways the approach adopted in this dissertation is too heavy handed for the particular needs of formal methods researchers in aviation safety. Rushby makes this case particularly well in his response to Ahrenbach and Goodloe [6]. Many of the benefits achieved by the approach pursued here can be had through more straight forward integrations of modal logic with existing formal verification frameworks. Rushby achieves analogous results with merely a modal operator for a pilot’s belief sufficing, with the belief treated as an input variable to the verification system just like any other.

However, our purposes here differ slightly from those of researchers focused on the aviation domain. An aviation-focused researcher asks ‘what is the best way to use this formal method in order to make aviation safer?’, while we ask, ‘What domains are suitable proving grounds for this new general formal method?’ We seek a level of generality beyond aviation, and so in some cases are unwilling to sacrifice the heavy machinery for the benefit of having a simpler pilot monitoring logic. We seek to develop a general logic of realistic agency, and we validate it in the proving ground of aviation safety. Just as a geneticist is interested in *Drosophila melanogaster* (fruit flies), or an artificial intelligence researcher is interested in chess, because these subject (or target) domains are areas to apply new techniques and test new approaches [35, 40]¹.

The new approach pursued here, as mentioned earlier, is to relax the assumptions

¹In his book about the building of *Deep Blue*, Hsu decries being called an AI researcher, and prefers his work to be seen as Very Large Scale Integration (VLSI) research, but the point stands that chess to him was a target domain useful as a proving ground for his primary research interest.

about agent rationality that other philosophical logics make, while augmenting the dynamic modality of action to include considerations of safety. If all goes well, this yields a logic for reasoning about realistic agents performing safe or unsafe actions, and about their beliefs and knowledge about those actions.

The logic, which I call Dynamic Agent Safety Logic (*DASL*), is based on the logical foundations of game theory, in which models of agency formally capture how knowledge, rationality, and action relate to each other. Game theory presents a model that, given a description of a scenario, allows one to deduce what actions are dictated by a given theory of rationality. The standard game-theoretic inference works as follows:

$$Knowledge_of_Situation \wedge Rationality \Rightarrow Good_Action.$$

One can read this as, “if an agent has knowledge of a situation (*e.g.* a game), and the agent is rational, then the agent executes a good action. In game theory, the important terms are suitably formalized for mathematical treatment. Knowledge is assumed to be perfect, rationality is defined as the maximization of some measure of utility, and good actions are those that bring about the outcomes with the most possible utility payoffs. The definitions make the inference an analytic truth.

Empirically, however, humans frequently deviate from the prescribed behavior. Looking at the above formula, we can ask a question: what can we infer when an agent fails to execute the prescribed action, as when pilots provide unsafe control inputs to their aircrafts? We can answer this question by examining the contrapositive of the above game-theoretic inference:

$$\neg Good_Action \Rightarrow \neg(Knowledge_of_Situation \wedge Rationality),$$

or equivalently,

$$\neg \textit{Good_Action} \Rightarrow \neg \textit{Knowledge_of_Situation} \vee \neg \textit{Rationality}.$$

With a bit more Boolean manipulation, we have the following:

$$\neg \textit{Good_Action} \wedge \textit{Rationality} \Rightarrow \neg \textit{Knowledge_of_Situation}.$$

This can be read, “If an agent is rational but executes a bad action, then the agent lacked knowledge of the situation.” Thus, embedded in the classical game-theoretic model of agency is a logical inference from bad action to missing knowledge. This makes intuitive sense upon reflection. If someone is rational, yet they commit an irrational (read: “bad”) action, then it must be the case that they didn’t know some crucial information. With this insight in hand, I identify a logic in which the above inference is sound, with details about which particular pieces of information are missing from an agent’s knowledge base when she executes a bad action. Again, it should not be surprising that such a logic exists, because classical game theory already posits a *logical* relationship between knowledge of particular propositions and particular actions.

I have formally captured such inferences with *DASL*, where a (mostly) rational agent executes a bad action, and from this we can infer which safety-critical information they are missing. By developing a logic that can model the information flow in these situations, we advance the project of formally reasoning about human agency. Credit for initiating this project belongs to many researchers over the years, especially philosophers in the analytic tradition concerned with analyzing the epistemology and metaphysics of agency in a rigorous fashion. I have been most influenced by the works of the Amsterdam school of modal logic, led by Professor Johan van Benthem, where the efforts center around rich combinations of modal logics in order to model human

agency. In [33], they develop a modal logic for reasoning about the knowledge and decision-making of agents in games, and in [54], van Benthem explores similar themes around information flow and interaction.

The history of research in this area dates back to the 1950s and '60s in the development of temporal logic (or tense logic) by Arthur Prior [34], a graph theoretic semantics by Saul Kripke [37], and logics for belief and knowledge due to work by Rescher [75], von Wright [80], and Hintikka [59]. These logics formalize reasoning about what was true, what will be true at some point, what is always going to be true, what is believed to be true, and what is known to be true. Each of these modifiers is a truth modality, and hence they each constitute a modal logic. The Amsterdam school and others built on these methods, especially the work by Patrick Blackburn [20], which illustrates the general features of any modal logic as a tool for reasoning about systems that can be modeled as graphs from an *internal* perspective, whereas first order logic reasons about such systems from an external perspective. A first order formula might say what is true of the object x , from a global perspective, but modal logic allows us to formalize truth from x 's perspective. If the first order formula is $\forall x, \exists y : R(x, y) \wedge P(y) \Rightarrow P(x)$, the corresponding modal formula would simply be $\Diamond_{Rp} \Rightarrow p$. They both say the same thing: For any node x in the graph, if it can reach a node y by relation R , and y is a node where the proposition “ y is P ” holds, then “ x is P ” holds. The former explicitly quantifies over the nodes, and the latter does so implicitly through the semantics, to be described later.

While these developments occurred in what might be called the philosophical branch of modal logic research, researchers in economics explored the mathematical foundations of game theory. Aumann, in [73], showed the axioms of epistemic logic that must be assumed in order for classical game theoretic results to hold. The agents in classical games, otherwise known as *homo economicus*, are ideally rational, with perfect knowledge of their situations. We will meet these axioms and modify them

for our purposes later.

Advances in dynamic modal logic come from theoretical computer science, where computer programs are modeled as state transition diagrams. As a program executes, the computer transitions from state to state, where each state is a collection of values assigned to variables, and each transition is a simple action executed. As this formalization lends itself to graph theoretic representation, it lends itself to formalization in modal logic, per Blackburn's insight. There are two main approaches to applying modal logic to the analysis of programs. The first is a static approach, where the entirety of the program's execution tree is modeled at once. Transitions from each state are captured by temporal logic. A program might be formalized in temporal logic, and the following theorem might be proven of it: *At the source of the execution graph, it will always be true that bad event B does not occur* [45]. The second approach is dynamic, where each simple transition action A or B gets a modal operator, which allows us to reason about what happens after every execution of subprocess A, or after some executions of subprocess B, etc [39].

One thing that is easy to do with dynamic modal logic but somewhat complicated in the static approach is to model actions and knowledge. An epistemic logic is modeled by a static Kripke structure, and in the dynamic case this structure changes as the agents act and learn different things. The static approach requires a grand two dimensional Kripke structure with one dimension capturing the epistemic relations at a moment, and one dimension capturing the temporal relation as actions move forward through time. For an example of this approach, see John Horty [60]. For the dynamic approach, see van Ditmarsch *et al.* [62].

Van Benthem and the Amsterdam school identified these various threads dispersed around campuses and saw how they related to each other, and how they might be fruitfully combined for various ends. As philosophers, they were mostly concerned with using the rich tools from economics and computer science to analyze human

agency robustly and accurately. Modal logics offer tremendous expressive power at often a lower cost than first- or second order logic, because modal logics take an internal view of the graphs they reason about and are defined by. This often means that a powerful, useful modal logic can be defined that is also sound, complete, and *decidable*. So, by carefully defining the modal operators for, say, knowledge, preference, and action, a modal logic of game theory can be developed; not just the epistemic aspect of the agents, but the games themselves, which van Benthem calls a Theory of Play.

Applications of the Theory of Play have thus far been limited to relatively simple, artificial examples, in the same way methods in genetics research are often developed on fruit flies. In this thesis, I continue this work by extending application of the methods to richer real-world cases of humans in cockpits, which for reasons mentioned earlier make good cases for early forays into the formal modeling of human behavior. Just as genetics methods mature and eventually apply to humans, so must modal logic methods mature and apply to real humans in the world.

Dynamic Agent Safety Logic allows us to reason about which critical information is not being delivered to the human's brain. Because we can deduce which safety-critical information is missing from her knowledge base, we can automatically act to correct this failure of delivery, and therefore build systems that have a high assurance that the information will be delivered. I apply this technique to aviation safety as a formal method, but in principle it could be applied to other domains of human agency that meet certain conditions. Some examples that strike me as plausible include doctors and nurses in emergency rooms, cybersecurity analysts monitoring network traffic alerts, power plant operators, and of course motor vehicle drivers. During crises, these environments can quickly become saturated with alarms, and humans quickly suffer from information overload. If actions can be properly related to instruments such that unsafe actions can be detected, then the agent's missing knowledge can be

deduced and rectified.

In order to adequately model human agency, we must depart from the familiar and comfortable logic that serves as the static base for most dynamic approaches: *S5*. This departure represents an innovation, which some in the formal methods community disagree with, so we must devote some time to motivating the move. We make the case that departing from *S5* as a static base representing human knowledge and belief makes sense for philosophical and technical reasons.

In what follows, we describe the relevant background material in Chapter 2, including the foundations of game theory and the logical models of agency informing our new developments. Chapter 3 motivates departure from *S5* as a static base and presents *DASL*, then proves that it is sound and complete using the Coq Proof Assistant. Chapter 4 applies *DASL* to three aviation mishaps, again formalized in Coq. Chapter 5 considers a problem facing logics that allow for reasoning about self-referential sentences, due to Gödel and Löb’s discoveries. We present Löb’s Obstacle that these logics face, and show how to relax *DASL* in order to avoid it, when the agents being modeled must be capable of self-reflection. In 6 we summarize and conclude.

Chapter 2

Background

This chapter describes the context in which this dissertation makes advances. Section 2.1 lays out the basics of game theory, which provides a model of agency that we target and make more realistic. Section 2.2 introduces modal logic, upon which the logic presented in this thesis is based. We also describe the various epistemic and doxastic logics that have influenced *DASL*'s static base. Section 2.3 presents Dynamic Epistemic Logic and the so-called dynamic turn, which combines modalities for reasoning about how models change due to actions. *DASL* extends this idea by including modalities for safe actions in addition to mere actions. Section 2.4 describes the formal tools we use to mechanically check that *DASL* is sound and complete, that its application to aviation safety is without errors, and to encode the inference of safety-critical information. We use Coq, which is both a programming language and a tool for verifying proofs.

2.1 Game Theory

This section presents the basics of game theory in order to motivate the model of agency that *DASL* seeks to formally capture. However, this dissertation does not seek to make advances in game theory itself. Future research could take the resulting logic and use it as a new foundation for game theoretic reasoning among more realistic agents, but that task is not taken on in the present work. We begin this section with a brief overview and an example game, and then present the model of agency required for the standard solution to the game.

Game theory is a mathematical model for strategic reasoning. Strategic reasoning refers to the way an agent reasons in situations where her payoffs depend on the actions of other agents in addition to her own, and in which she knows about these dependencies. For turn-based games, the mathematical structure employed is a *game tree*, where each node represents a player's turn, and each edge the transition via a player's action. The leaves of the tree represent the payoffs each player receives at the end of the game. This paper is not concerned with the games themselves, but rather with the underlying assumptions about agency that entail their solutions. We briefly illustrate these underlying assumptions with the following example.

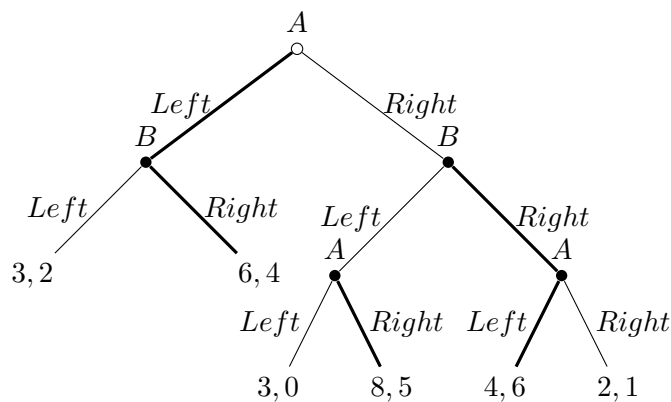


Figure 2.1: A game between players A and B

We can see in the figure below that the first player to act is Player A, at the root node. Her choices are to move *Left* or *Right*. Player B faces similar choices at

the resulting nodes, and the players alternate turns until the game ends and they receive their payoffs, listed (A,B) . Briefly glancing through the outcomes, it looks like A should aim for the node with payoff $(8,5)$, because 8 is the highest payoff available to A. However, the official solution to the game is for A to first go *Left*, and then for B to go *Right*, resulting in a payoff of $(6,4)$, where both get less than the intuitively appealing outcome! Why is this so?

Game theory makes strong assumptions about agent knowledge and rationality. The solution to this game is reached through an algorithm called *backward induction* [33]. The players reason by starting at each end node and looking to the immediate parent node, and asking what the deciding player will do at that node, assuming she will choose the path with the highest payoff. So, at the bottom right of the figure, player A is to act, and she can go *Left* for a payoff of 4, or *Right* for a payoff of 2. So, she will go *Left*, illustrated by the bold line. The end nodes not selected are subsequently eliminated. This process is repeated at each end node. Then it is recursively applied up the tree. So along the right branch, player B decides between *Left* for a payoff of 5, or *Right* for a payoff of 6, because B knows that A is rational, and he knows how she will act at each end node. A, at the root, then must choose between *Left* for a payoff of 6, or *Right* for a payoff of 4, because she knows that B is rational, and knows that B knows that she is rational. The explanation begins to illustrate the assumptions game theory makes about each player's knowledge. In fact, this only scratches the surface.

Game theory, and classical economics in general, makes the following assumptions about agent knowledge, formalized in epistemic logic [73].

Agency Model in Classical Game Theory.

- (1) $\mathbf{K}_i(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \psi)$
- (2) $\mathbf{K}_i \varphi \Rightarrow \varphi$
- (3) $\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi$

- (4) $\neg \mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \neg \mathbf{K}_i \varphi$
- (5) $\mathbf{C}_G((1) \wedge (2) \wedge (3) \wedge (4) \wedge (5))$.

This forms an idealized model of the knowledge component of classical game theory’s agents. \mathbf{K}_i is a modal operator for knowledge, and $\mathbf{K}_i \varphi$ reads, “agent i knows that φ .” \mathbf{C}_G is a modal operator for common knowledge, the fixpoint for “everyone in group G knows that everyone knows that...” The agents are logically omniscient due to (1), knowledge implies truth with (2), agents have *positive introspection* with (3), and *negative introspection* with (4). Assumptions (1), (3), (4), and (5) are somewhat dubious. The model also fails to formally represent other aspects of agency, like action and evaluation of outcomes. The model we propose makes weaker, more realistic assumptions about knowledge, includes a modal operator for belief, and formally represents action and the evaluation of actions as either safe or unsafe.

Recent work at the intersection of game theory and logic focuses on the information flow that occurs during games. Van Ditmarsch identifies a class of games called *knowledge games*, in which players have diverging information [56]. This slightly relaxes the assumption of classical game theory that players have common knowledge about each other’s perfect information. Similarly, it invites logicians to study the information conveyed by the fact that an action is executed. For example, if the action is that agent 1 asks agent 2 the question, “ p ?”, the information conveyed is that 1 does not know whether p , believes that 2 knows whether p , and after the action occurs, this information becomes publicly known. The logic modeling games of this kind is of particular interest to us, as we are concerned with identifying the knowledge and belief state of human pilots based on their actions.

The proceeding sections introduce the various logical systems that form a foundation for the work of this thesis, starting with modal logic in its traditional philosophical interpretation, and expanding to epistemic and doxastic logic. Then, we introduce dynamic logic, and its expansion into Dynamic Epistemic Logic and Public

2.2 Modal Logic

Aristotle noted a distinction between contingent truth and necessary truth, and some Medieval philosophers continued this line of inquiry. Necessary truths could not have been otherwise. The definitions of natural numbers and the addition function guarantee that in all possible worlds, $2 + 2 = 4$. On a plane, the truths of Euclidean geometry are necessarily true. “There is life on Saturn’s moon Enceladus” is a possible truth. “Enceladus has water on it” is a contingent truth. What about “Water is H₂O”? Modal logic was formalized in the early 1900’s by C. I. Lewis and has recently had a resurgence in multidisciplinary interest [55]. At its core, modal logic allows us to reason about necessary and possible truths through the use of modal operators. What follows is a brief illustration of modal logic’s concepts and formalisms.

When philosophers talk about necessity they usually mean *metaphysical* necessity. In addition to formalizing this notion for systematic reasoning, modal logic is used for clarifying what exactly this means [81]. Simple and intuitive examples are those from mathematics and notions of identity. Suppose p is the arithmetic expression “ $2 + 2 = 4$ ”. Obviously, p is true in the actual world. We can make the stronger claim that p is necessarily true, but what does this mean? As informal shorthands, initial attempts to define necessary truths might appeal, as I did above, to the claim that they could not possibly have been false. But then what do we mean by “possible”? The formal semantics for dealing with these questions comes from Arthur Prior and Saul Kripke, and we introduce that machinery in the next section [34, 37]. For now, we say a statement is necessarily true if and only if it is true in all worlds we consider possible. The modal operator for necessity makes the formal statement: $\Box p$. Consider the following inference, where \Rightarrow means “implies”: $\Box p \Rightarrow p$. This reads,

“Necessarily p implies p ,” or equivalently, “If p is necessarily true, then p is true.” If p is necessarily true, is p true? Intuitively, the answer is ‘yes’, and indeed a modal logic of metaphysical necessity includes this axiom for all formulas φ : $\Box\varphi \Rightarrow \varphi$.

What other inferences can we make, based on our intuitive notion of necessity and possibility? What about “if p is true, then p is possibly true”? To formalize this, we need the modal operator for possibility: \Diamond . So, the modal formula would be $p \Rightarrow \Diamond p$. This seems true as well, and indeed it is a theorem for metaphysical modal logic: $\varphi \Rightarrow \Diamond\varphi$.

It is obvious then that $\Box\varphi \Rightarrow \Diamond\varphi$ is a theorem. This states that if something is necessarily true then it is possibly true. What about the other direction: $\Diamond\varphi \Rightarrow \Box\varphi$? It turns out this is not a theorem under the typical notions of necessity and possibility. But this raises a question about how we would present a counterexample that disproves it. To do this, we need a semantics for the logic. The semantics we use are called *possible world semantics*, and they are usually attributed to Saul Kripke [37]. One would be hard-pressed to find a species of modal logic, whether in economics, computer science, or philosophy, that does not use possible world semantics in some form or another. Sometimes they are referred to as *Kripke semantics*.

In possible world semantics, a graph structure is created with worlds as nodes and accessibility relations among worlds as the edges in the graph [20]. Propositional formulas are true or false at each world. These graph structures are typically called Kripke structures. We can define the following Kripke structure, $\mathcal{M} = \{W, R, V\}$, where W is a finite set of worlds, $\{w, v\}$, R is a binary accessibility relation defined on those worlds $\{(w, v), (v, w)\}$, meaning w has access to v and v has access to w , and V is a *valuation* function, which maps propositions to sets of worlds at which they are true. For example, if p is true at w , then $w \in V(p)$. In our model we only care about the proposition p , which now stands for some contingent proposition, like “all swans are white”. Formally, we say $w \notin V(p)$ if not all swans are white in world

w , denoted by $\neg p$, while $w \in V(p)$ if they are, denoted by p . The following figure illustrates \mathcal{M} :

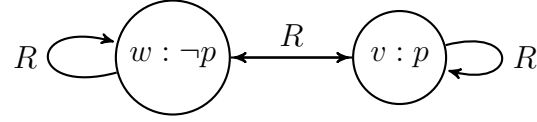


Figure 2.2: \mathcal{M} : A simple counterexample using possible world semantics.

According to possible world semantics, $\Box\varphi$ is true at a world w , written $w \models \Box\varphi$, if and only if for all worlds v such that $R(w, v)$ (v is related to w by the R relation), $v \models \varphi$. This says a formula is necessarily true at a world if and only if it is true at all worlds accessible by that world according to the underlying R relation. Similarly, $w \models \Diamond\varphi$ if and only if there is some world v such that $R(w, v)$ and $v \models \varphi$. In \mathcal{M} , w is R -accessible to itself, so there is a world accessible to w where p is false, and thus $w \models \neg\Box p$. However, since v is R -accessible to w , and $v \models p$, it is true that $w \models \Diamond p$. Thus, we have $w \models \Diamond p$ and $w \models \neg\Box p$, a negation of $\Diamond p \Rightarrow \Box p$, so it cannot be the case that $\Diamond\varphi \Rightarrow \Box\varphi$ is a theorem, for arbitrary formula φ .

There are many systems of modal logic. The way any is distinguished from another is based entirely on the definition of R . For the popular $\mathcal{S5}$, which we will examine later, R is a Euclidean and reflexive binary relation on worlds. Thus, that is how “possibility” is formally defined, and by extension, “necessarily”. That is also what determines the axioms and theorems comprising the logic. We formally define these notions in the next section, along with the syntax and semantics of propositional modal logic.

2.2.1 Modal Logic Syntax and Semantics

This section formally defines the syntax, what the logic looks like, and the semantics, what the truth conditions are for the logic.

Recall that Boolean logic is a simple logic for reasoning about basic propositions using the logical connectives ‘and’, ‘or’, ‘not’, and ‘if...then’. It forms the foundation of most logics and has applications ranging from philosophy to circuit design. Propositions are represented as constants p , q and well-formed formulas of the language are constants and any proper combination of constants using the above logical connectives, represented symbolically as,

$$\varphi \stackrel{def}{=} p \mid \varphi \wedge \varphi \mid \varphi \vee \varphi \mid \neg \varphi \mid \varphi \Rightarrow \varphi.$$

As illustrated in the previous section, modal logic adds to propositional logic with modal operators for necessary and possible truths. The syntax for propositional logic is extended in the following way to make modal logic:

The semantics for Boolean logic are simply truth tables for each connective, which I will not reproduce here. However, the operators in modal logic are not truth functional, and require more complex semantics, which we earlier mentioned are called possible world semantics. They are as follows.

Let $\mathcal{M} = \{W, R, V\}$ be a model such that W is a set of possible worlds, R is a binary relation on those worlds, and V is a valuation function mapping atomic propositions to the sets worlds satisfying them,

$$w \models p \text{ iff } w \in V(p)$$

$$w \models \neg \varphi \text{ iff } w \not\models \varphi$$

$$w \models \varphi \wedge \psi \text{ iff } w \models \varphi \text{ and } w \models \psi$$

$$w \models \Box \varphi \text{ iff } \forall v, R(w, v) \text{ implies } v \models \varphi.$$

$$w \models \Diamond \varphi \text{ iff } \exists v, R(w, v) \text{ and } v \models \varphi.$$

The character of a modal logic is determined by the binary relation on worlds underlying the modal operators. We mentioned $\mathcal{S5}$ in the previous section as having a R relation in which every world is accessible to itself by the binary relation. It is also one in which the R relation is Euclidean. These conditions are called a *frame* conditions, and the models that satisfy these conditions belong to said frame. All reflexive frames have the following frame conditions:

$$\forall x, R(x, x) \tag{2.1}$$

Likewise, all reflexive frames have the following axiom:

$$\Box\varphi \Rightarrow \varphi \tag{2.2}$$

It is not just a stipulation that all reflexive frames must have that axiom: They have that axiom *because* they have that frame condition. This is due to correspondence theory, which we explain later on.

Relaxing a frame condition changes the R relation, and in doing so changes the logic. If we remove the reflexivity condition, the above axiom is no longer an axiom for the logic defined on \mathcal{M} . This means that we can specify the axioms we want by specifying the frame condition on the accessibility relation. Each frame condition corresponds to a modal logic axiom. In addition to reflexivity, the other common frame conditions are as follows, with their corresponding modal logic axiom:

- Transitivity

$$\forall x, y, z R(x, y) \wedge R(y, z) \Rightarrow R(x, z) \tag{2.3}$$

$$\Box\varphi \Rightarrow \Box\Box\varphi \tag{2.4}$$

- Symmetry

$$\forall x, y \ R(x, y) \Rightarrow R(y, x) \quad (2.5)$$

$$\varphi \Rightarrow \Box \Diamond \varphi \quad (2.6)$$

- Euclidean

$$\forall x, y, z \ R(x, y) \wedge R(y, z) \Rightarrow R(x, z) \quad (2.7)$$

$$\Diamond \varphi \Rightarrow \Box \Diamond \varphi \quad (2.8)$$

- Seriality

$$\forall x \exists y, R(x, y) \quad (2.9)$$

$$\Box \varphi \Rightarrow \Diamond \varphi \quad (2.10)$$

These conditions are not exhaustive, but they represent some of the commonly combined conditions used to define axioms of different modal logics. By combining frame conditions, a modal operator is defined with the properties desired. For example, if we wish to define a modal operator for reasoning about the knowledge of ideally rational agents, as is done in Fagin *et. al.*[57], we impose the frame conditions of reflexivity and Euclidicity (Euclideaness), or one that is transitive and reflexive, as in Hintikka [59]. However, if we wish to develop a logic for belief, as the previously cited works also do, we must impose transitivity, Euclidicity, and seriality. The reasons for this are explored in the following sections.

2.2.2 Epistemic Logic

Epistemic logic began with philosophical concerns about knowledge [57, 59, 75]. Presently, it is likely to be studied in computer science as the logic for reasoning about ideally rational agents in well-structured environments. The agents are ideally rational in the sense that there is no bound on how much they can be said to know, nor on what propositions they are aware of at any one time. Thus, they have no problem conceiving of every possible sequence of moves a game may consist of, evaluating all possible outcomes, and deducing the optimal sequence of moves in order to maximize their own utilities. They are even stronger than contemporary computers in this way, which are bounded by time and space, and so are unable to actually compute the ideal strategy for an otherwise solvable game like Go. An ideally rational agent can solve Go and compute the game tree all the way to its end. These are the agents of game theory, which always select the optimal move when they have perfect information about the game and the other players.

The axioms that specify this level of knowledge are those introduced in Section ??, and recounted here without reference to common knowledge, which complicates things too much for our purposes:

Agency Model in Classical Game Theory.

- (1) $\mathbf{K}_i(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \psi)$
- (2) $\mathbf{K}_i \varphi \Rightarrow \varphi$
- (3) $\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi$
- (4) $\neg \mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \neg \mathbf{K}_i \varphi$.

These axioms make the knowledge operator an *S5* modal operator. Strictly speaking, including (3) is unnecessary because it follows as a theorem from (2) and (4), but we include it so that the reader is not surprised when we refer to it.

The first axiom holds for all *normal* modal logics, and under the epistemic interpretation, it states that if an agent knows that φ implies ψ , and she knows φ , then

she knows ψ . This is intuitive enough on the first pass, but can lead to a problem of global skepticism. If i knows that she is an embodied agent in the external world only if she knows that she is not a brain in a vat perceiving a simulated reality, then she knows she is an embodied agent only if she knows she is not a brain in a vat. It seems clear that she cannot know that she is not a brain in a vat, because her sensory experience is compatible with either the embodied scenario of the brain in a vat scenario. Therefore, she does not know she is an embodied agent in the external world. Thus, many philosophers question whether knowledge for human-like agents is closed under logical implication like this.

The second axiom states that a known proposition must be true, and finds ample support from philosophers devoted to studying the nature of knowledge. It can be considered the property of knowledge that most distinguishes it from belief, because beliefs can be false.

The third axiom, *positive introspection*, states that an agent knows something only if she knows *that* she knows it. This imposes a high standard on knowledge, under the assumption that an agent can identify the conditions that guarantee the truth of a known proposition, and in being able to do so, are justified to a sufficient degree. However, this property is largely rejected by philosophers, as this requirement is thought to be largely unattainable.

Finally, *negative introspection* states that an agent does not know something only if she knows that she does not know it. This has an undesirable effect of creating knowledge out of ignorance, for if i does not know that she does not know φ , this is sufficient for inferring that she knows φ , according to negative introspection. This is fine for ideally rational agents, and indeed it is a noble aspirational goal to know what one is ignorant about, but for human-like agents it is entirely unrealistic as an axiom.

These axioms are firmly rooted in the literature of formal epistemology and epis-

temic logic, and relaxing them in order to model more realistic agents requires a word or two. This section introduces and defends our relaxation of the classical axiomatization.

2.2.3 Why S5 Is Appealing

Before presenting our relaxed axiom schema for knowledge, we acknowledge the appealing properties of the S5 knowledge operator. Logicians who adopt the S5 knowledge operator have good reasons for doing so, both technical and intuitive. The technical appeal lies in the fact that the S5 operator is well-behaved from a logical point of view. Epistemic logics with S5 operators are decidable because S5 operators have semantics that satisfy the finite model property, meaning that every invalid proposition has a finite counter-model. The S5 knowledge operator allows a formula prefixed by any arbitrary combination of \mathbf{K}_i and $\langle \mathbf{K}_i \rangle$ symbols to be reduced to a formula prefixed by at most two. This keeps the models small.

On the intuitive side, the relations that define the semantics of an S5 knowledge operator are equivalence relations. They are defined as reflexive and Euclidean, which implies that they are also transitive and symmetric. Equivalence relations capture a notion of indistinguishability because each world in the equivalence class is indistinguishable from each other. All the modal formulas true at one of the worlds in an equivalence relation are true at each of the other worlds in the relation, so from a modal perspective, they are indistinguishable. When the modality is knowledge, this means that relative to the agent's knowledge, the worlds are indistinguishable.

These properties of $\mathcal{S5}$ epistemic logic are strong points in its favor. We suppose that any formal system seeking to model human-like knowledge must capture the intuitive idea that an agent does not know *whether* a proposition φ is true or false just in case she cannot distinguish the world she is in from the possible worlds in which either might be the case. Distinguishability seems to imply that there is some

evidence possessed by the agent that allows her to rule out possibilities, and therefore gain knowledge.

The technical benefit of the finite model property makes $\mathcal{S5}$ a friendly logic to work with, but it is not a reason to believe that the theory of knowledge described by $\mathcal{S5}$ is true of human agents. Therefore, we do not see the need to impose this condition as a requirement for realistic epistemic logics. However, it would be desirable to achieve the finite model property while maintaining realism.

Next we turn to a brief summary of doxastic logic, which shall play a role in \mathcal{DASL} 's static foundation.

2.2.4 Doxastic Logic

This section presents doxastic logic, the logic for reasoning about belief. Hintikka's formalization in [59] of the belief operator is standard, and we adopt it here. The operator \mathbf{B}_i for 'agent i believes ℓ ' is defined by a binary relation on worlds that is serial and Euclidean. It is a normal modal operator, and the property of transitivity follows from seriality and Euclidicity. So, the model of an agent's belief is defined by the following axioms,

Hintikka Belief Model.

- (1) $\mathbf{B}_i(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{B}_i\varphi \Rightarrow \mathbf{B}_i\psi)$
- (2) $\mathbf{B}_i\varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$
- (3) $\mathbf{B}_i\varphi \Rightarrow \mathbf{B}_i\mathbf{B}_i\varphi$
- (4) $\neg\mathbf{B}_i\varphi \Rightarrow \mathbf{B}_i\neg\mathbf{B}_i\varphi$.

Clearly the primary difference between this model of belief and the model of knowledge from classical game theory is axiom (2). Rather than guaranteeing truth, beliefs guarantee consistency. One could correctly argue that this is too strong of an assumption for humans, but we leave relaxations of the belief model to future work, and are content to relax the idealizations about knowledge in this dissertation.

Similarly, reasonable cases could be made against axioms (3) and (4).

The logic associated with these axioms is called $\mathcal{KD}45$. On the doxastic interpretation, it represents a slight idealization of how humans actually adopt beliefs. Clearly, humans can have contradictory beliefs, and they may believe something without believing that they believe it. Cases of cognitive dissonance like that are idealized away, not only because modeling such a psychology is complicated, but also because we seek a logic that maintains some level of normativity. If someone believes two contradictory propositions, she should be able to give one up if this contradiction is pointed out to her, and furthermore she is rational to do so. This reflects Hintikka's goal of developing logics of knowledge and belief that allow us to consistently identify when attributions of knowledge and belief are in some sense criticizable. Maintaining a notion of normativity in a logic is important to us, and underlies some the idealizations we accept in \mathcal{DASL} 's static base logic.

Knowledge involves a condition that is external to the agent's direct perception: the truth or falsity of the proposition. Belief, on the other hand, is entirely internal to the agent's perception. She can introspect on her evidence and deliberate on a proposition, and this is sufficient to see that she believes or disbelieves it, and this attitude is immediately available to her. For this reason, the positive and negative introspection axioms are more appropriate for belief than they are for knowledge.

2.2.5 Knowledge and Belief Combined

In addition to developing a stand-alone axiomatization of the belief operator, Hintikka combined the knowledge and belief operators into a single multi-modal system. Other philosophers and logicians have constructed similar systems. Hintikka considers the following combined axioms, which he calls conditions, as options:

Hintikka Combinations Considered.

- (1) $\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$

- (2) $\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi$
- (3) $\mathbf{B}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{B}_i \varphi$
- (4) $\mathbf{B}_i \varphi \Rightarrow \mathbf{K}_i \varphi$.

He subsequently rejects (4) immediately, and through deliberation (3), while he accepts (1) and (2). Thus, Hintikka's combined logic of knowledge and belief is,

$\mathbf{K}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \psi)$	Distribution of \mathbf{K}_i
$\mathbf{K}_i \varphi \Rightarrow \varphi$	Truth Axiom
$\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi$	Positive Knowledge Introspection
$\mathbf{B}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \psi)$	Distribution of \mathbf{B}_i
$\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$	Belief Consistency Axiom
$\neg \mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \neg \mathbf{B}_i \varphi$	Negative Belief Introspection
$\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi$	Knowledge implies Belief
$\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$	Undeniable Knowledge
From $\vdash \varphi$ and $\vdash \varphi \Rightarrow \psi$, infer $\vdash \psi$	Modus Ponens
From $\vdash \varphi$, infer $\vdash \mathbf{K}_i \varphi$	Necessitation of \mathbf{K}_i

Table 2.1: Hintikka's Combined Logic of Knowledge and Belief

We omit his defense of positive knowledge introspection, but point out that his defense of undeniable knowledge depends on his interpretation of the operators. Hintikka is concerned in his project with identifying statements about knowledge and belief that utterers cannot defensibly deny. Furthermore, his interpretation of knowledge and belief are not those propositions actively held in one's head, but those to which one could be led through a series of applications of *Modus Ponens* from what one already knows or believes. We can think of knowledge, for Hintikka, as those propositions implied by what one already knows, and beliefs as those propositions that one is logically committed to by one's beliefs. This avoids the problem of logical omniscience for both modal operators. And it sets up an argument like the following,

1. φ follows logically from what i knows
2. If (1) is pointed out to i , it is indefensible for i to deny that she believes that she knows φ .

3. Therefore, $\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$

For Hintikka, this suffices to accept the formula as valid, as well as Knowledge implies Belief, for which a similar argument can be presented.

Conversely, Hintikka rejects (3) because he does not find it plausible that $\mathbf{B}_i \varphi \Rightarrow \mathbf{K}_i \langle \mathbf{K}_i \rangle \varphi$, which (3) and the Truth Axiom jointly imply. He furthermore rejects arguments from introspection, which might ground a defense in considered formula (3), as fallacious.

Another potential combined logic of knowledge and belief is due to Kraus and Lehmann in [36].

$\mathbf{K}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \psi)$	Distribution of \mathbf{K}_i
$\mathbf{K}_i \varphi \Rightarrow \varphi$	Truth Axiom
$\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi$	Positive Knowledge Introspection
$\neg \mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \neg \mathbf{K}_i \varphi$	Negative Knowledge Introspection
$\mathbf{B}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \psi)$	Distribution of \mathbf{B}_i
$\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$	Belief Consistency Axiom
$\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi$	Knowledge implies Belief
$\mathbf{B}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{B}_i \varphi$	Known Belief
From $\vdash \varphi$ and $\vdash \varphi \Rightarrow \psi$, infer $\vdash \psi$	Modus Ponens
From $\vdash \varphi$, infer $\vdash \mathbf{K}_i \varphi$	Necessitation of \mathbf{K}_i

Table 2.2: Kraus-Lehmann Combined Logic of Knowledge and Belief

Knowledge is $\mathcal{S5}$ and belief is based on a serial relation. However, from this system the remaining $KD45$ formulas are derivable. Additionally, they have the interesting theorem $\mathbf{K}_i \varphi \equiv \mathbf{B}_i \mathbf{K}_i \varphi$.

Kraus and Lehmann identify the formula $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$ as interesting but exclude it, as it would cause a collapse of knowledge and belief into equivalence: $\mathbf{K}_i \varphi \equiv \mathbf{B}_i \varphi$. Obviously, a logic combining knowledge and belief should treat them as distinct operators, and we treat the avoidance of this collapse as a desirable feature of our own static base for \mathcal{DASL} . We also include the interesting formula above as an axiom, which imposes a condition on agents' beliefs that we consider a desirable

idealization. To say i believes φ only if she believes that she knows φ is to say she believes only those things that she has very good evidence for. This helps prevent agents from believing things without evidence, which we want to exclude from our model. The challenge with this axiom, as Kraus and Lehmann note, is to avoid the knowledge-belief collapse. We give this further treatment in a later chapter.

This concludes our discussion of static epistemic and doxastic logics, including their combination. We pick this thread back up later when we present the static base of \mathcal{DASL} . The next section presents the relevant background of dynamic extensions to epistemic logic.

2.3 Dynamic Epistemic Logic

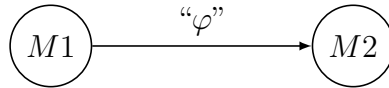
This section presents the so-called *dynamic turn* in modal logic, which incorporates elements of propositional dynamic logic into epistemic logic in order to model knowledge and action. See van Benthem [54, 55, 33], and van Ditmarsch, van der Hoek, and Kooi [62] for detailed examinations. The related Public Announcement Logic (PAL) uses similar techniques to model information flow among agents. See Baltag, Moss, and Solecki in [51] for more on this.

Dynamic Epistemic Logic (DEL) formalizes situations in which agents' epistemic states change over time, due to announcements or other informational events[62]. For example, if Alice truthfully and trustworthily communicates to Bob that φ , then after this informational event it is true that Bob knows φ . This situation cannot be modeled by the epistemic logic introduced in the previous section. To model it, we introduce the following formal machinery.

To capture informational events, we introduce the idea of relativizing a Kripke structure. In the previous example, if we model the Alice and Bob situation prior to Alice's communication, we can have a world w from which Bob considers φ - as

well as $\neg\varphi$ -worlds possible. However, after the informational event, Bob knows φ , so the model is *relativized* to a submodel in which only φ -worlds are accessible by Bob's epistemic possibility relation. Thus, after the informational event, the model transitions to a submodel with fewer edge relations.

The logic for reasoning about information flow in knowledge games is called Dynamic Epistemic Logic (DEL). As its name suggests, it combines elements of epistemic logic and dynamic logic. Epistemic logic is the static logic for reasoning about knowledge, and dynamic logic is used to reason about actions. In dynamic logic semantics, nodes are states of the system or the world, and relations on nodes are transitions via programs or actions from node to node. If we think of each node in dynamic logic as being a model of epistemic logic, then actions become relations on models, representing transitions from one multi-agent epistemic model to another. For example, if we have a static epistemic model $M1$ representing the knowledge states of agents Alice and Bob at a moment, then the action “ φ ” is a relation between $M1$ and $M2$, a new static epistemic model of Alice's and Bob's knowledge after the question is asked. All of this is captured by DEL.



The above figure illustrates the relationship between static epistemic models and dynamic logic models. As a purely dynamic model, the figure shows the action “ $p?$ ” transitioning between nodes $M1$ and $M2$. If we zoom in on the nodes, we see their structure as epistemic models, with their own nodes and edges, representing possible worlds and epistemic relations.

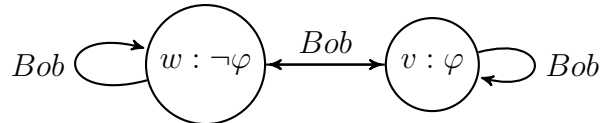


Figure 2.3: $M1$: The model before Alice announces “ φ ”.



Figure 2.4: $M2$: The model after Alice announces “ φ ”.

We are concerned with an additional element: the *safety* status of an action, and an agent’s knowledge and belief about that. To capture this, we extend DEL and call the new logic Dynamic Agent Safety Logic (\mathcal{DASL}), which we introduce in the next chapter. The next section lays out the state of formal methods involving human-machine systems.

2.4 Formal Methods Tools

This section describes the Coq Proof Assistant [13] and that this dissertation makes use of. We use Coq to mechanically check our metatheory proofs of soundness and completeness, and our object theory proofs about safety-critical information flow. Insofar as this thesis constitutes a work of logic or formal philosophy, we hope it can serve as an example for how to use formal tools to increase the rigor of research in these areas.

The Coq Proof Assistant is a combination of a dependently-typed functional programming language and a language of tactics for partially automating proof verification. Coq supports compilation to Haskell and OCaml, which is convenient because one can write a program and verify its correctness in Coq, then compile the code to a nicer implementation language. This dissertation does not take advantage of that feature. Instead, we embed \mathcal{DASL} in Coq and use it as a proof checker that allows us to write recursive functions and proof-automating tactics. This represents a use case of formal verification tools that one does not often see in the philosophy, math-

ematics, or logic departments of academia. Coq is a tool for us to use to increase the rigor of our research. Many who have ventured into the realm of mechanical theorem proving have been surprised initially by how complicated proofs ‘by routine induction’ can become once all of the details must be spelled out, and indeed one can discover embarrassing mistakes when forced to spell out each detail.

The field of formal methods grew out of attempts to verify the correctness of hardware and software. It was normally confined to systems deemed mission- or safety-critical, due to its high cost. Arguments for the use of formal methods normally involve an anecdote in which a system fails due to some sort of error missed by traditional testing methods but that would have been caught and prevented were formal methods used.

Formal methods increase confidence that hardware and software components function correctly [30]. They involve the application of mathematical techniques to the design and analysis of systems, usually the safety-critical components [26]. A standard approach is to develop an abstract specification of the software or hardware system, design the system to meet those specifications, and then use formal logic to prove that the designed system meets the desired safety specifications.

Recent concerns with AI safety from Yudkowsky [83] and Bostrom [12] in particular highlight the need for formal methods that address correctness of AI software. Among the many challenges associated with this endeavor, one is the ability to formally model the AI agents themselves. By hypothesis, checking the correctness of these systems after they have been built is insufficient, so a formal language for specifying correctness properties is also needed. We do not pretend to provide a solution to these problems in this thesis, but we take modest steps in that direction.

By carefully choosing our axiom schemas, we construct a model of a human-like reasoner, and we show that we can specify formal properties about these reasoners. We do so with mechanization in Coq, which future work might be able to lever-

age through compilation to Haskell, OCaml, or some other functional language that preserves correctness proved in Coq.

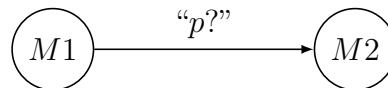
Additionally, the use of formal tools like Coq to model human-like agency offers opportunities for philosophers to formalize and explore the implications of their theories. Sometimes merely attempting to formalize a theory and failing will be a valuable contribution to say, epistemology. Rigorously formalizing arguments and checking their validity could be another use of this technology, and it has been used for example by Rushby in [48], where he mechanizes Anselm’s ontological argument for the existence of God, and in doing so reveals some specious aspects. We hope our chapter on the philosophical considerations underlying epistemic logics contribute to the further adoption of formal tools for philosophical analysis.

This concludes our discussion on formal methods as they relate to our thesis. In the next chapter, we present *DASL* and argue for its adoption over other epistemic logics for reasoning about human-like agents.

Chapter 3

Dynamic Agent Safety Logic

The logic for reasoning about information flow in knowledge games is called Dynamic Epistemic Logic (DEL). As its name suggests, it combines elements of epistemic logic and dynamic logic. Epistemic logic is the static logic for reasoning about knowledge, and dynamic logic is used to reason about actions. In dynamic logic semantics, nodes are states of the system (or of the world), and relations on nodes are transitions via programs or actions from node to node. If we think of each node in dynamic logic as being a model of epistemic logic, then actions become relations on models, representing transitions from one multi-agent epistemic model to another. For example, if we have a static epistemic model $M1$ representing the knowledge states of agents 1 and 2 at a moment, then the action “ $p?$ ” is a relation between $M1$ and $M2$, a new static epistemic model of 1’s and 2’s knowledge after the question is asked. All of this is captured by DEL.



The above figure illustrates the relationship between static epistemic models and dynamic logic models. As a purely dynamic model, the figure shows the action “ $p?$ ”

transitioning between nodes $M1$ and $M2$. If we were to zoom in on the nodes, we would see their structure as epistemic models, with their own nodes and edges, representing possible worlds and epistemic relations.

We are concerned with an additional element: the *safety* status of an action, and an agent's knowledge and belief about that. To capture this, we extend DEL and call the new logic Dynamic Agent Safety Logic (\mathcal{DASL}). The remainder of this section presents \mathcal{DASL} 's syntax, semantics, and proves its soundness.

3.1 Syntax and Semantics

3.1.1 Syntax

\mathcal{DASL} has the following syntax.

$$\varphi ::= p \mid \neg\varphi \mid \varphi \wedge \varphi \mid \mathbf{K}_i \varphi \mid \mathbf{B}_i \varphi \mid [\mathbf{A}, a]_i \varphi \mid [\mathbf{A}, a]_i^S \varphi,$$

where $p \in AtProp$ is an atomic proposition letter drawn from a finite set of such letters, i refers to $i \in Agents$, a is the name of an action, called an action token, belong to a set of such tokens, $Actions$, and \mathbf{A} refers to an action model. The knowledge operator \mathbf{K}_i indicates that “agent i knows that ...” Similarly, the operator for belief, \mathbf{B}_i can be read, “agent i believes that...” The notion of action tokens and structures will be defined in the semantics. The operators $[\mathbf{A}, a]_i$ and $[\mathbf{A}, a]_i^S$ are the dynamic operators for action a from action structure A occurring in the former case, and happening safely in the latter case. One can read the action operators as “after a from A occurs (*safely*), φ holds in all resulting worlds.” We define the dual modal operators $\langle \mathbf{K}_i \rangle$, $\langle \mathbf{B}_i \rangle$, $\langle \mathbf{A}, a \rangle_i$, and $\langle \mathbf{A}, a \rangle_i^S$ in the usual way.

The semantics of \mathcal{DASL} involve two structures that are defined simultaneously, one for epistemic models, and one for action structures capturing the transition re-

lation among epistemic models. Additionally, we define numerous helper functions that straddle the division between metalanguage and object language.

3.1.2 Metalanguage

The metalanguage defining \mathcal{DASL} consists of traditional Kripke models, and action models, with functions for action pre- and post-conditions and the product update of Kripke and action models [62].

Kripke (Relational) Structure

A Kripke structure, called model M , is a tuple $\langle W, R_k^i, R_b^i, w, V \rangle$. It is a set of worlds W , epistemic and doxastic relations on worlds for agents, a world denoting the actual world, and a valuation function V mapping atomic propositions to the set of worlds satisfying them.

Action Model

An action structure \mathbf{A} is a tuple $\langle Actions, \chi_k^i, \chi_b^i, pre, post \rangle$. It is a set of action tokens, sets of epistemic and doxastic relations on action tokens for agents, action pre- and post-condition functions, and safe action pre-condition functions.

An action model captures the agents' subjective perspectives of an event's occurrence. For example, consider a situation in which Alice flips a coin and Bob calls whether it is *heads* or *tails*. The action model imposes no precondition for occurring, represented by the value \top , and each action token's post-condition maps \top to *heads* or *tails* respectively. Before Alice peeks to see what it is, the Action Model \mathbf{A} looks like (omitting the safety pre-conditions and doxastic relations, for simplicity):

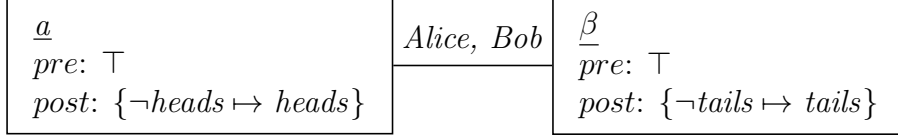


Figure 3.1: Action Model for “Alice flips a coin”

We use rectangles for action tokens to distinguish them from possible worlds in Kripke Structures. Each action token is marked by a token name, a , β , etc. The labels $Alice$, Bob indicate the agents’ epistemic (—) relations. When the relations lack direction, it indicates bi-directionality, and every action token has a looping epistemic relation to itself that we omit in order to keep the pictures simpler. In Figure 3.1, neither Alice nor Bob sees whether the coin lands *heads* or *tails*, but the post-condition of the action is that the proposition expressing the state of the coin becomes *heads* or *tails*. Note that Action Models themselves do not have propositions true and false at them, but rather have pre-conditions and post-conditions for execution of their tokens, and modal relations over the token pairs. The action model for “Alice peeks” is:

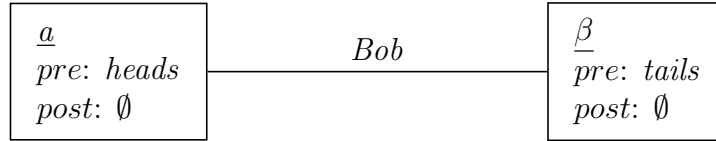


Figure 3.2: Action Model for “Alice peeks”

In Figure 3.2, the epistemic and doxastic relations for Alice have vanished. In the event that Alice peeks and sees *heads*, she sees that the coin landed *heads*, and similarly for *tails*. Bob, however, still considers it possible that Alice peeked and saw *heads* or that she peeked and saw *tails*. However, Bob knows that Alice knows the coin’s state. To represent Alice surreptitiously peeking, we use the following action model:

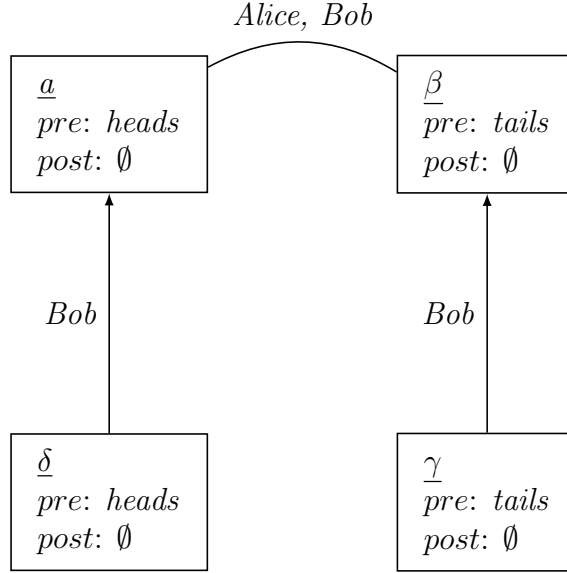


Figure 3.3: Action Model for “Alice secretly peeks”

In Figure 3.3, Alice knows the state of the coin, Bob has the false impression that Alice is, like him, totally ignorant about the coin’s state, and Alice knows this about him. But again, the propositions are not true or false in the Action Model itself; it just describes the pre- and post-conditions and the subjective perspectives of the action. Once the Action Model is used during the application of the *Update Function* (described below), in conjunction with a Kripke Structure, it produces worlds with true and false propositions.

Model Relation

Just as R_k^i denotes a relation on worlds, $\llbracket (A, a)_i \rrbracket$ denotes a relation on Kripke model-world pairs. It represents the relation that holds between M, w and M', w' when agent i engages in action (A, a) at M, w and causes the world to transition to M', w' .

Precondition Function

The Precondition function, $pre :: Actions \mapsto \varphi$, maps an action to the formula capturing the conditions under which the action can occur. For example, if we assume agents tell the truth, then an announcement action has as a precondition that the announced proposition is true, as with regular Public Announcement Logic [21]. Compare this function with the weakest precondition from Hoare logic [25]. It returns a formula that must be true prior to an action's execution, just as the weakest precondition of Hoare logic is a formula that must be true prior to a program's execution. We impose the condition on the precondition function that it output formulas that are composed entirely of literals in disjunctive normal form.

Postcondition Function

The Postcondition function, $post :: Actions \times \varphi \mapsto (AtProp \mapsto \mathcal{P}(W))$, takes an action token and a formula in disjunctive normal form, and maps to a new valuation function from an atomic proposition to sets of worlds satisfying the proposition. It is similarly a partial function, undefined on worlds where the input DNF formula is false.

This functions is a substitution function, as in the logic of DEL with assignment [53]. A similarity is again seen in the strongest postcondition of Hoare logic. If a consequence of a coin flip action is that $h \equiv \text{the coin shows heads}$, then first the precondition function filters out worlds at which the coin flip is not possible, and then the postcondition function assigns all remaining worlds to $V(h)$.

Update Function

The Update function, $update :: (Model \times ActionStruct \times W \times Actions \times Agents) \mapsto (Model \times W)$, takes a Kripke model M , an action structure A , a world from the Kripke

model, an action token from the action structure, and an agent i from $Agents$, and returns a new Kripke model-world pair. It represents the effect actions have on models, and is more complicated than other DEL semantics in that actions can change the facts on the ground in addition to the knowledge and belief relations. It is a partial function that is defined iff a model-world pair satisfies the action's preconditions.

$update(M, A, w, a, i) = (M', w')$ where :

1. $M = \langle W, \{R_k^i\}, \{R_b^i\}, w, V \rangle$
2. $A = \langle Actions, \{\chi_k^i\}, \{\chi_b^i\}, a, pre, post \rangle$
3. $M' = \langle W', \{R_k'^i\}, \{R_b'^i\}, w', V' \rangle$
4. $W' = \{(w, a) | w \in W, a \in Actions, \text{ and } w \models pre(a)\}$
5. $R_k'^i = \{((w, a), (v, b)) | wR_k^i v \text{ and } a\chi_k^i b\}$
6. $R_b'^i = \{((w, a), (v, b)) | wR_b^i v \text{ and } a\chi_b^i b\}$
7. $w' = (w, a)$
8. $V'(p) = post(a, pre(a))(p)$

Safety Precondition Function

The Safety Precondition Function, $pre_s :: Actions \mapsto \varphi$, is a more restrictive function than pre . Where pre returns the conditions that dictate whether the action is possible, pre_s returns the conditions that dictate whether the action is safely permissible. This function is the key reason the dynamic approach allows for easy inference from action to safety-critical information. It represents an innovation in dynamic logic that this thesis introduces, because it introduces modality to the transition relation among models. Just as in static modal logics two worlds can be related by an epistemic relation but not by a doxastic relation, one static model can dynamically transition to another via a *mere* action even if it cannot do so *safely*. One could extend the analogy to Hoare logic and introduce a weakest secure precondition in addition to a

weakest precondition. The weakest secure precondition is stronger than the weakest precondition, because it imposes more restrictions on the state of the system prior to a program's execution.

3.1.3 Semantics

\mathcal{DASL} has the following relational semantics.

$M, w \models p \Leftrightarrow w \in V(p)$
$M, w \models \neg\varphi \Leftrightarrow M, w \not\models \varphi$
$M, w \models \varphi \wedge \psi \Leftrightarrow M, w \models \varphi \text{ and } M, w \models \psi$
$M, w \models \mathbf{K_i} \varphi \Leftrightarrow \forall v, w R_k^i v \text{ implies } M, v \models \varphi$
$M, w \models \mathbf{B_i} \varphi \Leftrightarrow \forall v, w R_b^i v \text{ implies } M, v \models \varphi$
$M, w \models [\mathbf{A}, a]_i \varphi \Leftrightarrow \forall M', w', (M, w) \llbracket (A, a)_i \rrbracket (M', w')$ $\text{implies } M', w' \models \varphi$
$M, w \models [\mathbf{A}, a]_i^S \varphi \Leftrightarrow \forall M', w', (M, w) \llbracket (A, a)_i \rrbracket^S (M', w')$ $\text{implies } M', w' \models \varphi$

Table 3.1: \mathcal{DASL} semantics

The definitions of the dynamic modalities make use of a relation between two model-world pairs, which we now define.

$(M, w) \llbracket (A, a)_i \rrbracket (M', w') \Leftrightarrow M, w \models \text{pre}(a)$ $\text{and } \text{update}(M, A, w, a, i) = (M', w')$
$(M, w) \llbracket (A, a)_i \rrbracket^S (M', w') \Leftrightarrow M, w \models \text{pre}_s(a)$ $\text{and } \text{update}(M, A, w, a, i) = (M', w')$

Table 3.2: \mathcal{DASL} 's dynamic model relations.

3.1.4 Hilbert System

\mathcal{DASL} is axiomatized by the following Hilbert system.

$\mathbf{K}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \psi)$	Distribution of \mathbf{K}_i
$\mathbf{K}_i \varphi \Rightarrow \varphi$	Truth Axiom
$\mathbf{B}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \psi)$	Distribution of \mathbf{B}_i
$\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$	Belief Consistency Axiom
$\neg \mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \neg \mathbf{B}_i \varphi$	Negative Belief Introspection
$\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi$	Knowledge implies Belief
$\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$	Evidential Restraint
$[\mathbf{A}, a]_i p \Leftrightarrow post(a, pre(a))(p)$	Atomic Consequence
$[\mathbf{A}, a]_i \neg \varphi \Leftrightarrow (pre(a) \Rightarrow \neg [\mathbf{A}, a]_i \varphi)$	Action Negation
$[\mathbf{A}, a]_i (\varphi \wedge \psi) \Leftrightarrow ([\mathbf{A}, a]_i \varphi \wedge [\mathbf{A}, a]_i \psi)$	Action Conjunction
$[\mathbf{A}, a]_i \mathbf{K}_i \varphi \Leftrightarrow (pre(a) \Rightarrow \bigwedge_{a\chi_k^i \beta} \mathbf{K}_i [\mathbf{A}, \beta]_i \varphi)$	Action and Knowledge
$[\mathbf{A}, a]_i \mathbf{B}_i \varphi \Leftrightarrow (pre(a) \Rightarrow \bigwedge_{a\chi_b^i \beta} \mathbf{B}_i [\mathbf{A}, \beta]_i \varphi)$	Action and Belief
$[\mathbf{A}, a]_i^S p \Leftrightarrow post(a, pre_s(a))(p)$	Safe Atomic Consequence
$[\mathbf{A}, a]_i^S \neg \varphi \Leftrightarrow (pre_s(a) \Rightarrow \neg [\mathbf{A}, a]_i^S \varphi)$	Safe Action Negation
$[\mathbf{A}, a]_i^S (\varphi \wedge \psi) \Leftrightarrow ([\mathbf{A}, a]_i^S \varphi \wedge [\mathbf{A}, a]_i^S \psi)$	Safe Action Conjunction
$[\mathbf{A}, a]_i^S \mathbf{K}_i \varphi \Leftrightarrow (pre_s(a) \Rightarrow \bigwedge_{a\chi_k^i \beta} \mathbf{K}_i [\mathbf{A}, \beta]_i^S \varphi)$	Safe Action and Knowledge
$[\mathbf{A}, a]_i^S \mathbf{B}_i \varphi \Leftrightarrow (pre_s(a) \Rightarrow \bigwedge_{a\chi_b^i \beta} \mathbf{B}_i [\mathbf{A}, \beta]_i^S \varphi)$	Safe Action and Belief
$[\mathbf{A}, a]_i \varphi \Rightarrow [\mathbf{A}, a]_i^S \varphi$	Inevitability
$\langle \mathbf{A}, a \rangle_i \varphi \Rightarrow \mathbf{B}_i \langle \mathbf{A}, a \rangle_i^S \varphi$	Minimum Rationality
From $\vdash \varphi$ and $\vdash \varphi \Rightarrow \psi$, infer $\vdash \psi$	Modus Ponens
From $\vdash \varphi$, infer $\vdash \mathbf{K}_i \varphi$	Necessitation of \mathbf{K}_i
From $\vdash \varphi$, infer $\vdash [\mathbf{A}, a]_i \varphi$	Necessitation of $[\mathbf{A}, a]_i$

Table 3.3: Hilbert System of \mathcal{DASL}

Above are the axioms characterizing the logic, including the reduction axioms translating formulas with dynamic modalities into purely static formulas.

The reduction axioms are recursively defined on φ , terminating with propositional atoms. When a dynamic operator, either $[\mathbf{A}, a]_i$ or $[\mathbf{A}, a]_i^S$, applied to a propositional atom p , is translated to a static formula, first the pre (or pre_s) function is called on the action a , which if defined, is then passed through the $post$ function to substitute the atoms that change as a result of the action. If this translation process results in

an implication that p , then the dynamic formula is true; otherwise, false.

We proceed by defending the axiom schema for the static foundation, then continue with the dynamic extensions.

3.1.5 Static Base

Like other dynamic logics for knowledge and belief, \mathcal{DASL} consists of a static base of axiom schemas for the knowledge and belief operators. This section explains and defends \mathcal{DASL} 's static base, as they represent a novel axiomatization of formal epistemology and a contribution to the literature by this thesis.

The standard axiom schema for the knowledge operator, whether for a wholly static logic or for the static base of a dynamic system, is that of the $\mathcal{S5}$ operator, described earlier in section 2.2.2. This axiom schema for the knowledge operator's role in classical economics was first formalized by Aumann in [73]. In this model, the belief operator is not necessary, because agents do not form false beliefs. They either know a proposition, or they know that they do not know it. This is not the case for humans in most epistemic situations.

The static base of \mathcal{DASL} logic attends to the formalization of agents who can make decisions based on false beliefs, but whose beliefs are all justified. It is a weaker model than $\mathcal{S5}$, while still making an idealization about its agents.

Recall the static base of \mathcal{DASL} .

$\mathbf{K}_i(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K}_i\varphi \Rightarrow \mathbf{K}_i\psi)$	Distribution of \mathbf{K}_i
$\mathbf{K}_i\varphi \Rightarrow \varphi$	Truth Axiom
$\mathbf{B}_i(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{B}_i\varphi \Rightarrow \mathbf{B}_i\psi)$	Distribution of \mathbf{B}_i
$\mathbf{B}_i\varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$	Belief Consistency Axiom
$\neg \mathbf{B}_i\varphi \Rightarrow \mathbf{B}_i\neg \mathbf{B}_i\varphi$	Negative Belief Introspection
$\mathbf{K}_i\varphi \Rightarrow \mathbf{B}_i\varphi$	Knowledge implies Belief
$\mathbf{B}_i\varphi \Rightarrow \mathbf{B}_i\mathbf{K}_i\varphi$	Evidential Restraint
From $\vdash \varphi$ and $\vdash \varphi \Rightarrow \psi$, infer $\vdash \psi$	Modus Ponens
From $\vdash \varphi$, infer $\vdash \mathbf{K}_i\varphi$	Necessitation of \mathbf{K}_i

Table 3.4: Static Base of \mathcal{DASL}

As far as we can tell, this axiom scheme is novel, and disagrees with other attempts to model the interaction of knowledge and belief. The primary point of departure is that the knowledge operator is no longer an $\mathcal{S5}$ operator, but merely a T operator. Thus, it is severely weakened as a modal operator, but the corollary to this fact is that knowledge no longer has the imposing necessary conditions of positive and negative introspection.

This section presents the desirable features that a static logic should have, and along the way argues against $\mathcal{S5}$'s appropriateness as an epistemic base. Before beginning, we briefly establish a consistent informal interpretation of knowledge and belief.

Defense of $\mathcal{S5}$ for knowledge rests on an interpretation of the knowledge operators not as capturing what agents *actively* know at a moment, but rather what their knowledge allows them to potentially infer. Thus, the Distribution axiom is not a license for logical omniscience at a moment, but rather describes what is logically inferable from what is known by an agent. Similarly for positive and negative introspection. However, this defense does not quite provide positive reasons for accepting the $\mathcal{S5}$ axiom schemas. We present arguments against the positive and negative introspection axiom schemas that show them to be inappropriate even on an *inferable*

or *subjunctive* interpretation of knowledge, where an agent knows φ if and only if her present knowledge allows her to come to know φ after sufficient reflection or inference power.

Doxastic logics similarly tend to adopt a subjunctive interpretation, where an agent believes φ if and only if her current beliefs would allow her eventually infer φ . Both of these interpretations come straight out of Hintikka in [59]. These explications suffer an obvious problem of circularity that we do not wish to suffer here. We wish to start from a more fundamental concept, and use it to build the notions of knowledge and belief. The more fundamental concept we use is *evidence*.

We shall say that an agent knows that φ when all of her objectively observable evidence entails φ . Thus, we are still friendly to the subjunctive interpretations, because it reduces to an inferrability condition. However, we lack the problem of circularity. Similarly, for belief. We shall say that an agent believes that φ when all of her subjectively available evidence entails that φ . We shall spell out the details to these components below.

The formal theory of epistemology espoused by this dissertation departs from those typical of the field. We reject $\mathcal{S}5$, we require the mixing of knowledge and belief operators, and we depend on more fundamental notions of evidence. Now we proceed to lay out what we take to be the desirable features of any epistemic logic.

3.1.6 Desiderata of a Static Base

In constructing our static base, we must identify the desirable features to be achieved. A model of knowledge and belief for human-like agents can be assessed by how well it satisfies these desirable features.

First, the model should be well-balanced in the trade off of normativity *vs.* realism so that it can be used to analyze real-world cases as well as idealized formal cases. It scores well on realism to the extent that it represents the reasoning of human-like

agents in environments of relative messiness, like the real world. It scores well on normativity to the extent that its departures from realism are improvements over reality to be strived for. Not all idealizations in models are normative in character, as we shall see.

Second, the model should make sense from an epistemological standpoint, *viz.* the philosophical considerations of epistemologists. To construct a system of epistemic logic is to construct a formal theory of epistemology for *some* type of agent, whether it is a perfect deductive reasoner like a computer over a relatively small database ontology or *homo economicus* in well-structured strategic situations, or an imperfect reasoner like a human in the real world. So far, the latter type of agent has been largely neglected by formal efforts, because attention has been on properties of formal systems, see below, rather than primarily on their philosophical foundations.¹

Third, the model should satisfy some basic formal criteria. It should be a sound and complete logical system to serve as a static foundation for an extension that treats actions. The logic must avoid a collapse of belief and knowledge into an equivalent operator. The logic must abide by the principle of indistinguishability, where an agent cannot distinguish between possible worlds it cannot rule out. A particularly challenging

1. Normative-Realism Trade-Off A model of rationality is situated in a spectrum trading off realism for normative evaluation and prescription. A totally idealized model like that of game theory or rational choice theory is entirely normative. There is no doubt that the agents being modeled are not anything like human beings. The model presumes normativity by showing that acting like an ideally ra-

¹One significant exception to this claim is the study of bounded rationality, which imposes limitations on the number of iterations of *i* knows that *k* knows that *i* knows . . . We regard this approach as valuable, but taking a different approach than ours. Typically, bounded rationality approaches still assume that knowledge is $\mathcal{S}5$, while we relax this assumption because we regard it is philosophically untenable.

tional agent results in optimal payoffs possible. So, for humans who want to optimize their decisions for payoffs, acting like ideally rational agents is best. However, ideally rational agents have properties that real humans lack, so unless a human can run a simulated agent on a correctly formalized model of a situation, she is unlikely to identify the correct action. After the fact, with a correctly formalized model of the situation, her decision can be analyzed in a post-mortem, but post-mortems never do any good for their subjects.

Theories of bounded rationality achieve more realism, but the trade off of realism and normativity is not clearly balanced. By keeping the $\mathcal{S5}$ knowledge operator and limiting the agents' epistemic depth², theories of bounded rationality seek to model actual human behavior in the rarefied world of games of perfect information, where the agents know the payoff structure and the game tree or game matrix. This is fine for modeling humans in behavioral economics experiments, where they are sitting at a table with a formal description of the game in front of them. But for the messier world of cockpits, emergency rooms, power plants, and automobiles, this assumption of $\mathcal{S5}$ knowledge fails. The difference is between *closed* worlds and *open* worlds. A formal game environment is a closed world, meaning the human is aware of all of the relevant propositions and their truth values, *i.e.*, the game structure. The real world's environment is an open one, where the set of relevant propositions is a mystery, and the truth values of the propositions known to be relevant is a matter of uncertainty.

Likewise, theories of bounded rationality with $\mathcal{S5}$ knowledge do not offer clear normative prescriptions or evaluative guidelines. Rather, it appears to be an attempt to make sense of the following finding of behavioral economics: In some situations, human agents depart from the Nash equilibrium, and they both end up with higher payoffs. The centipede game is a good example. The Nash equilibrium states that player *A* should immediately take the money on the table, whereas experiments show

²*i* knows that *j* know that *i* knows that... depth

that humans regularly let the pot build up until a few steps before the end, one of the players realizes what will happen, and takes the majority of the pot. From the perspective of maximizing expected payoffs, it is hard to criticize either player, as each achieved higher payoff from this outcome than they would have in the Nash equilibrium. A theory of bounded rationality takes this at face value, and seeks to build a model for it. With such a model, the hope is to offer general normative guidance for agents with bounded rationality. However, to examine the particular outcome of a centipede game, supposing player B is the one with the minority payoff, we can truly say that B will regret not having taken the majority pot one step earlier, and properly that he *should* have, because then he would have had a higher payoff, and so this outcome is not in subgame perfect equilibrium. If the game is analyzed as a sequence of decisions, the final decision the minority-payoff-receiving player makes is the wrong one, normatively. Bounded rationality with $S5$ knowledge explains this as bad luck determined by the structure of the game and the first (or second) mover advantage, depending on the players' bounds on epistemic depth. But it is nonetheless true that B could have and should have taken the majority pot one step earlier, and this is not a feature of games with outcomes determined by first (or second) mover advantage, like the game of nim with a 3-4-5 setup (or nim with a 1-3-5-7 setup).

Other attempts to relax the assumptions of classical game theory abandon the assumption of $S5$ knowledge and turn toward probability theory in order to capture uncertainty. Uncertainty of various sorts can be adequately modeled in this fashion, including of the game's structure, and the type of the other agents, *e.g.* whether they are a type of agent capable of reasoning 3 epistemic levels deep or merely 2 levels deep. This is realistic for translation to the real world, and normative in that the probabilistic reasoning must abide by the axioms of probability theory. However, humans are not perfect Bayesian reasoners, so this model is unable to formally represent the mistakes humans frequently make when reasoning probabilistically without

additional adjustments. A realistic model should be able to model mistakes systematically in such a way that valid inferences can be made from mistakes in belief or action.

The model presented here abandons the *S5* knowledge assumption in order to model how beliefs and knowledge logically relate for a new type of agent, only slightly improved over a normal human agent. This agent knows some things, believes more than she knows, and has false beliefs. She is different from normal humans in that she believes a proposition only if she has a good justification to do so, like the observation of evidence that supports the proposition, or because she believes she has made a valid inference to that proposition. The agent is realistic to the extent that a very careful human reasoner can achieve this level some of the time, when it really matters to do so, and this model is normative to the extent that an agent with more true beliefs about the environment will make better decisions after the dynamic extension of actions is layered on top.

Along similar lines but in pursuit of different theoretical interests is the so-called Belief-Desire-Intention (BDI) model of action. This model seeks to explain the modules responsible for actions among humans and presumably other potential agents. The model holds that beliefs about the world, desires for how the world could be, and intentions to realize those desired states, are jointly sufficient and individually necessary conditions for actions. This model is related to the model defended here in that its agents can make mistakes due to false beliefs. However, the model is primarily a theory of action and motivation, not a theory of rational action, which we seek to establish.

2. Philosophical Foundations A formal epistemic theory should strive for adequate philosophical grounding in good epistemology. Hintikka took this very seriously in his 1967 book on *Knowledge and Belief*[59], which serves as a seminal

work in the field of formal epistemology. Hintikka dedicates a great portion of the book to exploring how his formal system handles the intuitive judgments of philosophers regarding ordinary language statements, which was the primary method at the time. An epistemic logic divorced from a philosophical foundation is no longer an epistemic logic for reasoning about human-like knowledge. Some idealized epistemic logics can serve as logics for reasoning about machine knowledge in closed worlds, but as machines come to interact with larger and more complex parts of the world, these idealized logics diminish in adequacy, as computational limitations prevent the machine from implementing even a perfect Bayesian reasoning system. Thus, careful consideration of an epistemic logic’s philosophical foundation is in order.

Epistemologists primarily seek a theory of justification that coheres with our considered intuitions about cases, and a theory of knowledge that can solve the problem of skepticism. A system of formal epistemology should respect these questions by coherently addressing them. We consider it required for a system of epistemic logic to be evaluated as a theory of epistemology. The epistemic logic of $\mathcal{S5}$ with common knowledge does not strive for this. We briefly recount some broad positions of interest from epistemology that our model should take a considered position on.

The Space of Justification Theories Epistemology is the study of knowledge, and knowledge is widely held to require truth and justified beliefs in some way, with perhaps additional conditions being required for a jointly sufficient set of conditions. Our model defines truth through the semantic \models relation between worlds and formulas, and knowledge explicitly requires belief via the axiom *Knowledge implies Belief*. Justification remains for us to elaborate, although we do not model it explicitly with a justification modal operator. A theory of epistemic justification must respond to the infinite regress argument which concludes that justified beliefs are impossible. The infinite regress argument against justification runs as follows:

1. A belief that φ is justified if and only if the belief that φ is justified by other justified beliefs.
2. A justification chain of beliefs either: (a) terminates with unjustified base beliefs, or (b) is circular, or (c) is infinite and non-circular.
3. If (a), then no beliefs are justified.
4. If (b), then no beliefs are justified.
5. If (c), then no beliefs are justified.
6. Therefore, justification is impossible for beliefs.

A theory of justification is shaped by how it responds to this regress argument. Denying premise (1) requires a theory to explain what, other than beliefs, justifies beliefs. Popular candidates include experiential mental states. This involves justification flowing from something that is not true or false to something true or false, and thus is not easily formalized in a logic. Denying premise (3) pushes one toward a *foundationalist* theory, where the task then is to describe how a foundation of basic beliefs can be justified, and how the principle of premise (1) does not properly apply to them. Denying premise (4) pushes one toward *coherentism*, where the task is to justify how a circular chain of beliefs is not viciously circular. Denying premise (5) pushes one toward *infinitism*, where the task is to justify how an incomplete infinite chain can have justified constituent beliefs.

In addition to responding to the regress argument, a theory of epistemic justification must present an analysis of what constitutes justification. In addressing the regress argument, a theory takes a position about the structure of the justification relation and whether beliefs are of a single type or of multiple types. A theory must then take a position about the definition of the justification relation so as to avoid the problem of skepticism.

Skepticism We take it as a basic assumption that human knowledge is possible, and therefore that the epistemic position of general skepticism is false. Therefore, in constructing a formal model of epistemic logic, that model should not lead to general skepticism when applied to humans.

In addition to addressing the infinite regress argument, theories of justification must address an argument from logical closure. The logical closure argument runs as follows:

1. i knows that φ only if i knows \neg *skeptical hypothesis*, where *skeptical hypothesis* is some proposition that entails $\neg\varphi$.
2. i does not know \neg *skeptical hypothesis*.
3. Therefore, i does not know φ .

Denying (1) denies that knowledge is closed under logical entailment, and denying (2) requires care so that the theory will generalize and avoid all potential skeptical scenarios. Denying (1) presents a significant problem for formal epistemologists hoping to create a normal modal logic of knowledge. Common attempts to deny (2) include placing knowledge in a direct justification relation with the real world via causation, or defending a position that some hypotheses can be justifiably ignored from consideration.

A theory of justification is internalist if it holds that the believer is justified in believing some proposition φ if and only if she has awareness or access to the basis for believing φ . Externalist theories deny this, and hold that external facts, *e.g.* the causal process that produced the belief, are what justify beliefs and produce knowledge. Theories of justification sometimes apply to both knowledge and justified belief, and sometimes to just one or the other.

An example internalist view due to Prichard is that, through reflection, one can know that she knows a proposition, and know that she believes but does not know

a proposition [19]. This corresponds to including both introspection axioms for the knowledge operator. A weaker internalist view that we endorse is *evidentialism*, which holds that a belief is justified if and only if the believer has evidence supporting the belief. The task of an *evidentialist* theory is to spell out what evidence is.

An example externalist position is that the process that produced the belief that φ be a *reliable* one, in some sense to be defined by a given *reliabilist* theory. This avoids the problem of skepticism by placing the justification relation between the world and the knowledge of the agent, with or without the agent's awareness of this justification being present. Thus, if φ is true, and i 's belief that φ is produced by a reliable process in the real world, then i knows that φ and therefore knows that the *skeptical hypothesis* is false, because inferring that from their being inconsistent is also a reliable process. This approach clearly rejects both positive and negative introspection axioms, as introspection plays no role in the justification of knowledge.

The problem for an epistemic logician constructing a formal theory is to select an appropriate target for formalization. Most recent work avoids this problem by instead focusing on technical innovations applied to toy problems. We deny ourselves that route, and in doing so make decisions regarding the above problems in a way that we can philosophically defend.

Moore's Paradox As with Hintikka's work, another test of an epistemic logic's adequacy is its ability to correctly handle touchstone cases from the philosophical literature, or better yet, to adequately explain problems that currently lack convincing explanations. For example, the so-called Moore's sentence is a problem pervasive in epistemology and the philosophy of language, attributed to G.E. Moore [43] by Wittgenstein [82]. The standard Moore's paradox to be solved is a statement of the form,

φ , but I don't believe that φ .

The alleged paradox of Moorean sentences is that they are not formal contradictions, but they nonetheless seem like they should be. Hintikka solves this problem by first asserting that his logic is meant for actual communicated statements, not merely sentences, and he distinguishes three new types of implication under study: virtual, epistemic, and doxastic. All implications are established by checking for a model where the antecedent is satisfiable while the consequent is falsifiable. If no such model is possible, then the implication is established. Virtual implication refers to checking for satisfiability of the negated implication in a model. Epistemic implication refers to checking for satisfiability of the negated implication with a knowledge operator distributed over it, and similarly for doxastic implication. Thus, Hintikka solves Moore's paradox by searching for a model that satisfies the above statement. Such a statement is falsifiable, and so virtual implication from φ to $\mathbf{B_i} \varphi$ fails. Hintikka turns then to doxastic implication, which involves distributing $\mathbf{B_i}$ over the implication under inspection:

$$\mathbf{B_i} \varphi \wedge \neg \mathbf{B_i} \mathbf{B_i} \varphi.$$

This is justified by the principle that people normally believe what they assert as statements. Hintikka's doxastic logic has a transitive doxastic relation, so it takes positive introspection about belief as an axiom, and therefore the above formalization is not satisfiable in a model (and so is doxastically implied).

3. Formal Properties A system of logic ought to be sound and complete. An epistemic logic with both belief and knowledge operators must avoid an equivalence collapse of the knowledge and belief operator, which can occur if one is not careful in choosing axioms. However, the logic must satisfy the intuition that for v to be epistemically possible from w , v and w should in some sense be indistinguishable to the agent. Finally, for a logic of sufficient expressive power, namely one with

agents who can reason about their own reasoning, it must avoid a special obstacle to an agent trusting its own conclusion, dubbed the Löbian Obstacle by researchers in the foundations of artificial intelligence. Modal logics with reflective reasoners must address this obstacle. We hold off on discussing this desideratum until Chapter 5.

Soundness and Completeness Soundness and completeness are familiar properties for students of logic. They are briefly mentioned here as a reminder.

Definition 3.1.0.1 (Soundness). *Logic \mathcal{L} is sound if and only if $\models \varphi$ implies $\vdash \varphi$,*

where $\models \varphi$ is semantic validity of φ and $\vdash \varphi$ is axiomatic derivability.

Definition 3.1.0.2 (Completeness). *Logic \mathcal{L} is complete if and only if $\vdash \varphi$ implies $\models \varphi$.*

Avoid KB Collapse Knowledge and belief are distinct. There are two distinguishing characteristics between them. The first is that belief and knowledge sit in a part-whole relationship to each other. i knows that φ only if i believes that φ , with other necessary conditions that are jointly sufficient. The second is that knowledge of a proposition requires that the known proposition be true, while belief requires no such truth condition. If an epistemic logic with operators for belief and knowledge collapses them into equivalence, such that $\mathbf{B}_i \varphi \Leftrightarrow \mathbf{K}_i \varphi$, it has gone wrong, and is not acceptable.

An offending theorem flagged in the literature as responsible for equivalence collapse is $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$, which is our axiom *Evidential Restraint*. We must show that our system avoids **KB** collapse.

Indistinguishability Many find it intuitive to speak of the epistemic relation as an equivalence or indistinguishability relation, which requires it to be symmet-

ric, transitive, and reflexive. The epistemic relation in S5 is Euclidean and reflexive, which together yield symmetry and transitivity, and therefore an equivalence relation. Equivalence relations satisfy an indistinguishability relation because each world in the equivalence class is indistinguishable from each other relative to the relation defining the class. We must show that our axiom schema can satisfy the intuition behind epistemic indistinguishability.

We proceed now to our arguments that *DASL*'s static base fares well on these counts.

3.2 Interpreting K_i and B_i

Epistemic logic can produce some counterintuitive results if the knowledge operator is interpreted inappropriately for the system. For example, the standard Distribution Axiom of knowledge, $K_i(\varphi \Rightarrow \psi) \Rightarrow (K_i\varphi \Rightarrow K_i\psi)$ states that knowledge is closed under implication. This presents several problems when the operator is interpreted as “agent i actively knows right now that φ ”. For empirical knowledge, the threat of global skepticism looms, because presumably i knows that she exists as an embodied agent in the external world only if she is not a mind in a simulated reality. Therefore, she knows she is an embodied agent in the external world only if she knows she is not a mind in a simulated reality. By hypothesis, the skeptic would say, she does not know she is not a mind in a simulated reality. Therefore, she does not know she is an embodied agent in the external world. For statements of pure logic, the Distribution Axiom states that agents know all logical truths.

The standard interpretation in the literature for an epistemic logic with just a knowledge operator is that “agent i *implicitly* knows that φ ,” meaning something like the claim that φ is a live candidate for knowledge for i , should i ever examine the proposition. For doxastic logic, the same concern arises for belief, as it has the

Distribution Axiom as well. The interpretation here is usually something like “agent i ’s beliefs *commit* i to φ ,” or “ i is disposed to believe that φ . In both cases, the interpretation moves away from an active, explicit occurrence of the proposition in the agent’s mind in order to avoid the weirdness of closure. We adopt this interpretation, but further analyze the notions of epistemic and doxastic possibility in order to avoid circularity of our definitions and to avoid the problem of logical omniscience for our agents.

Our interpretation is comparable to a debate that occurs concerning the interpretation of probabilities among frequentists and Bayesians. Frequentists hold that probabilities represent objective frequencies that occur in the world, which can be counted through controlled experiment, and which should directly map to the probability for an event. For Bayesians, probabilities represent subjective uncertainty in the mind of a reasoner. Rather than constructing a model based on one or the other of these notions, we construct a model that includes both the notion of objectively available evidence and subjectively available evidence.

Definition 3.2.0.1 (Objectively Available Evidence). *Evidence is objectively available for proposition φ if and only if some source of information exists in the actual world that supports the proposition φ .*

What is it for a source of information to support a proposition? We take it to be a primitive notion of an event with high mutual information with the proposition, and we leave it at that. Future work can pursue the relationship between information theory and epistemic possibility relations, but we are satisfied to leave things at this level of abstraction. We do, however, provide some brief examples.

Definition 3.2.0.2 (Subjectively Available Evidence). *Evidence is subjectively available to agent i for proposition φ if and only if i has direct mental access to some information, accurate or inaccurate, that supports the proposition φ .*

When a coin lands with *heads* facing up, the light reflecting off the coin's surface is a source of information that supports the propositions that "the coin landed *heads* up," and "the coin is unfair and has two *heads*," and "the coin landed *tails* down." This is because evidence underdetermines the theories it supports, and so multiple propositions are always supported by the same source of information. The light reflecting off the surface is objectively available evidence for these propositions. If anyone notices the light reflecting off the coin, then this evidence becomes subjectively available to them. From the actual world, doxastic relations spider-web outward to possible worlds at which those propositions are true, and she considers those worlds possible. The event itself caused those worlds to become epistemically possible.

Thus, from the actual world, epistemically possible worlds are created by objectively available evidence, and doxastically possible worlds are created by subjectively available evidence. Seeing a coin land *heads* is subjectively available evidence to the agent that sees it, so doxastically possible *heads*-worlds proliferate and possibly overtake all *tails*-worlds, in which case the agent believes that the coin is *heads*.

In the case of false beliefs, suppose Bob is an early bird in 1948 America, and he reads the headlines that *DEWEY DEFEATS TRUMAN*. Here we have objectively available evidence, a source of information in the actual worlds that supports a proposition. One of the propositions supported is that Dewey defeated Truman in the 1948 United States Presidential election. This proposition happens to be false, but it nonetheless has objectively available evidence supporting it. Thus, objectively available evidence does not necessarily guarantee truth for the propositions it supports. It merely lends support to them. When Bob reads the newspaper, he may or may not form the belief that Dewey won the election, depending on the other evidence he has perceived.

Epistemically possible worlds are tied to the actual world in some sense, and doxastically possible worlds are tied to an agent's mental consideration. When objective

evidence is perceived and interpreted, it becomes subjective evidence, and when an agent reflects on the evidence itself, it becomes subjectively objective. This is the intuitive gloss of the model-theoretic notion of composing the doxastic and epistemic relations. When world $wR_b^i v$ and $vR_k^i u$, agent i is reflecting on the objectively available evidence at w , but because it is mediated by her reflection, it must be behind a subjective layer.

Thus, we have provided an interpretation of the possibility relations underlying knowledge and belief, which is not often done in the literature. We now present our theory’s notion of justification and its response to skepticism.

3.2.1 Justification

Justification in our logic is not represented explicitly with a justification operator, as in Artemov and Nogina [1]. There, they formalize justification and knowledge with a formula like $t : \varphi$, where t is a justification (a proof-like object) for φ polynomial in its length. They adopt $\mathcal{S}4$ as the static base epistemic logic, and interpret knowledge in the usual implicit or subjunctive way. They include a connective axiom schema $t : \varphi \Rightarrow \mathbf{K} \varphi$, read “if φ is justified then φ is knowable”. We seek a weaker notion of justification here, under the notion that false propositions can nonetheless be justified by evidence. Rather than representing justification explicitly in the object language, we bury it in the underlying model theory, relating evidence to the possibility relations over worlds. We capture the notion of justification in the object language with the axiom schema *Evidential Restraint*.

Recall *Evidential Restraint*: $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$. An agent believes φ only if *her reflection* on all the objectively available evidence yields only possible worlds where φ is true. It places a constraint on the beliefs agents form such that they have direct awareness of solid evidence for the proposition. Yet they can still be mistaken, because it is their reflection on the evidence that supports the proposition, not the

objectively available evidence itself. Unpacking this, we accomplish several desiderata simultaneously. First, we make a normative idealizing assumption about agents, that all their beliefs are constrained by subjectively available evidence. Second, we capture an internalist theory of justification, where an agent's belief is justified by the evidence available to her³. This theory is philosophically defensible, even if it is not unanimously held (few philosophical positions are). The species of internalism we endorse, evidentialism, can be adopted as a coherentist version or a foundationalist version, and we are happy to adopt the foundationalist version, where perceptual seemings and beliefs about evidence can serve as justifiers of beliefs. Finally, the distinction this theory draws between objective and subject perspectives helps account for powerful intuitions about knowledge.

As a balance between the normative and descriptive, this theory achieves a plausible and desirable trade off. If an agent forms a belief that φ while having defeating evidence directly accessible to her, we are right to criticize her for forming and continuing to hold the belief. This follows the tradition set by Hintikka that the intuitive standard for evaluating formulas of epistemic logic is rational defensibility. On this ground, the axiom is intuitively plausible as a normative claim.

Recall the sentences used to test intuitive validity, called Moorean sentences, of the form, " φ , but I do not know that φ ." The apparent absurdity of such a sentence is belied by the lack of a formal contradiction. Hintikka resolves this issue by defending that claim that when honest agents proclaim sentences, minimally, they believe them to be true. So when such a sentence is uttered, it is equivalent to a formula $\mathbf{B}_i \varphi \wedge \mathbf{B}_i \neg \mathbf{K}_i \varphi$. From this, we can deduce a contradiction if we assume *Evidential Restraint*.

$\mathbf{B}_i \neg \mathbf{K}_i \varphi$ implies $\langle \mathbf{B}_i \rangle \neg \mathbf{K}_i \varphi$, which is equivalent to $\neg \mathbf{B}_i \mathbf{K}_i \varphi$. Thus, we have $\mathbf{B}_i \varphi$ and $\neg \mathbf{B}_i \mathbf{K}_i \varphi$. However, from the former, and the axiom of *Evidential Restraint*,

³The specific internalist theory we adopt is called evidentialism; cf Feldman and Conee in [15].

$\mathbf{B}_i \mathbf{K}_i \varphi$ holds. Thus, we have a contradiction. The fact that the axiom schema helps explain the absurdity of the Moore sentence lends philosophical support to its being a normative constraint on belief. The fact that the constraint is a subjective one, *i.e.* relative to an agent's subjectively available evidence, makes it a realistically achievable idealization, unlike positive and negative introspection of knowledge.

Recall that the problem of logical omniscience is that all propositional tautologies are theorems of the logic, and by the rule of \mathbf{K}_i Necessitation, all theorems are known. The Distribution Axiom of \mathbf{K}_i over implication says if a theorem is known, and i knows that the theorem φ implies ψ , then i knows ψ . Since knowledge implies belief, $\mathbf{B}_i \psi$. Thus, i knows and believes every mathematical theorem. Our solution depends on a function that partitions the set of atomic propositions into an *active* set and an *inactive* set. This activation function constructs a model with atomic propositions that the agent considers relevant. It has the effect of constraining *AtProp* to a finite subset of its original self.

Definition 3.2.0.3 (Active Set). For a model M with set of atomic propositions *AtProp*, the active set *active* $\stackrel{def}{=} \{p \mid p \in \text{AtProp} \text{ and } \text{activate}(p) = \text{True}\}$.

Definition 3.2.0.4 (Inactive Set). For a model M with set of atomic propositions *AtProp*, the inactive set *inactive* $\stackrel{def}{=} \{p \mid p \in \text{AtProp} \text{ and } p \notin \text{active}\}$.

These sets formalize the notion that a human-like agent can actively consider only a limited number of atomic propositions at a time. When an atomic proposition is *inactive*, it will be excluded from the model. Just as we build simple Kripke structures with only a few atomic propositions included, so do humans actually reason this way, with only a few atomic propositions under active consideration at a time.

Now we can ameliorate the problem of skepticism by appealing to our notion of evidence as the source of possible worlds, and the fact that models include only active atomic propositions. For logic and mathematics, evidence consists in proofs.

A proof is a set of sentences realizing the sequence of deductions from axioms and theorems to the proved theorem. A proof is tokenized when it is instantiated in the form of writing, speech, or via mechanization software. We say a tokenized proof constitutes objectively available evidence for a theorem, and perceiving that tokenized proof constitutes subjectively available evidence. Thus, even though there may be a Platonic proof that $P \neq NP$ out there waiting to be instantiated, until it is instantiated, no such objective evidence is available to human-like reasoners. That's why humans truly and without contradiction say that they do not know whether $P=NP$.

Furthermore, we can deny the problematic part of the problem of logical omniscience by denying that agents know all tautological logical implications. They know the ones with tokenized proofs. The ones without tokenized proofs are actually unknown to us. In order to identify an inconsistency, a critic would have to instantiate a proof of a previously unknown theorem, and in so doing, the logic under consideration here would update to include it as objectively available evidence for formulas consisting of atomic propositions that are active. We do not formalize these ideas in this thesis, but leave it for future work to flesh out.

3.3 Indistinguishability Satisfied

Many find it intuitive to speak of the epistemic relation as an equivalence or indistinguishability relation, which requires it be an equivalence relation. Euclidicity, combined with the reflexivity underlying the Truth Axiom and transitivity of positive introspection, yields such an equivalence relation. It makes sense to say that world v is epistemically possible from world w for agent i if i cannot tell the two worlds apart, *i.e.* they are indistinguishable. The power of this intuitive appeal presents a challenge for us. We must show that our axiom schema can satisfy the intuition

behind epistemic indistinguishability without imposing symmetry on the epistemic relation.

The key to achieving this is through the composition of the doxastic relation with the epistemic relation. From a first person perspective, our model does collapse belief and knowledge together. That means that an agent cannot distinguish her mere beliefs from her genuine knowledge. Let $R_{k/b}^i(w)$ denote the extension of the epistemic/doxastic relation of w . We also define the following mathematical notion on the relations.

Definition 3.3.0.1 (Doxastic Capture). *For all $w, v, u \in W$, the doxastic capture of w , $dc(w)$ is defined recursively as follows:*

1. $v \in dc(w)$ iff $wR_b^i v$

We say world w is *path connected* to v if there is some number of relational steps, R_b^i and R_k^i interchangeably, from w to v . This yields the following theorem.

Theorem 3.3.1 (Doxastic Path Capture). *For all $w, v, u \in W$, if $wR_b^i v$ and v is path connected to u , then $v, u \in dc(w)$.*

Proof. $v \in dc(w)$ trivially from the base case of the definition of $dc(w)$. For u , we iteratively simplify the path between v and u in the following way. Since u is path connected to v , there is some number n of relational steps from v to u . If it is by n R_b^i steps, then immediately $vR_b^i u$, since R_b^i is transitive. Then $wR_b^i u$ as well, by transitivity through v .

If there are n steps from v to u interchanging R_b^i and R_k^i , then let $x \in W$ be the first world path connected to v by a R_k^i step, with $m < n$ R_b^i steps between v and x . If $m = 1$, then $w(R_k^i \circ R_b^i)x$, and therefore $wR_b^i x$, yielding $x \in dc(w)$. Otherwise, there are $m - 1$ R_b^i steps between v and x . Examine the world at step $m - 1$, call it v' . We have that $vR_b^i v'$ by transitivity, and therefore $wR_b^i v'$. So $v' \in dc(w)$. Since $v'R_k^i x$, $w(R_k^i \circ R_b^i)x$, and so $wR_b^i x$ and $x \in dc(w)$. This process repeats, simplifying

R_b^i steps by transitivity, and simplifying $(R_k^i \circ R_b^i)$ steps by the subset relation to R_b^i , until reaching u . \square

We also have the following.

Theorem 3.3.2 (Doxastic Capture is Equivalence Relation). *An equivalence relation is reflexive, transitive, and Euclidean. The doxastic capture of w is such a relation for both R_b^i and R_k^i .*

Proof. It suffices to show that R_b^i is reflexive and R_k^i is transitive and Euclidean under $dc(w)$. Suppose $v \in dc(w)$.

R_b^i -Reflexivity. From the supposition, $wR_b^i v$. Since R_b^i is serial, there is some u such that $vR_b^i u$. By transitivity, $wR_b^i u$. By Euclidean, $uR_b^i v$. Since R_k^i is reflexive, $vR_k^i v$. So, through u , $v(R_k^i \circ R_b^i)v$, and therefore $vR_b^i v$.

R_k^i -Transitivity. With $wR_b^i v$, let $x, y \in dc(w)$ be such that $vR_k^i x$ and $xR_k^i y$. From the definition of $dc(w)$, $wR_b^i y$. The Euclidean property yields $vR_b^i y$, and since $R_b^i \subseteq R_k^i$, it follows that $vR_k^i y$.

R_k^i -Euclidean. We must show $vR_k^i x$ and $vR_k^i y$ implies $xR_k^i y$. Assume $v, x, y \in dc(w)$, $vR_k^i x$, and $vR_k^i y$. By the definition of $dc(w)$, $wR_b^i v$ and $wR_b^i y$, so $vR_b^i y$ by the Euclidean property of R_b^i , and by $R_b^i \subseteq R_k^i$, it follows that $vR_k^i y$. \square

Since R_k^i and R_b^i are part of the same bimodal equivalence relation under $dc(w)$, for all $w \in W$, the following theorems follow in a straightforward manner, establishing notions of weak and strong indistinguishability, which we call *BK Indistinguishability*, highlighting the fact that behind the belief operator's doxastic relation, belief and knowledge are indistinguishable from each other, and in a mutual equivalence relation with each other.

Theorem 3.3.3 (Weak BK Indistinguishability). *For all $w, v, w' \in W$, $w(R_k^i \circ R_b^i)v$ and $wR_b^i w'$ implies $R_b^i(w') = R_k^i(w') = R_k^i(v) = R_b^i(v)$.*

Proof. Suppose $w(R_k^i \circ R_b^i)v$ and wR_b^iw' for some arbitrary $w, w', v \in W$. It suffices to show that $w', v \in dc(w)$. Since wR_b^iw' , $w' \in dc(w)$. Since $w(R_k^i \circ R_b^i)v$, the composed relation being a subset of R_b^i , it follows that wR_b^iv and therefore $v \in dc(w)$. \square

Theorem 3.3.3 does not guarantee that the actual world w is in this bimodal composed equivalence relation, because it may not be the case that wR_b^iw . This fails when i has false beliefs, which is frequently the case with humans. However, the following is still the case: any branch departing w via the R_b^i relation will be a bimodal composed equivalence class of epistemic and doxastic possibility, and thus indistinguishable from each other. Any path that departs from w solely by the R_k^i relation will not partake in the equivalence class. The way to interpret this is as follows. For some reason, the epistemically possible worlds outside the bimodal composed equivalence class are not registering with her cognizance. It could be that she is overlooking some information. It could be that she is motivated in her skepticism for some psychological reason. It could be that she is under the cloud of some cognitive bias. Regardless, when φ is true in the bimodal composed equivalence class but false in some epistemically possible world (*e.g.* the actual world), her belief that φ overpowers the lone epistemic relation, and she instead believes that she knows φ .

If the actual world is in its own doxastic relation, we have the following.

Theorem 3.3.4 (Strong BK Indistinguishability). *For all $w, v \in W$, if wR_b^iw then for all $v \in W$, wR_b^iv implies $R_b^i(w) = R_k^i(w) = R_k^i(v) = R_b^i(v)$.*

Proof. Suppose for an arbitrary $w, v \in W$ that wR_b^iw and wR_b^iv . By definition, both are in $dc(w)$, which establishes the bimodal equivalence relation. \square

So, when i has no false beliefs, her belief and knowledge collapse into an $\mathcal{S5}$ modal operator. When this condition holds, she makes decisions on par with those of an ideally rational agent who does not make mistakes. However, she is not guaranteed to remain in this state, as the model (world) might change in such a way that the

doxastic relation comes to exclude the actual world, in which case i will have false beliefs again.

Thus, in the static epistemic-doxastic logic of \mathcal{DASL} , the intuitively appealing property of indistinguishability is maintained from a subjective, internal standpoint, *i.e.*, from i 's own perspective. This should be all that is needed, because the motivating intuition behind the indistinguishability requirement is that an agent should not be able to tell which of her possible worlds she is in, and that is the case here with respect to her subjective evidence.

3.4 Avoiding The Collapse

It is important that we maintain the distinction between knowledge and belief if we aim to model humans. To show that our logic maintains this distinction, we present the following counterexample.

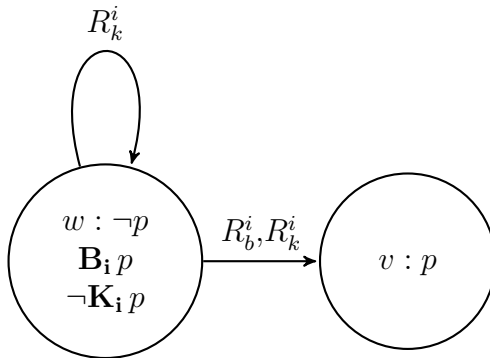


Figure 3.4: A counterexample to $\mathbf{B}_i \varphi \Rightarrow \mathbf{K}_i \varphi$.

The counterexample is based on the key difference between belief and knowledge, namely that knowledge must be true, while belief does not need to be true. Here, at w , the proposition p is not true. However, because i only considers p possible, and neglects the possibility that $\neg p$, she believes i at w . Since p is not true at w , and

knowledge entails truth (via the reflexivity frame condition), it is not true at w that i knows p . Therefore, at w , i believes p but she does not know it.

This suffices to show that the two modalities do not collapse. We have thus far defended the static base in terms of indistinguishability and collapse avoidance. We turn now to soundness and completeness.

3.5 Soundness and Completeness

This section establishes the soundness and completeness of the logic. We present mechanized formalizations of each proof in the Coq Proof Assistant. The proof of soundness shows that axioms of \mathcal{DASL} are each valid for all \mathcal{DASL} frames. The proof of completeness shows that each axiom schema is a Sahlqvist formula, and therefore forms a complete logic built by adding them to the axiom K of the minimal modal logic, by Sahlqvist's Theorem [76].

Both proofs make use of a special property of modal logics. Every modal logic formula is equivalent to a first- or second-order logic formula, in the sense that it is valid on the same frames as the *corresponding* first- or second order formula. The frame conditions of \mathcal{DASL} are first order logic formulas, and they each correspond to an axiom schema of the modal logic. This correspondence suffices for soundness. When the modal formulas correspond to first order formulas, and the formulas have a certain special structure, the resulting logic is guaranteed to be complete with respect to the first order frame relations.

We begin this section with the soundness proof, and follow up with the completeness proof.

3.5.1 Soundness

A logic has the property of soundness *if and only if* all of its derivable theorems are valid in its model theory. \mathcal{DASL} 's model theory is the class of relational frames with two binary relations for each agent i , R_k^i and R_b^i , defined over possible worlds. The R_k^i relations are reflexive, and the R_b^i relations are serial, euclidean, and supersets of R_k^i , and of $(R_k^i \circ R_b^i)$. Any formula satisfied by all models of this class of frames, no matter the valuation assignment of propositional variables, is a valid formula. Thus, to show that \mathcal{DASL} is sound is to show that all of its theorems are satisfied by all such models.

Formally, this amounts to:

Theorem 3.5.1 (Soundness). *Dynamic Agent Safety Logic is sound for Kripke structures with*

- (1) reflexive R_k^i relations,
- (2) serial, Euclidean R_b^i relations,
- (3) which are partially ordered $(R_k^i \circ R_b^i) \subseteq R_b^i$, and $R_b^i \subseteq R_k^i$.

In order to formalize this in Coq for mechanical proof checking, we must define the notions of a frame, model, proposition, and theorem. We follow the Master's thesis of Paulien De Wind [14] for the model theory mechanization and the works of Lescanne and Puisségur [66, 72] and Maliković and Čubrilo [68, 69] in setting up the proof theory.

A frame is a record in Coq that consists of a set of possible worlds and relations on agent-world-world tuples to capture the fact that each binary relation is parameterized by each agent in the model.

```

Record frame : Type := {
  W : Set;
  Rk : DASL.Agents → W → W → Prop;
  Rb : DASL.Agents → W → W → Prop
}.

```

A model is a frame combined with a valuation function assigning worlds to atoms. We likewise create a type of agents here, imported from a library *DASL*, defining them as,

```

Inductive Agents: Type := Pilot | CoPilot | AutoPilot.

```

but any target domain's agents will suffice here.

```

Record model : Type := {
  F : frame;
  Val : (W F) → Atoms → Prop;
  Agents: DASL.Agents
}.

```

We say a proposition is an **Inductive** type of atoms, implications, negations, falsum, and the doxastic and epistemic modal operators.

```

Inductive prop : Type :=
  | atm : Atoms → prop
  | imp : prop → prop → prop
  | negp : prop → prop
  | falsum : prop
  | K : Agents → prop → prop
  | B : Agents → prop → prop.

```

The type of `Atoms` is domain specific for the ground truth facts on which propositions in the system are built. For *DASL* these are instrument readings in the cockpit, but for other domains they will be different.

A theorem type is defined in Coq so as to include all propositional tautologies, the modal axioms at the static base of *DASL*, and closure through Modus Ponens and Necessitation.

```

Inductive theorem : prop → Prop :=
| Hilbert_K: forall p q : prop, theorem (p ⇒ q ⇒ p)
| Hilbert_S: forall p q r : prop,
  theorem ((p ⇒ q ⇒ r) ⇒ (p ⇒ q) ⇒ (p ⇒ r))
| Classic_NOTNOT : forall p : prop, theorem ((NOT (NOT p)) ⇒ p)
| MP : forall p q : prop, theorem (p ⇒ q) → theorem p → theorem q
| K_Nec : forall (a : DASL.Agents) (p : prop),
  theorem p → theorem (K a p)
| K_K : forall (a : DASL.Agents) (p q : prop),
  theorem (K a p ⇒ K a (p ⇒ q) ⇒ K a q)
| K_T : forall (a : DASL.Agents) (p : prop), theorem (K a p ⇒ p)
| B_K : forall (a : DASL.Agents) (p q : prop),
  theorem (B a p ⇒ B a (p ⇒ q) ⇒ B a q)
| B_Serial : forall (a : DASL.Agents) (p : prop),
  theorem (B a p ⇒ NOT (B a (NOT p)))
| B_5 : forall (a : DASL.Agents) (p : prop),
  theorem (NOT (B a p) ⇒ B a (NOT (B a p)))
| K_B : forall (a : DASL.Agents) (p : prop), theorem (K a p ⇒ B a p)
| B_BK : forall (a : DASL.Agents) (p : prop),
  theorem (B a p ⇒ B a (K a p)).

```

We denote the **theorem** type judgment with $| -$. Our mechanical proof strategy is to follow closely what a pen-and-paper proof of soundness consists in. We will define the frame conditions of \mathcal{DASL} 's relations, and we will prove that each is sufficient for deriving that its corresponding modal formula is valid for frames with that condition. An example will illustrate this.

```

Definition reflexive_Rk_frame (F : frame) : Prop :=
  forall (w : (W F)) (ags : DASL.Agents), (Rk F ags w w).

```

First we define the property of epistemic reflexivity on a frame. We then prove a lemma showing it is sufficient for the Truth Axiom's being valid on reflexive frames.

```

Lemma K_is_refl : forall (phi : prop) (F : frame) (a : DASL.Agents),
  (reflexive_Rk_frame F) →
  F || = ((K a phi) ⇒ phi).
Proof.
  intros.
  unfold reflexive_Rk_frame in H.
  unfold Frame_validity.
  intros.
  unfold Model_satisfies.
  intros. pose proof H w; clear H. pose proof H0 a; clear H0.
  unfold satisfies.
  intros. pose proof H0 w; clear H0.
  simpl in H1. pose proof H1 H; clear H1.
  auto.
Qed.

```

We add this lemma to our `hints` so that Coq's engine will automatically try to use it on calls to `auto`.

```

Hint Resolve K_is_refl.

```

We define a *DASL* frame as a simple conjunction over the above properties.


```

Definition DASL_Frame (F : frame) : Prop :=
  reflexive_Rk_frame F ∧
  serial_Rb_frame F ∧
  euclidean_Rb_frame F ∧
  Rb_subset_Rk F ∧
  Rb_subset_Rb_compose_Rk F.

```

Through similar definitions and lemmas of the remaining frame conditions and axiom schemas, the soundness theorem is stated and proven as follows.

```

Theorem DASL_Soundness : forall (phi : prop) (F : frame) (a : DASL.Agents),
  DASL_Frame F →
  |-- phi →
  F || = phi.
Proof.

  intros phi F.
  unfold DASL_Frame.
  intros. destruct H; destruct H1; destruct H2; destruct H3.
  induction H0; eauto.
Qed.

```

The proof proceeds by instantiating `phi`, `F`, `a`, unfolding the definition of `DASL_Frame` into its constituent frame conditions, and running and induction over `|-- phi`. The call to `eauto` suffices due to each of the helper lemmas added to the `hints`. We turn now to mechanizing completeness.

3.5.2 Completeness

For completeness, our proof depends on the notion of a Sahlqvist formula, which itself depends on syntactical properties of modal formulas, which we must also define. The completeness proof is only *mostly* mechanized. The overall theme of the proof is to rely on the *Sahlqvist Completeness Theorem* (SComT) and Sahlqvist Correspondence Theorem (SCorT) of modal logic.

We first explain what a Sahlqvist formula is. A Sahlqvist formula is composed of Sahlqvist implications combined in a manner to be described later.

A *Sahlqvist implication* is a formula with a *Sahlqvist antecedent* and a *positive consequent*.

A formula is *positive* if all of its propositional letters are in the scope of an even number of negation signs. A formula is *negative* if one or more of its propositional letters are in the scope of an odd number of negation signs.

A *Sahlqvist antecedent* is a formula built from a) propositional letters prefixed by $n \geq 0$ \Box 's and b) negative formulas, c) using only \vee , \wedge , and \Diamond to build the antecedent from (a) and (b) components. We can now state the key theorems.

The SCorT is as follows:

Theorem 3.5.2 (Sahlqvist Correspondence Theorem[20]⁴). *If φ is a Sahlqvist formula, then φ corresponds to a first order frame condition χ , and χ is effectively computable from φ via the Sahlqvist-van Benthem algorithm.*

The SComT is:

Theorem 3.5.3 (Sahlqvist Completeness Theorem[20]⁵). *Every Sahlqvist formula is canonical for the first-order property it defines. Hence, given a set of Sahlqvist axioms Σ , the logic $\mathbf{K}\Sigma$ is strongly complete with respect to the class of frames \mathbf{F}_Σ (that is, the first-order class of frames defined by Σ .)*

⁴Theorem 3.54 of Chapter 3 in Blackburn *et. al.*

⁵Theorem 4.42 of Chapter 4 in Blackburn *et. al.*

This states that every modal formula with a syntactic structure, defined by Sahlqvist, is complete for the class of frames defined by its corresponding first order formula. Our completeness proof first mechanically proves that each static axiom schema in *DASL* (excluding those valid for all normal modal logics) is a Sahlqvist formula. We then apply the Sahlqvist-van Benthem Algorithm to each formula to show that the *DASL* frames we defined are correspond to the static axiom schemas.

To mechanize these notions and reason about axiom schemas rather than actual formulas of the object language, we define an inductive type in Coq that more closely resembles the formula schemas one typically uses when discussing a logic, as in $\mathbf{K}_i \varphi \Rightarrow \varphi$, where φ is a metalanguage variable standing in for any well-formed formula of the language. Because our inductive Coq type for well-formed formulas of the language is `prop`, we must create an inductive type that takes `prop` as basic.

```
Inductive schema : Type :=
  | SProp : prop → schema
  | SAnd : schema → schema → schema
  | SOr : schema → schema → schema
  | SImp : schema → schema → schema
  | SNeg : schema → schema
  | SK : DASL.Agents → schema → schema
  | SB : DASL.Agents → schema → schema.
```

We use the following notation symbols.

```
Notation “\ p” := (SNeg p) (at level 70, right associativity ).
Infix “=s⇒” := SImp (right associativity, at level 85).
Infix “| s|” := SOr (right associativity , at level 75).
Infix “&s&” := SAnd (right associativity, at level 75).
```

We define a **negative** formula as one whose proposition letters are all under the scope of an odd number of negation signs.

```
Fixpoint negative_formula (phi : schema) : Prop :=
  match phi with
  | SProp p  $\Rightarrow$  False
  | SAnd phi1 phi2  $\Rightarrow$  (negative_formula phi1)
     $\wedge$  (negative_formula phi2)
  | SOr phi1 phi2  $\Rightarrow$  (negative_formula phi1)
     $\wedge$  (negative_formula phi2)
  | SImp phi1 phi2  $\Rightarrow$  (negative_formula phi1)
     $\wedge$  (negative_formula phi2)
  | SNeg phi'  $\Rightarrow$  not (negative_formula phi')
  | SK a phi'  $\Rightarrow$  negative_formula phi'
  | SB a phi'  $\Rightarrow$  negative_formula phi'
  end.
```

We define a **positive** formula appropriately as one whose proposition letters are all under the scope of an even number of negation signs..

```

Fixpoint positive_formula (phi : schema) : Prop :=
  match phi with
  | SProp p  $\Rightarrow$  True
  | SAnd phi1 phi2  $\Rightarrow$  (positive_formula phi1)
     $\wedge$  (positive_formula phi2)
  | SOr phi1 phi2  $\Rightarrow$  (positive_formula phi1)
     $\wedge$  (positive_formula phi2)
  | SImp phi1 phi2  $\Rightarrow$  (positive_formula phi1)
     $\wedge$  (positive_formula phi2)
  | SNeg phi'  $\Rightarrow$  not (positive_formula phi')
  | SK a phi'  $\Rightarrow$  positive_formula phi'
  | SB a phi'  $\Rightarrow$  positive_formula phi'
  end.

```

To construct Sahlqvist antecedents, we define the notions of a **boxed** formula and a **s_a_component** (for Sahlqvist antecedent component).

```

Fixpoint boxed_formula (phi : schema) : Prop :=
  match phi with
  | SProp p  $\Rightarrow$  True
  | SAnd phi1 phi2  $\Rightarrow$  False
  | SOr phi1 phi2  $\Rightarrow$  False
  | SImp phi1 phi2  $\Rightarrow$  False
  | SNeg phi'  $\Rightarrow$  False
  | SK a phi'  $\Rightarrow$  boxed_formula phi'
  | SB a phi'  $\Rightarrow$  boxed_formula phi'
  end.

```

This states that a `boxed_formula` can be a formula variable on its own, or any number of modal boxes prefixing a formula variable, but any other structure imposed on the formula schema means it is not boxed.

A `s_a_component` is a formula schema built up of `boxed_formula`'s and `negative` formulas connected by conjunction, disjunction, and $\langle \mathbf{K}_i \rangle$, $\langle \mathbf{B}_i \rangle$.

```

Fixpoint s_a_component (phi : schema) : Prop :=
  match phi with
  | SProp p => True
  | SAnd phi1 phi2 => (s_a_component phi1) ∧ (s_a_component phi2)
  | SOr phi1 phi2 => (s_a_component phi1) ∧ (s_a_component phi2)
  | SImp phi1 phi2 => not (s_a_component phi1) ∧ (s_a_component phi2)
  | SNeg phi' => match phi' with
    | SProp p => True
    | SAnd p1 p2 => positive_formula p1 ∧ positive_formula p2
    | SOr p1 p2 => positive_formula p1 ∧ positive_formula p2
    | SImp p1 p2 => negative_formula p1 ∧ positive_formula p2
    | SNeg p' => s_a_component p'
    | SK a p' => not (s_a_component p')
    | SB a p' => not (s_a_component p')
  end
  | SK a phi' => boxed_formula phi'
  | SB a phi' => boxed_formula phi'
end.

```

Next, we define what it is for a subformula to be a Sahlqvist antecedent.

```
Fixpoint sahlqvist_antecedent (phi : schema) : Prop :=  
  s_a_component phi.
```

We combine the above components to define what it is for a formula to be a Sahlqvist implication.

```
Definition sahlqvist_implication (phi psi : schema) : Prop :=  
  sahlqvist_antecedent (phi)  $\wedge$  positive_formula (psi).
```

According to Blackburn *et. al.*[20], a Sahlqvist formula is built up from “Sahlqvist implications by freely applying boxes and conjunctions, and by applying disjunctions only between formulas that do not share any proposition letters.” So we must define a function that parses implications to identify when a proposition letter appears in both antecedent and consequent. We do that with the following in Coq.

```

Fixpoint share_prop_letter (phi psi : schema) {struct phi} : Prop :=
  match phi with
  | SProp phi' => match psi with
    | SProp psi' => phi' = psi'
    | SAnd psi1 psi2 => (prop_in_schema phi' psi1)
      ∨ (prop_in_schema phi' psi2)
    | SOr psi1 psi2 => (prop_in_schema phi' psi1)
      ∨ (prop_in_schema phi' psi2)
    | SImp psi1 psi2 => (prop_in_schema phi' psi1)
      ∨ (prop_in_schema phi' psi2)
    | SNeg psi' => (prop_in_schema phi' psi')
    | SK a psi' => (prop_in_schema phi' psi')
    | SB a psi' => (prop_in_schema phi' psi')
  end
  | SAnd phi1 phi2 => (share_prop_letter phi1 psi)
    ∨ (share_prop_letter phi2 psi)
  | SOr phi1 phi2 => (share_prop_letter phi1 psi)
    ∨ (share_prop_letter phi2 psi)
  | SImp phi1 phi2 => (share_prop_letter phi1 psi)
    ∨ (share_prop_letter phi2 psi)
  | SNeg phi' => (share_prop_letter phi' psi)
  | SK a phi' => (share_prop_letter phi' psi)
  | SB a phi' => (share_prop_letter phi' psi)
end.

```

With these components, we define a Sahlqvist formula as follows.


```

Fixpoint sahlqvist_formula (phi : schema) : Prop :=
  match phi with
  | SProp phi' => True
  | SAnd phi1 phi2 => (sahlqvist_formula phi1)
    ^ (sahlqvist_formula phi2)
  | SOr phi1 phi2 => not (share_prop_letter phi1 phi2)
    ^ (sahlqvist_formula phi1) ^ (sahlqvist_formula phi2)
  | SImp phi1 phi2 => (sahlqvist_implication phi1 phi2)
  | SNeg phi' => (positive_formula phi')
  | SK a phi' => sahlqvist_formula phi'
  | SB a phi' => sahlqvist_formula phi'
  end.

```

For completeness, we must show that each static *DASL* axiom schema is a Sahlqvist formula. The SComT allows us to prove *DASL* is complete with respect to a class of frames corresponding to the axiom schemas, and by the SCorT it holds that this class of frames is effectively computable by the Sahlqvist-van Benthem Algorithm. So, if we prove that each static axiom schema is a Sahlqvist formula, and use the Sahlqvist-van Benthem algorithm to show that the frame conditions for which *DASL* is complete are indeed the *DASL* frame conditions we have defined, it follows that those frame conditions are the ones for which the axioms schemas are complete.

We can represent a theory of modal logic as a list of axiom schemas, and represent *DASL* in particular as the following list.

```

Definition DASL_Axioms (p q r : prop) (a : DASL.Agents) :=
  (SK a (SProp p) =s⇒ SK a ((SProp p) =s⇒ (SProp q)) =s⇒ SK a (SProp q))
:: (SK a (SProp p) =s⇒ (SProp p))
:: (SB a (SProp p) =s⇒ SB a ((SProp p) =s⇒ (SProp q)) =s⇒ SB a (SProp q))
:: (SB a (SProp p) =s⇒ \ (SB a (\ (SProp p))))
:: (\ (SB a (\ (SB a (SProp p)))) =s⇒ SB a (SProp p))
:: (SK a (SProp p) =s⇒ SB a (SProp p))
:: (SB a (SProp p) =s⇒ SB a (SK a (SProp p)))
:: nil.

```

We define a function over lists of axiom schemas to determine whether each schema is a Sahlqvist implication.

```

Fixpoint Complete_via_Sahlqvist (l : list formula) : Prop :=
  match l with
  | nil ⇒ True
  | (l' :: els) ⇒ sahlqvist_formula (l') ∧
    Complete_via_Sahlqvist (els)
  end.

```

Then we prove that \mathcal{DASL} is complete by calling `Complete_via_Sahlqvist` on its representative list of axiom schemas.

```

Theorem DASL_Axioms_Complete : forall (p q r : prop) (a : DASL.Agents),
  Complete_via_Sahlqvist (DASL_Axioms p q r a).

```

We have only mechanized the proof that \mathcal{DASL} is complete with respect to *some* set of frame relations. To show that \mathcal{DASL} is complete with respect to the very

same relations as those for which we proved soundness, we apply the Sahlqvist-van Benthem Algorithm, which takes a starting point .

Our proof continues by applying van Benthem's standard translation algorithm, as presented by Blackburn *et. al.* in [20] in order to yield a second order translation of each \mathcal{DASL} axiom schema, and then applies the Sahlqvist-van Benthem algorithm, presented also by [20], to derive the corresponding first order formulas. The transformation takes us from a modal formula through monadic second order logic, and applies lambda calculus to transform second order formulas into first order formulas. We do not mechanize all of this machinery, but do provide the important details here.

The Sahlqvist-van Benthem Algorithm translates a modal logic formula into a first order frame condition. It does so by first applying the standard translation of a modal formula $ST(\varphi)$ to generate a second order logic formula, and then applying some choice generalizations and instantiations of predicates in order to produce a minimal valuation satisfying the formula's antecedent. Monotonicity results preserve validity from the minimal valuation to all expansions of it. Thus, the result is the smallest first order formula defining a class of frames that can be expanded by any valuation function.

We proceed by describing the Standard Translation as an inductive function, and then applying it to the static \mathcal{DASL} axiom schemas.

The Standard Translation algorithm consists of the following recursive rewrite rules.

Definition 3.5.3.1 (Standard Translation). . *For binary modal operators \mathbf{B}_i and*

\mathbf{K}_i , we define the function ST on a given variable x and proposition letter p :

$$ST(x, p) = P(x)$$

$$ST(x, \perp) = (x \neq x)$$

$$ST(x, \neg\varphi) = \sim (ST(x, \varphi))$$

$$ST(x, \varphi \vee \psi) = ST(x, \varphi) + ST(x, \psi)$$

$$ST(x, \mathbf{K}_i \varphi) = \forall y, R_k^i(x, y) \implies ST(y, \varphi)$$

The resulting predicate and objects corresponding to modal propositions and variables are arbitrary, and so there are infinitely many standard translations from a modal formula to a first order formula.⁶ By the van Benthem Characterization Theorem, each such translation is equivalent to the other. An intuitive interpretation is that the translated formula in classical logic represents the satisfaction predicate for the proposition semantics of the modal logic, with the *current* world w assigned to the free variable x . However, because we quantify over predicates, the resulting formula is second order, and we can instantiate the predicates to minimal valuations. We restrict our quantification over monadic predicates, so that we are limited to monadic second order logic⁷.

The monadic second order standard translation is fed into the Sahlqvist-van Benthem algorithm for producing a corresponding frame condition.

The Sahlqvist-van Benthem algorithm is described in the following way.

Step 1. Apply the following equivalences, from left to right, to the antecedent (in

⁶See chapter 2, example 2.46 in [20].

⁷monadic second order logic is restricted to predicates over sets.

this order):

$$(\exists x, \Phi(x) \cdot \psi) \iff \exists x, (\Phi(x) \cdot \psi) \quad (3.1)$$

$$(\exists x, \Phi(x) \implies \psi) \iff \forall x, (\Phi(x) \implies \psi), \quad (3.2)$$

where $\Phi(x)$ is a sentence with one free variable x .

Step 2. Construct minimal valuations. For each proposition letter instance in the formula's antecedent, select a minimal valuation that makes the antecedent true. The resulting formula will quantify over these variables, so they must be distinct from the free variable in the standard translation. Our intuitive interpretation is that the variables are worlds from the frame. An appropriate lambda function can be created for each predicate such that the smallest frame possible validates it.

For a proposition letter p appearing once in a formula's antecedent, either as an atom or behind a diamond, $\lambda u. u = x$ is the minimal valuation satisfying it.

For example, in the formula $p \Rightarrow \langle \mathbf{B}_i \rangle p$, p appears once in the antecedent. The Standard Translation yields $\forall P, (P(x) \implies \exists y, (R_b^i(x, y) \cdot P(y)))$. We instantiate P using the above lambda function, and use beta reduction to simplify, yielding $x = x \implies \exists y, (R_b^i(x, y) \cdot y = x)$.

If p appears twice in the antecedent, either as an atom or behind a diamond, each instance appears in the translated sentence in some form like $\exists x_1 (R_b^i(w, x_1) \cdot P(x_1))$ (for diamonds) $\dots \cdot P(x_2)$. The lambda function to serve as P 's instance is $\lambda u. (u = x_1 + u = x_2)$. For more than two appearances as atoms or behind diamonds, this pattern is repeated.

For proposition letters behind boxes (*e.g.* $\mathbf{B}_i p$), lambda functions of the form $\lambda u. (\bigvee_{j=0}^n (R_b^i(x_j, u)))$ are used, for the number of proposition letters n in the formula. When a proposition letter appears behind boxes and atomically, the atomic instances use the $=$ relation in the lambda function.

Step 3. Perform the instantiations. This step substitutes the predicate variables for the minimal valuations constructed in *Step 2*. For example, if $P(x)$ appears in the second order formula's antecedent as a singular instance of P , with the minimal valuation substitution $\lambda u.u = x$, the instantiated predicate instance yields the sentence $(\lambda u.u = x)(x)$, which beta reduces to $x = x$. If the proposition letter instance appears behind a box in the modal formula, *e.g.* $\mathbf{B}_i p$, the second order formula is $\forall y(R_b^i(x, y) \implies P(y))$ with the lambda function $\lambda u.(R_b^i(x, u))$, which yields the formula $\forall y(R_b^i(x, y) \implies \lambda u.(R_b^i(x, u)))(y)$. This beta reduces to $\forall y(R_b^i(x, y) \implies R_b^i(x, y))$, which is equivalent to $\forall y, R_b^i(x, y)$. For more examples, see Blackburn *et. al.*[20] chapter 3.

The result is a first order logic frame condition corresponding to the axiom schema. We apply this algorithm to each \mathcal{DASL} axiom schema and show that it is equivalent to a frame condition comprising \mathcal{DASL} frames.

Lemma 3.5.4 (Schema to First Order Frame Condition Correspondence). *Each static \mathcal{DASL} axiom schema corresponds to a first order frame condition for \mathcal{DASL} frames.*

Proof. We must show a translation from the following axiom schemas to the appropriate first order frame conditions:

1. $\mathbf{K}_i \varphi \Rightarrow \varphi$
2. $\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$
3. $\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi$
4. $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$

$\mathbf{K}_i p \Rightarrow p$. Applying the Standard Translation rules yields $\forall P(\forall y, R_k^i(x, y) \implies P(y)) \implies P(x)$. We have a single predicate with one instance in the antecedent behind a box. This yields the lambda function $\lambda u.(R_k^i(x, u))$. We substitute this

instance of P into the formula to yield:

$\forall y, (R_k^i(x, y) \implies \lambda u. (R_k^i(x, u))(y)) \implies \lambda u. (R_k^i(x, u))(x)$, which beta reduces and simplifies to:

$\forall y, R_k^i(x, y) \implies R_k^i(x, x)$, which is equivalent to $R_k^i(x, x)$ ⁸.

$\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$. We skip ahead to this more interesting example.

$$\begin{aligned}
& \rightsquigarrow \forall P, (\forall y, R_b^i(x, y) \implies P(y)) \implies (\forall z, (R_k^i \circ R_b^k)(x, z) \implies P(z)) \\
& \rightsquigarrow \forall P, \forall y, z, (R_k^i \circ R_b^k)(x, z) \cdot (R_b^i(x, y) \implies P(y)) \implies P(z) \\
& \rightsquigarrow \forall y, z, (R_k^i \circ R_b^k)(x, z) \cdot (R_b^i(x, y) \implies \lambda u. (R_b^i(x, u))(y)) \implies \lambda u. (R_b^i(x, u))(z) \\
& \rightsquigarrow \forall y, z, (R_k^i \circ R_b^k)(x, z) \cdot R_b^i(x, y) \implies R_b^i(x, z)
\end{aligned}$$

$\forall z, (R_k^i \circ R_b^k)(x, z)$ is equivalent to $\forall z, z', R_b^i(x, z') \cdot R_k^i(z', z)$. Thus, the above translation simplifies to $\forall z, (R_k^i \circ R_b^k)(x, z) \implies R_b^i(x, z)$, which is the subset relation of the frame condition.⁹

$\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$.

$$\begin{aligned}
& \rightsquigarrow \forall P, y, (R_b^i(x, y) \implies P(y)) \implies \exists z, (R_b^i(x, z) \cdot P(z)) \\
& \rightsquigarrow \forall y, (R_b^i(x, y) \implies \lambda u. (R_b^i(x, u))(y)) \implies \exists z, (R_b^i(x, z) \cdot \lambda u. (R_b^i(x, u))(z)) \\
& \rightsquigarrow \forall y, (R_b^i(x, y)) \implies \exists z, (R_b^i(x, z))
\end{aligned}$$

which is equivalent to frame condition for seriality.¹⁰

⁸Proof: Assume $R_k^i(x, x)$. Then by monotonicity, $\forall y, (R_k^i(x, y) \implies R_k^i(x, x))$. For the other direction, by contraposition, assume $\sim R_k^i(x, x)$. Then $\exists y, \sim R_k^i(x, y)$, which is equivalent to $\sim \forall y, R_k^i(x, y)$

⁹Because $\forall y, z, z', R_b^i(x, y) \cdot R_b^i(x, z') \cdot R_k^i(z', z)$ can rewrite z' with y and simplify to $R_b^i(x, y) \cdot R_k^i(y, z)$.

¹⁰In one direction by monotonicity, in the other by instantiating the consequent and then the antecedent, and then existentially generalizing.

$$\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi.$$

$$\begin{aligned}
& \rightsquigarrow \forall P, y, (R_k^i(x, y) \Rightarrow P(y)) \Rightarrow \forall z, (R_b^i(x, z) \Rightarrow P(z)) \\
& \rightsquigarrow \forall y, (R_k^i(x, y) \Rightarrow \lambda u. (R_k^i(x, u))(y)) \Rightarrow \forall z, (R_b^i(x, z) \Rightarrow \lambda u. (R_b^i(x, u))(z)) \\
& \rightsquigarrow \forall y, (R_k^i(x, y)) \Rightarrow \forall z, (R_b^i(x, z) \Rightarrow R_k^i(x, z)) \\
& \equiv \forall y, z, (R_k^i(x, y) \Rightarrow R_b^i(x, z) \Rightarrow R_k^i(x, z)) \\
& \equiv \forall y, z, (R_k^i(x, y) \cdot R_b^i(x, z) \Rightarrow R_k^i(x, z)) \\
& \equiv \forall z, (R_b^i(x, z) \Rightarrow R_k^i(x, z))^{11}
\end{aligned}$$

which is the frame condition for Knowledge implies Belief. \square

Lemma 3.5.4, in conjunction with the proof that each static \mathcal{DASL} axiom schema is a Sahlqvist formula, prove that \mathcal{DASL} is complete for \mathcal{DASL} frames as we have defined them.

Theorem 3.5.5 (\mathcal{DASL} complete on \mathcal{DASL} frames). *If \mathbf{F} is a \mathcal{DASL} frame, then $\mathbf{F} \models \varphi$ implies $\vdash \varphi$ in \mathcal{DASL} .*

Proof. By the SComT, \mathcal{DASL} is complete with respect to the class of frames its axiom schemas define. We mechanically proved that each static \mathcal{DASL} schema is a Sahlqvist formula, which from SCorT it follows that each axiom schema corresponds to a first order frame condition. By the Sahlqvist-van Benthem algorithm, each axiom schema translates to a corresponding \mathcal{DASL} frame condition. Therefore, the axiom schemas define the \mathcal{DASL} frame conditions. Thus, \mathcal{DASL} is complete with respect to the \mathcal{DASL} class of frames. \square

This concludes our discussion of completeness. We include some additional discussion of our mechanization in the next section.

3.6 Discussion

To help validate our mechanization of the `sahlqvist_formula` function in Coq, we run it against some modal formulas whose status as Sahlqvist formulas, in the positive and negative cases, is known.

Ló b’s formula $\Box(\Box\varphi \Rightarrow \varphi) \Rightarrow \Box\varphi$, the characteristic axiom of provability logic, is not a Sahlqvist formula, and lacks a first order correspondent[11].

To remain in the mechanized axiom schema logic we’ve produced in Coq, we use the knowledge operator as the \Box operator of the formula, and prove:

```
Example Lob_not_sahlqvist : forall (phi : prop) (a : DASL.Agents),
  not sahlqvist_formula (
    SK a (SK a (SProp phi) =s=> (SProp phi)) =s=> SK a (SProp phi)
  ).
Proof.
intros. unfold not. not_sahlqvist.
Qed.
```

It makes use of a custom tactic we use in the Coq file.

```
Ltac not_sahlqvist :=
  try (unfold sahlqvist_formula;
    unfold sahlqvist_implication;
    simpl; intuition).
```

We define a similar tactic for verifying that formulas are Sahlqvist formulas.

```

Ltac sahlqvist_reduce :=
  simpl; try (unfold sahlqvist_implication; split);
  try (unfold positive_formula; simpl; intuition);
  try (unfold sahlqvist_antecedent; simpl; intuition; unfold normal_form).

```

We make use of this tactic in, for example,

```

Lemma B_5_is_sahlqvist : forall (phi : prop) (a : DASL.Agents),
  sahlqvist_formula (\ (SB a (\ SB a (SProp phi))) =s=> (SB a (SProp phi))).
Proof.
  intros; sahlqvist_reduce.
Qed.

```

We prove that each of the modal axioms schemas of \mathcal{DASL} are Sahlqvist formulas using the above techniques, which can be seen in the appendix. Another example Sahlqvist formula that we prove is the Church-Rosser formula, $\Diamond\Box\varphi \Rightarrow \Box\Diamond\varphi$, again using the knowledge operator:

```

Example Church_Rosser_is_sahlqvist : forall (phi : prop) (a : DASL.Agents),
  sahlqvist_formula (
    \ (SK a (\
      (SK a (SProp phi)))) =s=>
      (SK a (\ (SK a (\ (SProp phi)))))).
Proof.
  intros; sahlqvist_reduce.
Qed.

```

However, if one rearranges the \Diamond and \Box in Church-Rosser, this yields its converse,

the McKinsey formula, which is not a Sahlqvist formula: $\Box\Diamond\varphi \Rightarrow \Diamond\Box\varphi$, and we likewise can prove this.

```

Example McKinsey_not_sahlqvist : forall (phi : prop) (a : DASL.Agents),
  not sahlqvist_formula (
    SK a (\ (SK a (\ (SProp phi)))) =s=>
    \ (SK a (\ (SK a (SProp phi)))).
Proof.
  intros; not_sahlqvist.
Qed.

```

The sharp-eyed reader might notice that we have not defined the `prop` type that appears in the above Coq formulas. Our mechanization of \mathcal{DASL} consists in two levels. The top level is that of the axiom schema, which allows us to prove theorems about the logic using a metavariable for propositions, just as one sees logicians do in papers. The metavariable represents propositions in the object language, which bottoms out in atomic propositions about airplane instruments in this case. With this architecture, one could define alternative object languages in Coq for different domains, perhaps one for driving automobiles, one for power plant operators, and one for submarine operators. Each of these domains involves well-structured relationships between the instrument readings and the actions' safety status, but each concerns different atomic propositions. However, the top level metalanguage remains the same.

Our mechanization of `prop` is discussed in the following chapter, where we present aviation mishap case studies mechanized in the object language.

3.7 Dynamic Extension

The previous section established \mathcal{DASL} as a sound and complete logic, which is a desirable formal property of the static base, since the dynamic extension is reducible to the static base in finitely many steps, it is also sound and complete.

Below are the axioms characterizing the reduction laws from the dynamic logic to a purely static logic through recursive application.

$[\mathbf{A}, a]_i p \Leftrightarrow (post(a, pre(a)) \Rightarrow p)$	Atomic Consequence
$[\mathbf{A}, a]_i \neg\varphi \Leftrightarrow (pre(a) \Rightarrow \neg[\mathbf{A}, a]_i \varphi)$	Action Negation
$[\mathbf{A}, a]_i (\varphi \wedge \psi) \Leftrightarrow ([\mathbf{A}, a]_i \varphi \wedge [\mathbf{A}, a]_i \psi)$	Action Conjunction
$[\mathbf{A}, a]_i \mathbf{K}_i \varphi \Leftrightarrow (pre(a) \Rightarrow \bigwedge_{a\chi_k^i\beta} \mathbf{K}_i [\mathbf{A}, \beta]_i \varphi)$	Action and Knowledge
$[\mathbf{A}, a]_i \mathbf{B}_i \varphi \Leftrightarrow (pre(a) \Rightarrow \bigwedge_{a\chi_b^i\beta} \mathbf{B}_i [\mathbf{A}, \beta]_i \varphi)$	Action and Belief
$[\mathbf{A}, a]_i^S p \Leftrightarrow (post(a, pre_s(a)) \Rightarrow p)$	Safe Atomic Consequence
$[\mathbf{A}, a]_i^S \neg\varphi \Leftrightarrow (pre_s(a) \Rightarrow \neg[\mathbf{A}, a]_i^S \varphi)$	Safe Action Negation
$[\mathbf{A}, a]_i^S (\varphi \wedge \psi) \Leftrightarrow ([\mathbf{A}, a]_i^S \varphi \wedge [\mathbf{A}, a]_i^S \psi)$	Safe Action Conjunction
$[\mathbf{A}, a]_i^S \mathbf{K}_i \varphi \Leftrightarrow (pre_s(a) \Rightarrow \bigwedge_{a\chi_k^i\beta} \mathbf{K}_i [\mathbf{A}, \beta]_i^S \varphi)$	Safe Action and Knowledge
$[\mathbf{A}, a]_i^S \mathbf{B}_i \varphi \Leftrightarrow (pre_s(a) \Rightarrow \bigwedge_{a\chi_b^i\beta} \mathbf{B}_i [\mathbf{A}, \beta]_i^S \varphi)$	Safe Action and Belief

Table 3.5: The reduction axioms of \mathcal{DASL} .

Theorem 3.7.1. *Translating a dynamic formula from the dynamic extension of \mathcal{DASL} via the reduction axiom schemas terminates in finitely many steps with an equivalent static formula*

Proof. It is clear from observing each of the reduction axiom schemas that they result in formulas with smaller dynamic components, *i.e.* smaller parse trees for those components. Eventually an iterative application of the reduction axiom schemas results in a purely static formula that has preserved the truth of the original formula. \square

The truth-preservation of the reduction axiom schemas has not yet been established, and we do so here.

Proof. Atomic Consequence : \Rightarrow . Assume $M, w \models [\mathbf{A}, a]_i p$. We must show that $M, w \models (post(a, pre(a)) \Rightarrow p)$. By the semantics of $[\mathbf{A}, a]_i$, for all (M', w') , if $M, w \models pre(a)$ and $update(M, A, w, a) = (M', w')$, then $M', w' \models p$. By definition of $post(a, p)$, if $update(M, A, w, a) = (M', w')$ and $M', w' \models p$, then $post(a, p) = p$. So, if $M, w \models pre(a)$ is defined, then $post(a, pre(a)) = p$, and thus $post(A, p) \Rightarrow p$.

\Leftarrow . Assume $M, w \models (post(a, pre(a)) \Rightarrow p)$. By the definition of $post(a, p)$, if $post(a, p) = p$ then $update(M, A, w, a) \models p$. So, if $M, w \models pre(a)$, then $update(M, A, w, a) \models p$. Therefore, $M, w \models [\mathbf{A}, a]_i p$.

Atomic Negation : \Rightarrow . Assume $M, w \models [\mathbf{A}, a]_i \neg\varphi$. It suffices to show that $M, w \models pre(a) \Rightarrow \langle \mathbf{A}, a \rangle_i \neg\varphi$. From the assumption and the semantics, for all (M', w') , if $(M, w) \Vdash (A, a)_i (M', w')$ then $M', w' \models \neg\varphi$. So, if $M, w \models pre(a)$ and $update(M, A, w, a) = (M', w')$, then $M', w' \models \neg\varphi$. Assume $M, w \models pre(a)$, and it follows that $update(M, A, w, a)$ is defined, so there exists a M', w' such that $(M, w) \Vdash (A, a)_i (M', w')$ and $update(M, A, w, a) = (M', w')$ and $M', w' \models \neg\varphi$. Therefore, $M, w \models pre(a) \Rightarrow \langle \mathbf{A}, a \rangle_i \neg\varphi$.

\Leftarrow . Assume $M, w \models pre(a) \Rightarrow \neg[\mathbf{A}, a]_i \varphi$. This is equivalent to $M, w \models pre(a) \Rightarrow \langle \mathbf{A}, a \rangle_i \neg\varphi$. By the semantics, if $M, w \models pre(a)$, then there exists a (M', w') such that $(M, w) \Vdash (A, a)_i (M', w')$ and $M', w' \models \neg\varphi$. The relation $\Vdash (A, a)_i$ is functional, so \exists implies \forall . So, for all (M', w') , if $(M, w) \Vdash (A, a)_i (M', w')$, then $M', w' \models \neg\varphi$, and therefore $M, w \models [\mathbf{A}, a]_i \neg\varphi$.

Action Conjunction is obvious.

Action and Knowledge. For this proof, assume for simplicity, without loss of generality, that $Actions = \{a\}$.

\Rightarrow . Assume $M, w \models [\mathbf{A}, a]_i \mathbf{K}_i \varphi$. Unfolding the semantics, for all (M', w') , if $(M, w) \models pre(a)$ and $update(M, A, w, a) = (M', w')$, then $M', w' \models \mathbf{K}_i \varphi$. $M', w' \models \mathbf{K}_i \varphi$ iff for all $v \in W$, if $w R_k^i v$ and $M, v \models pre(a)$ and $update(M, A, v, a) = (M', v')$ and $a \chi_k^i a$, then $M', v' \models \varphi$. That is, $M, w \models \mathbf{K}_i [\mathbf{A}, a]_i \varphi$.

\Leftarrow . Assume $M, w \models pre(a) \Rightarrow \mathbf{K}_i [\mathbf{A}, a]_i \varphi$. We must show $M, w \models [\mathbf{A}, a]_i \mathbf{K}_i \varphi$. Thus, we must show $M, w \models pre(a)$ and $update(M, A, w, a) = (M', w')$ implies

$M', w' \models \mathbf{K}_i \varphi$. So it suffices to show that if $update(M, A, w, a) = (M', w')$ and $M, w \models \mathbf{K}_i [\mathbf{A}, a]_i \varphi$, then $M', w' \models \mathbf{K}_i \varphi$. Assume $update(M, A, w, a) = (M', w')$ and $M, w \models \mathbf{K}_i [\mathbf{A}, a]_i \varphi$. Then for all v , if $wR_k^i v$, then $M, v \models [\mathbf{A}, a]_i \varphi$. It follows that $M, v \models pre(a)$ and $update(M, A, v, a) = (M', v')$ implies $M', v' \models \varphi$. Since $wR_k^i v$ and $a\chi_k^i a$, it holds that $w'R_k^{i'} v'$. Thus, $M', w' \models \mathbf{K}_i \varphi$.

Proofs for **Action and Belief** through **Safe Action and Belief** follow the above proofs exactly analogously. □

Thus, the dynamic extension inherits soundness and completeness from the static base. We proceed now to the philosophical justification of the dynamic extension, specifically the axiom schemas of *Inevitability* and *Minimum Rationality*.

3.7.1 Inevitability and Minimum Rationality

Recall the two axiom schemas under consideration.

$[\mathbf{A}, a]_i \varphi \Rightarrow [\mathbf{A}, a]_i^S \varphi$	Inevitability
$\langle \mathbf{A}, a \rangle_i \varphi \Rightarrow \mathbf{B}_i \langle \mathbf{A}, a \rangle_i^S \varphi$	Minimum Rationality

Table 3.6: \mathcal{DASL} axiom schemas of dynamic non-reductive character.

The axiom schema for *Inevitability* states that if engaging in action a from action structure A guarantees that φ is true as a result, then safely engaging in that action also guarantees that φ will result. One way to think about this is that the postcondition of a mere action is the same function as the postcondition for a safe action. The converse of *Inevitability* does not hold, however, because pre is defined on a subset of worlds on which pre_s is defined for a given action, representing the notion that it is *easier* to merely engage in an action than to do so safely. The subset

relation here represents number of constraints. Thus there may be worlds on which the safety precondition for an action is undefined but on which the precondition is defined. When the safety precondition is undefined it is trivially true that engaging in the action safely guarantees φ , for all φ . However, φ may not be guaranteed given mere engagement in the action.

The axiom schema for *Minimum Rationality* states that if after engaging in a φ results, then i believes that after safely engaging in a φ results. It is more intuitive when the implication is read as an “only if”. This axiom schema is clearly an idealization meant to contribute normative value to the logic rather than descriptive value. It can be interpreted as an assumption that the human agents being modeled do not engage in actions that they believe to be unsafe. For commercial pilots, this assumption is reasonably secure. For other domains, this assumption should be carefully noted and considered.

Chapter 4

Case Studies

4.1 Case Study and Mechanization

In this section we apply the logic just developed to the formal analysis of the Air France 447 aviation incident, then mechanize the formalization in the Coq Proof Assistant. Our mechanization follows similar work by Maliković and Čubriilo [68, 69], in which they mechanize an analysis of the game of Cluedo using Dynamic Epistemic Logic, based on van Ditmarsch’s formalization of the game [56]. It is commonly assumed that games must be adversarial, but this is not the case. Games need only involve situations in which players’ payoffs depend on the actions of other players. Similarly, knowledge games need not be adversarial, and must only involve diverging information. Thus, it is appropriate to model aviation incidents as knowledge games of sorts, where players’ payoffs depend on what others do, specifically the way the players communicate information with each other. The goal is to achieve an accurate situational awareness and provide flight control inputs appropriate for the situation. Failures to achieve this goal result in disaster, and often result from imperfect infor-

mation flow. A formal model of information flow in these situations provides insight and allows for the application of formal methods to improve information flow during emergency situations.

4.1.1 Air France 447

This case study is based on the authoritative investigative report into Air France 447 performed and released by France’s Bureau d’Enquêtes et d’Analyses pour la Sécurité de l’Aviation Civile (BEA), responsible for investigating civil aviation incidents and issuing factual findings[32]. The case is mechanized by instantiating, in Coq, the above logic to reflect the facts of the case. One challenge associated with this is that the readings about inputs present in aviation are often real values on a continuum, whereas for our purposes we require discrete values. We accomplish this by dividing the continuum associated with inputs and readings into discrete chunks, similar to how fuzzy logic maps defines predicates with real values[65].

Air France flight 447 from Rio de Janeiro, Brazil to Paris, France, June 1, 2009. The Airbus A330 encountered adverse weather over the Atlantic ocean, resulting in a clogged Pitot-static system. Consequently, the airspeed indicators delivered unreliable data concerning airspeed to the pilot flying, resulting in confusion. A chain of events transpired in which the pilot overcorrected the plane’s horizontal attitude again and again, and continued to input nose up pitch commands, all while losing airspeed. Perhaps most confusing to the pilot was the following situation: the aircraft’s angle of attack (AOA) was so high it was considered invalid by the computer, so no stall warning sounded until the nose pitched down into the valid AOA range, at which point the stall warning would sound. When the pilot pulled up, the AOA would be considered invalid again, and the stall warning would cease. The aircraft entered a spin and crashed into the ocean. Palmer [71] argues that had the pilot merely taken no action, the Pitot tubes would have cleared in a matter of seconds,

and the autopilot could have returned to Normal Mode.

This paper will formalize an excerpted instance from the beginning of the case, involving an initial inconsistency among airspeed indicators, and the subsequent dangerous input provided by the pilot. Formalized in the logic, the facts of the case allow us to infer that the pilot lacked negative introspection about the safety-critical data required for his action. This demonstrates that the logic allows information about the pilot's situational awareness to flow to the computer, via the pilot's actions. It likewise establishes a safety property to be enforced by the computer, namely that a pilot should maintain negative introspection about safety-critical data, and if he fails to do so, it should be re-established as quickly as possible.

According to the official report, at 2 hours and 10 minutes into the flight, a Pitot probe likely became clogged by ice, resulting in an inconsistency between airspeed indicators, and the autopilot disconnecting. This resulted in a change of mode from Normal Law to Alternate Law 2, in which certain stall and control protections ceased to exist. The pilot then made inappropriate control inputs, namely aggressive nose up commands, the only explanation for which is that he mistakenly believed that the aircraft was in Normal Law mode with protections in place to prevent a stall. This situation, and the inference regarding the pilot's mistaken belief, is modeled in the following application and mechanization of the logic.

We first introduce a lemma relating belief to knowledge, and then formalize the critical moment.

Lemma 4.1.1 (Belief is epistemically consistent). $\mathbf{B}_i \varphi \Rightarrow \neg \mathbf{K}_i \neg \varphi$.

Proof. From the fact that the belief modality is serial, it holds that

$$\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi,$$

which is equivalent to

$$\mathbf{B}_i \varphi \Rightarrow \neg \mathbf{B}_i \neg \varphi.$$

Due to axiom EP1, it follows that

$$\mathbf{B}_i \varphi \Rightarrow \neg \mathbf{K}_i \neg \varphi. \quad \square$$

We now formalize the critical moment.

1. $\neg(\text{Mode} = \text{Normal}) \dots$ — configuration.
2. $\langle \text{pilot} \rangle_i \text{hardnoseuptrue}$ — pilot input.
3. $\mathbf{B}_{\text{pilot}}(\text{Mode} = \text{Normal})$ — from axiom PR, pres_s .
4. $\neg \mathbf{K}_{\text{pilot}}(\text{Mode} = \text{Normal})$ — from axiom K-Reflexive.
5. $\mathbf{B}_{\text{pilot}} \mathbf{K}_{\text{pilot}}(\text{Mode} = \text{Normal})$ — from (3), axiom EP2.
6. $\neg \mathbf{K}_{\text{pilot}} \neg \mathbf{K}_{\text{pilot}}(\text{Mode} = \text{Normal})$ — from (5), Lemma 4.1.1.
7. $\neg \mathbf{K}_{\text{pilot}}(\text{Mode} = \text{Normal}) \wedge \neg \mathbf{K}_{\text{pilot}} \neg \mathbf{K}_{\text{pilot}}(\text{Mode} = \text{Normal})$ — from (4), (6).
8. $\neg(\neg \mathbf{K}_{\text{pilot}}(\text{Mode} = \text{Normal}) \Rightarrow \mathbf{K}_{\text{pilot}} \neg \mathbf{K}_{\text{pilot}}(\text{Mode} = \text{Normal}))$ — from (7).

The crux of the case is that inconsistent information was being presented to the pilot, along with a cacophony of inconsistent alarms, and the pilot's control inputs indicated a lack of awareness of safety-critical information. A detailed analysis, using the Coq Proof Assistant and the logic developed by my research, will make explicit these failures of information flow, both from the computer to the pilots, between

the pilot and co-pilot, and from the pilot to the computer. This will motivate the description of a prototype safety monitor that identifies and corrects information flow failures like those found in Air France 447.

4.1.2 Mechanization in Coq

Mechanical theorem proving divides into two categories, automated and interactive. Automated theorem proving combines a search algorithm with a proof checking algorithm to fully automate the process. The problem itself is undecidable in general, and human control is limited to the injection of hints prior to the algorithm’s execution. Interactive theorem proving, however, combines human-directed search with a proof checking algorithm, allowing the human to have more control over the procedure. Coq is a tool that facilitates interactive theorem proving.

The underlying logic of Coq is called the Calculus of Inductive Constructions, a dependently-typed constructive logic. One uses Coq by formalizing the target logic and its semantics in Coq and using what are called *tactics* to manipulate proof objects. My project will implement the previously described Safe Dynamic Agency Logic in Coq and formally model the Air France 447 case, demonstrating the logic’s ability to dynamically model safety-critical information flow in a real-world scenario. This will require translating the logic as it appears here and its metatheory into Coq, complete with tactics appropriate for the desired proofs and fully instantiated semantics, a process called *mechanization*.

The following mechanization demonstrates progress from the artificially simply toy examples normally analyzed in the literature to richer real-world examples. However, it does not represent the full richness of the approach. The actions and instrument readings mechanized in this paper are constrained to those most relevant to the case study. The approach is capable of capturing the full richness of all instrument reading configurations and actions available to a pilot. To do so, one needs to consult

a flight safety manual and formally represent each action available to a pilot, and each potential instrument reading, according to the following scheme.

We first formalize the set of agents.

```
Inductive Agents: Type := Pilot | CoPilot | AutoPilot.
```

Next we formalize the set of available inputs. These themselves are not actions, but represent atomic propositions true or false of a configuration. When a pilot manipulates an actuator via an action, the actuator being manipulated changes its state. We assume the pilot changes only one control surface actuator at a time, and so an input consists of the new state of whichever actuator is being manipulated.

```
Inductive Inputs : Type :=
  HardThrustPlus | ThrustPlus
  | HardNoseUp    | NoseUp
  | HardWingLeft  | WingLeft
  | HardThrustMinus | ThrustMinus
  | HardNoseDown  | NoseDown
  | HardWingRight | WingRight.
```

We represent readings by indicating which *side* of the panel they are on. Typically, an instrument has a left-side version, a right-side version, and sometimes a middle version serving as backup. When one of these instruments conflicts with its siblings, the autopilot will disconnect and give control to the pilot.

```
Inductive Side : TYPe := Left | Middle | Right.
```

We divide the main instruments into chunks of values they can take, in order to provide them with a discrete representation in the logic. For example, the read-

ing *VertUp1* may represent a nose up reading between 0° and 10°, while *VertUp2* represents a reading between 11° and 20°.

```

Inductive Readings (s : Side) : Type :=
  VertUp1 | VertUp2 | VertUp3 | VertUp4
  | VertDown1 | VertDown2 | VertDown3 | VertDown4
  | VertLevel | HorLeft1 | HorLeft2 | HorLeft3
  | HorRight1 | HorRight2 | HorRight3 | HorLevel
  | AirspeedFast1 | AirspeedFast2 | AirspeedFast3
  | AirspeedSlow1 | AirspeedSlow2 | AirspeedSlow3
  | AirspeedCruise | AltCruise | AltClimb | AltDesc | AltLand.

```

The difference between an input and a reading is that the inputs are values of actuator states, while readings are values of instrument states. In a properly functioning aircraft, the input values and readings values would correspond. However, in emergency situations, these tend to diverge.

We define a set of potential modes the aircraft can be in.

```

Inductive Mode : Type := Normal | Alternate1 | Alternate2.

```

We define a set of global instrument readings representing the mode and all of the instrument readings, left, right, and middle, combined together. This represents the configuration of the instrumentation.

```

Inductive GlobalReadings : Type := Global (m: Mode)
  (rl : Readings Left)
  (rm : Readings Middle)
  (rr : Readings Right).

```

The set of atomic propositions we are concerned with are those representing facts about the instrumentation.

```

Inductive Atoms : Type :=
  | M (m : Mode)
  | Input (a : Inputs)
  | InstrumentL (r : Readings Left)
  | InstrumentM (r : Readings Middle)
  | InstrumentR (r : Readings Right)
  | InstrumentsG (g : GlobalReadings).

```

Next we again follow Lescanne and Puisségur [66, 72] and Maliković and Čubrilo [68, 69] in defining *prop* of propositions of type **Type**. The definition provides constructors for atomic propositions consisting of particular instrument reading predicate statements, implications, propositions beginning with a knowledge modality, and those beginning with a belief modality.

```

Inductive prop : Type :=
  | atm : Atoms → prop
  | imp: prop → prop → prop
  | negp: prop → prop
  | falsum : prop
  | K : Agents → prop → prop
  | B : Agents → prop → prop

```

These should be familiar, as they are similar to the Backus-Naur form for the static modal logic from 3,

$$\varphi ::= p \mid \varphi \Rightarrow \varphi \mid \neg\varphi \mid \perp \mid \mathbf{K}_i \varphi \mid \mathbf{B}_i \varphi$$

We use the following notation for implication.

```
Infix "=>" := imp (right associativity, at level 85).
```

We likewise follow Lescanne and Puisségur [66, 72] by defining an inductive type *theorem* representing a theorem of \mathcal{DASL} . Mechanized, it takes a previously defined *prop* and returns a built-in Coq type *Prop* for propositions in Coq’s core logic. The constructors correspond to the Hilbert system, either as characteristic axioms, or inference rules. The first three represent axioms for propositional logic, then the rule Modus Ponens, then the axioms for the epistemic operator plus its Necessitation rule, then the doxastic operator and its Necessitation rule. Do not confuse the Necessitation rules with material implication in the object language. The final constructors capture the axioms relating belief and knowledge. The axioms for dynamic modal operators are defined separately, and are not included here.


```

Inductive theorem : prop → Prop :=
| Hilbert_K: forall p q : prop, theorem (p ⇒ q ⇒ p)
| Hilbert_S: forall p q r : prop,
theorem ((p ⇒ q ⇒ r) ⇒ (p ⇒ q) ⇒ (p ⇒ r))
| Classic_NOTNOT : forall p : prop, theorem ((NOT (NOT p)) ⇒ p)
| MP : forall p q : prop, theorem (p ⇒ q) → theorem p → theorem q
| K_Nec : forall (a : Agents) (p : prop), theorem p → theorem (K a p)
| K_K : forall (a : Agents) (p q : prop),
theorem (K a p ⇒ K a (p ⇒ q) ⇒ K a q)
| K_T : forall (a : Agents) (p : prop), theorem (K a p ⇒ p)
| B_Nec : forall (a : Agents) (p : prop), theorem p → theorem (B a p)
| B_K : forall (a : Agents) (p q : prop),
theorem (B a p ⇒ B a (p ⇒ q) ⇒ B a q)
| B_Serial : forall (a : Agents) (p : prop),
theorem (B a p ⇒ NOT (B a (NOT p)))
| B_4 : forall (a : Agents) (p : prop), theorem (B a p ⇒ B a (B a p))
| B_5 : forall (a : Agents) (p : prop),
theorem (NOT (B a p) ⇒ B a (NOT (B a p)))
| K_B : forall (a : Agents) (p : prop), theorem (K a p ⇒ B a p)
| B_BK : forall (a : Agents) (p : prop), theorem (B a p ⇒ B a (K a p)).

```

We use the following notation for *theorem*:

Notation " $| \vdash p$ " := (theorem p) (at level 80).

We encode actions as records in Coq, which allow for the definition of an inductive type along with projection functions. An action consists of the acting pilot, the observability of the action (whether it is observed by other agents or not), the input

provided by the pilot, and the preconditions for the action and the safety preconditions for the action, both represented as global atoms.

```
Record Action : Set := act {Ai : Agents; Aj : Agents; pi : PI;
input : Inputs; c : GlobalReadings;
c_s : GlobalReadings}.
```

The variable c holds the configuration representing the precondition for the action, while the variable c_s holds the configuration for the safety precondition. We encode the precondition and safety precondition functions as follows.

```
Function pre (a:Action) : prop := atm (InstrumentsG (c a)).
Function pre_s (a : Action) : prop := atm (InstrumentsG (c_s a)).
```

In the object language, the dynamic modalities of action and safe action are encoded as follows.

```
Parameter aft_ex_act : Action → prop → prop.
Parameter aft_ex_act_s : Action → prop → prop.
```

Many standard properties of logic, like the simplification of conjunctions, hypothetical syllogism, and contraposition, are encoded as Coq axioms. As an example, here is how we encode simplifying a conjunction into just its left conjunct.

```
Axiom simplifyL : forall p1 p2,
|== p1 & p2 → |== p1.
```

We formalize the configuration of the instruments at 2 hour 10 minutes into the flight as follows.

```

Definition Config_1 := (atm (M Alternate2)) &
(atm (InstrumentL (AirspeedSlow3 Left))) &
(atm (InstrumentM (AirspeedSlow3 Middle))) &
(atm (InstrumentR (AirspeedCruise Right))).

```

The mode is Alternate Law 2, and the left and central backup instruments falsely indicate that the airspeed is very slow, while the right side was not recorded, but because there was a conflict, we assume it remained correctly indicating a cruising airspeed.

The pilot's dangerous input, a hard nose up command, is encoded as follows.

```

Definition Input1 := act Pilot Pilot Pri HardNoseUp
(Global Alternate2 (AirspeedSlow3 Left)
(AirspeedSlow3 Middle)
(AirspeedCruise Right))
(Global Normal (AirspeedCruise Left)
(AirspeedCruise Middle)
(AirspeedCruise Right)).

```

The action is represented in the object language by taking the dual of the dynamic modality, $\neg[i, (A, a)]_i \neg True$, equivalently $\langle i \rangle_i (A, a) True$, indicating that the precondition is satisfied and the action token is executed.

```

Definition Act_1 := NOT (aft_ex_act Input1 (NOT TRUE)).

```

The actual configuration satisfies the precondition for the action, but it is inconsistent with the safety precondition. The safety precondition for the action indicates that the mode should be Normal and the readings should consistently indicate cruis-

ing airspeed. However, in Config_1, the conditions do not hold. Thus, the action is unsafe. From the configuration and the action, DASL allows us to deduce that the pilot lacks negative introspection of the action's safety preconditions.

Negative introspection is an agent's awareness of the current unknowns. To lack it is to be unaware of one's unknown variables, so lacking negative introspection about one's safety preconditions is to be unaware that they are unknown.

```
Theorem NegIntroFailMode :
  |-- (Config_1 ==>
    Act_1 ==>
      ((NOT (K Pilot (pre_s(Action1)))) &
        (NOT (K Pilot (NOT (K Pilot (pre_s(Action1)))))))).
```

In fact, in general it holds that if the safety preconditions for an action are false, and the pilot executes that action, then the pilot lacks negative introspection of those conditions. We have proven both the above theorem, and the more general theorem, in Coq.

```

Theorem neg_intro_failure :
forall (A Ao : Agents) (pi : PI) (inp : Inputs)
(m : Mode)
(rl : Readings Left) (rm : Readings Middle) (rr : Readings Right)
(ms : Mode)
(rls : Readings Left) (rms : Readings Middle) (rrs : Readings Right)
phi,
|-- (NOT
(aft_ex_act
(act A Ao pi inp (Global m rl rm rr) (Global ms rls rms rrs))
(NOT phi)) ==>
NOT (atm (InstrumentsG (Global ms rls rms rrs))) ==>
(NOT (K A (atm (InstrumentsG (Global ms rls rms rrs)))) &
(NOT (K A (NOT (K A (atm (InstrumentsG (Global ms rls rms rrs)))))))).

```

This indicates that negative introspection about safety preconditions is a desirable safety property to maintain, consistent with the official report's criticism that the Airbus cockpit system did not clearly display the safety critical information. The logic described in this research accurately models the report's findings that the pilot's lack of awareness about safety-critical information played a key role in his decision to provide unsafe inputs. Furthermore, the logic supports efforts to automatically infer which safety-critical information the pilot is unaware of and effectively display it to him.

The next section formalizes additional case studies in *DASL*.

4.2 Additional Case Studies

To illustrate the flexibility of this approach, we now formalize additional case studies in the logic. We analyze Copa Airlines flight 201 and Asiana Airlines flight 214.

4.2.1 Copa 201 and Asiana 214

Copa flight 201 departed Panama City, Panama for Cali, Colombia in June, 1992. Due to faulty wiring in the captain's Attitude Indicator, he incorrectly believed he was in a left bank position. In response to this, he directed the plane into an 80 degree roll to the right, which caused the plane to enter a steep dive. A correctly functioning backup indicator was available to the captain, and investigators believe that the captain intended to direct the backup indicator's readings to his own, but due to an outdated training module, the flip he switched actually sent his own faulty readings to the co-pilot's indicator. Approximately 29 minutes after takeoff, the plane crashed into the jungle and all passengers and crew perished. We formalize the moment at which the pilot provides the hard right roll input.

We begin with the present configuration of the instruments, simplified.

```
Definition Config_2 := (atm (M Normal))
    & (atm (InstrumentL (HorLeft2 Left)))
    & (atm (InstrumentM (HorLevel Middle)))
    & (atm (InstrumentR (HorLevel Right))).
```

We define the pilot's input and corresponding action.

```

Definition Input2 := act Pilot Pilot Pri HardWingRight
  (Global Normal (HorLeft2 Left) (HorLevel Middle) (HorLevel Right))
  (Global Normal (HorLeft2 Left) (HorLeft2 Middle) (HorLeft2 Right)).

Definition Act_2 := NOT (aft_ex_act Input2 (NOT TRUE)).

```

Finally, we state the theorem but omit the proof for space. It is available on the *DASL* Github.

```

Theorem NegIntroFailHorLevel :
  |-- (Config_2 ==> Act_2 ==>
    ((NOT (K Pilot (pre_s(Input2)))) &
      (NOT (K Pilot (NOT (K Pilot (pre_s(Input2))))))).

```

Asiana flight 214 from South Korea to San Francisco departed in the evening of July 6, 2013 and was schedule to land just before noon that morning [70]. The weather was good and air traffic control cleared the pilots to perform a visual approach to the runway. The plane came in short and crashed against an embankment in front of the runway, resulting in the deaths of three passengers and 187 injured. The National Transportation Safety Board (NTSB) investigation found that the captain had mismanaged the approach and monitoring of the airspeed, resulting in the plane being too high for a landing. Upon noticing this, the captain selected a flight mode (flight level change speed) which unexpectedly caused the plane to climb higher. In response to this, the captain disconnected the autopilot and pulled back on the thrust. This caused an autothrottle (A/T) protection to turn off, so when the captain pitched the nose down, the plane descended faster than was safe, causing it to come down too quickly and collide with the embankment in front of the runway. We will formalize

the moment at which the pilot pitches the nose down.

```

Definition Config_3 := (atm (M Alternate1))
    & (atm (InstrumentL (AirspeedSlow3 Left)))
    & (atm (InstrumentM (AirspeedSlow3 Middle)))
    & (atm (InstrumentR (AirspeedSlow3 Right))).

```

```

Definition Input3 := act Pilot Pilot Pri HardThrustMinus
    (Global Alternate1 (AirspeedSlow3 Left)
      (AirspeedSlow3 Middle)
      (AirspeedSlow3 Right))
    (Global Normal (AirspeedSlow3 Left)
      (AirspeedSlow3 Middle)
      (AirspeedSlow3 Right)).

```

```

Definition Act_3 := NOT (aft_ex_act Input3 (NOT TRUE)).

```

The corresponding theorem states that the pilot engaging in the action under those conditions is unaware of that the safety precondition is false.

```

Theorem NegIntroFailATOff :
  |-- (Config_3 ==> Act_3 ==>
    ((NOT (K Pilot (pre_s(Input3))))
    & (NOT (K Pilot (NOT (K Pilot (pre_s(Input3)))))))).

```

The above formalizations follow the same format as that of Air France 447. A pilot provides an input whose safety precondition conflicts with one of the instruments in the configuration, and we infer that the pilot lacks negative introspection of the

safety precondition. This is distinct from but related to the property that the pilots, in engaging in an unsafe action, are unaware of the unsafe instrument readings. We can capture this in the form of safety properties, which I turn to in the next section.

4.3 Safety Properties

Definition. Safety Negative Introspection (SNI). If a safety precondition does not hold, then agent knows that he does not know it to hold.

$$\neg pre_s(a) \Rightarrow \mathbf{K_i} \neg \mathbf{K_i} pre_s(a)$$

Definition. Safety-Critical Delivery (SCD). If a safety precondition is false, then agent knows that it is false.

$$\neg pre_s(a) \Rightarrow \mathbf{K_i} \neg pre_s(a)$$

Our above formalizations show that SNI is false when a pilot provides an unsafe input. Notice that SCD implies SNI.

Lemma 4.3.1 (SCD implies SNI). *Safety-Critical Delivery (SCD) implies Safety Negative Introspection (SNI).*

Proof. It suffices to show that $\mathbf{K_i} \neg \varphi \Rightarrow \mathbf{K_i} \neg \mathbf{K_i} \varphi$. Assume $\mathbf{K_i} \neg \varphi$ holds. From EP1, it follows that $\mathbf{B_i} \neg \varphi$, and because knowledge is a normal modality, it follows that $\mathbf{K_i} \mathbf{B_i} \neg \varphi$ holds. From lemma 1, and again the fact that knowledge is normal modality, it follows that $\mathbf{K_i} \neg \mathbf{K_i} \neg \varphi$, or equivalently, that $\mathbf{K_i} \neg \mathbf{K_i} \varphi$. \square

However, the converse does not hold. We can satisfy SNI when the safety precondition is false, the agent knows that he doesn't know it, but doesn't know that it is false. A counterexample consists of a model with three worlds: $\{u, v\}$. Let φ be

the safety precondition, with the following truth assignment: {False, True}. Let the epistemic relation include (u,v), (v,u), and the reflexive relations. Then at world u φ is false, and $\mathbf{K}_i \neg \mathbf{K}_i \varphi$ is true, but $\mathbf{K}_i \neg \varphi$ is false.

The formalizations show that from the pilot's unsafe action, it follows that he lacks negative introspection of the safety precondition.

$$\langle \mathbf{A}, a \rangle_i \text{true} \wedge \neg \text{pre}_s(a) \Rightarrow \neg \mathbf{K}_i \neg \mathbf{K}_i \text{pre}_s(a) \quad (4.1)$$

This situation violates SNI, because the pilot doesn't know that he doesn't know the safety precondition. Since SCD implies SNI, SCD is also violated.

$$\langle \mathbf{A}, a \rangle_i \text{true} \wedge \neg \text{pre}_s(a) \Rightarrow \neg \mathbf{K}_i \neg \text{pre}_s(a) \quad (4.2)$$

So, from an unsafe action, we can also infer that the pilot does not know that the safety precondition is false, a stronger conclusion.

Thus, by restoring knowledge that the safety precondition is false, it follows that either the safety precondition is true, or the unsafe action is not executed.

$$\mathbf{K}_i \neg \text{pre}_s(a) \Rightarrow \neg \langle \mathbf{A}, a \rangle_i \text{true} \vee \text{pre}_s(a) \quad (4.3)$$

The pilot's knowledge in the antecedent implies that the safety precondition is false, so this simplifies to:

$$\mathbf{K}_i \neg \text{pre}_s(a) \Rightarrow \neg \langle \mathbf{A}, a \rangle_i \text{true}. \quad (4.4)$$

This squares with the standard game theoretic inference, wherein a rational agent with knowledge of the situation executes a good action. Because our model of knowledge and rationality is weaker, we make the weaker claim that a minimally rational pilot with knowledge of the safety-critical information does not execute a bad action.

4.4 Decision Problem

An important extension of the foundational work provided by this paper is the construction of a system that takes advantage of the logic as a runtime safety monitor. It will monitor the pilot's control inputs and current flight configurations, and in the event that an action's safety preconditions do not hold, infer which instrument readings the pilot is unaware of and act to correct this. In order to avoid further information overload, the corrective action taken by the computer should be to temporarily remove or dim the non-safety-critical information from competition for the pilot's attention, until the pilot's unsafe control inputs are corrected, indicating awareness of the safety-critical information. Construction of a prototype of this system is underway. Here we define the decision problem and prove that it is NP-Complete. The decision problem we face is formalized as follows.

SD. Input: $a : Action$, $C : Configuration$, $pre_s : Action \mapsto Configuration$, $k : Int$.

Output: Is there a set of instruments I of size k from configuration C that falsifies $pre_s(a)$?

Theorem 4.4.1 (NP). *SD is NP-Complete.*

Proof. First, we prove that SD is in NP by defining the decision problem SDV that takes input to SD and a *certificate* and verifies that the certificate falsifies $pre_s(a)$ in polynomial time. Let the certificate be a set of instruments I of size k from the configuration C . Note that $pre_s(a)$ has the form $(c_1 \wedge c_2 \wedge \dots \wedge c_n) \vee (c'_1 \wedge c'_2 \wedge \dots \wedge c'_n) \vee \dots$. The negation of this has the form $(\neg c_1 \vee \neg c_2 \vee \dots \vee \neg c_n) \wedge (\neg c'_1 \vee \neg c'_2 \vee \dots \vee \neg c'_n) \wedge \dots$, which is in Conjunctive Normal Form. Computing the negation of pre_s can be done in linear time. Next we take I and treat it as an assignment of values to instruments of the form $(i_1 \wedge i_2 \wedge \dots \wedge i_n)$ for each $i \in I$. Then we check whether I satisfies $\neg pre_s(a)$ in polynomial time, treating I as the certificate for the

SAT problem.

Second, we prove that *SD* is in NP-Hard. We do this by providing a polynomial time reduction from *nSAT* to *SD*. Taking the input from *nSAT*, we let the size of the configuration $|C| = n$, that is, we let n be the number of instruments on the flight deck. The maximum size of any solution I to *SD* is therefore n . We let the *nSAT* formula be the negation of $pre_s(a)$. We iterate over $k \in \{1..n\}$. Thus, for *nSAT*, the *SD* problem is run at most n times. Any input to *nSAT* will be at least of size n , so the number of times *SD* is run is polynomial in the length of *nSAT*'s input. \square

4.4.1 Summary

This chapter mechanized an instance of *DASL* appropriate for reasoning about safety-critical information flow and action in the context of aviation mishaps. It illustrates how *DASL* is useful as a formal method for modeling human agents and the logical relationship between their actions and their mental states. Similar instances of *DASL* could occur in a variety of situations, where a human agent makes decisions and it is important to rapidly identify when unsafe actions occur, and which safety-critical information would improve the safety of those actions.

We proved that when a rational pilot engages in unsafe action, it is because she is unaware of some safety-critical information. We explored the implications of this insight by identifying the information assurance properties of *Safety Negative Introspection* (SNI) and *Safety-Critical Delivery* (SCD), which formalize the notions that an agent should be aware when an action's safety precondition is unknown, and that she should know when an action's safety precondition is false, respectively. We proved that the latter implies the former, and that the aviation mishaps involve violations of SNI (and therefore SCD as well). This shows that both negative introspection and direct knowledge of a safety precondition's being false are sufficient for preventing unsafe actions. This formalizes the standard game-theoretic inference, weakened for our

correspondingly weaker agents, that knowledge implies the *avoidance* of bad actions.

Finally, we showed that capitalizing on this insight via an algorithm is an NP-Complete problem. The decision problem takes as input an action, configuration, safety precondition, and integer k , and returns true or false whether there is a set of instrument readings of size k that falsifies the safety precondition for the given action.

The next chapter presents a theoretic worry facing \mathcal{DASL} , and indeed most modal logics of belief and knowledge developed in the literature so far.

Chapter 5

Löb’s Obstacle

Thus far, the dynamic modal logic we have constructed and applied has followed closely the standard formalizations in the literature, departing primarily in how it deals with an agent’s static knowledge. The static operator for an agent’s belief remains *KD45*, as it is normally formalized in the literature. Our dynamic operators follow closely those defined in [62] and [21], with the novel addition of a safety precondition to capture normative constraints on action in addition to logical constraints. The dissertation thus far, and the logics of knowledge and agency in the literature, ignore a problem facing agents of a certain reflective type, identified by Smullyan in [78].¹ This chapter confronts the problem of reflective reasoners facing modal logics of knowledge and agency, including *DASL*. We respond with an adjustment to *DASL*’s static foundation so that the problem can be avoided, if the agents being modeled are reflective reasoners.

In Smullyan’s, “Logicians who reason about themselves,” he considers epistemic problems related to undecidability results in mathematics. He identifies “a complete parallelism between logicians who believe propositions and mathematical systems that

¹Smullyan refers to these as reflexive reasoners, but we use the term ‘reflective’ in order to avoid ambiguity with the reflexive frame condition on worlds.

prove propositions.” In provability logic, the formula $\varphi \Leftrightarrow \neg \mathbf{B}_i \varphi$ expresses the Gödel proposition, “This proposition is not provable in system i .” In a doxastic interpretation, the same formula expresses the reflective belief, “agent i does not believe this proposition.” This means that for any doxastic or epistemic system, if the agents it models are reflective, then it must not have the theorem $\mathbf{B}_i \varphi \Rightarrow \varphi$, under pain of contradiction.

A *reflective* agent is one that can form beliefs and knowledge about self-referential sentences and propositions. Such propositions are of the form $\varphi \Leftrightarrow (\mathbf{B}_i \varphi \Rightarrow \psi)$. The reflective proposition “agent i does not believe this proposition” is of that form, where ψ is replaced with \perp in the generic form. More specifically, an agent is *reflective* just in case for every proposition ψ , there is a proposition φ such that $\varphi \Rightarrow (\mathbf{B}_i \varphi \Rightarrow \psi)$.

The above uses the \mathbf{B}_i operator for belief, but just as easily we could have used \mathbf{K}_i . The problem for epistemic logic is immediately apparent. Most humans are capable of reasoning about self-referential propositions, so a logic for human knowledge ought to include such propositions. However, with the Truth Axiom of the knowledge operator, this seems to yield inconsistency. It seems to yield the very same inconsistency that mathematical system i would face if it could prove its own soundness.

In what follows, we examine the mathematical logic that underpins this issue. Just as Smullyan identified the problem as one facing reflective reasoners with Axioms K and 4 and the Rule of Necessitation in [78], Löb identified these same conditions as the ones that allow a mathematical system to derive his theorem in [28]. We highlight the obstacle this presents to epistemic and doxastic logics, and show how to avoid it. We present a relaxed *DASL* for reasoning about safety-critical information flow to reflective agents.

The potential problem is due to Löb in [28] from the mathematical logic perspective, and Smullyan in [78] from the epistemic/doxastic logic perspective. Contemporary research in artificial intelligence foundations addresses the problem as it

pertains to agents being confident in their own conclusions. See for example research at the Machine Intelligence Research Institute², who kindly invited the author to a workshop introducing him to the problem. Some of their preliminary results on the problem are available on the footnoted website. They have named the problem Löb’s Obstacle, or the Löbian Obstacle. We describe it, and the problem it presents to modal logic of agency, including our own, below.

5.1 Löb’s Theorem

Löb’s Theorem takes its name from Martin Hugo Löb, who tackled a question of mathematical logic posed by Leon Henkin in the years following the results of Gödel. Henkin asked what could be said of propositions asserting their own provability, as opposed to unprovability in the case of Gödel sentences.[79] Löb answered by showing that in a consistent system, proof that proof of a proposition implies that the proposition is true (soundness) only for propositions that are actually provable. This is surprising, because it generalizes Gödel’s Second Incompleteness Theorem: no formal system can prove its own soundness (as opposed to consistency).

Löb’s Theorem in provability logic is,

$$\Box(\Box\varphi \Rightarrow \varphi) \Rightarrow \Box\varphi. \tag{5.1}$$

The \Box is interpreted as “provability” in some formal system, particularly one at least as powerful as Peano arithmetic. If a modal operator involves the reasoning abilities of a human-like agent, then *a fortiori* it is a formal system at least as powerful as Peano arithmetic. This presents the following problem. If that agent, in its own reasoning

²See their website at <https://www.intelligence.org>

system, can deduce the soundness of its own system, then its reasoning system is unsound. This is because φ can be any formula, including \perp . The particular modal logic at risk, in our mind, is epistemic logic, or any doxastic logic with accurate reasoners.

Löb identified the following conditions of a formal system that would allow derivation of his theorem (which we present in modal form).

1. $\Box(\varphi \Rightarrow \psi) \Rightarrow (\Box\varphi \Rightarrow \Box\psi)$ Axiom K

2. $\Box\varphi \Rightarrow \Box\Box\varphi$ Axiom 4

3. From $\vdash \varphi$, infer $\Box\varphi$ Rule of Necessitation

Items (1) and (3) are constants for all normal modal logics. There is a suppressed condition that Löb did not mention because it was unnecessary for the domain of arithmetic, but we must mention it here. The system must either admit of self-referential sentences or involve modal fixed points. Because systems of human-like reasoning must include self-referential sentences, at least in the form of Peano arithmetic, in order to remain human-like, this condition is satisfied for our concerns.

Here we give a template derivation of Löb's Theorem, which we shall refer to below when describing how Löb's Obstacle corrupts various epistemic logics.

Proof.

(1.1) $\Box(\Box\varphi \Rightarrow \varphi)$ Assumption

(1.2) $\Box(\psi \Leftrightarrow (\Box\psi \Rightarrow \varphi))$ Löb Sentence³

(1.3) $\Box(\Box\psi \Leftrightarrow \Box(\Box\psi \Rightarrow \varphi))$ Axiom K

(1.4) $\Box(\Box\psi \Rightarrow \Box(\Box\psi \Rightarrow \varphi))$ (1.3) Simplification of \Leftrightarrow

³Sometimes referred as a Curry sentence after logician Haskell Curry.

(1.5)	$\Box(\Box\psi \Rightarrow (\Box\Box\psi \Rightarrow \Box\varphi))$	(1.4) Axiom K
(1.6)	$\Box(\Box\psi \Rightarrow \Box\Box\psi)$	Axiom 4
(1.7)	$\Box(\Box\psi \Rightarrow \Box\varphi)$	(1.5), (1.6)
(1.8)	$\Box(\Box\psi \Rightarrow \varphi)$	(1.7), (1.1)
(1.9)	$\Box\psi$	(1.3), (1.8)
(1.10)	$\Box\Box\psi$	(1.9), Axiom 4
(1.11)	$\Box\Box\psi \Rightarrow \Box\varphi$	(1.8), Axiom K
(1.12)	$\Box\varphi$	(1.10), (1.11)

QED

□

Mathematical and Provability logicians refer to the key components of this proof as Löb Conditions[11]. Identifying them in the proof above helps us identify which epistemic logics collide with Löb's Obstacle. Conversely, understanding how the Löb Conditions interact helps us construct epistemic logics that avoid the Obstacle.

The Conditions are:

1. The Löb Sentence. A self-referential sentence, also formalizable as a modal fixed point.
2. Axiom K. The standard distribution axiom of normal modal logics.
3. Axiom 4. The axiom corresponding to a transitive frame relation.
4. The rule of necessitation. Likewise a standard feature of normal modal logics.

The Löb Sentence is sometimes not mentioned as a Condition, because Löb's Theorem is typically studied in the context of mathematical logic or provability logic, where such self-referential expressiveness is known to exist. We point out, however,

that humans are capable of reasoning about self-referential sentences, as well, and any advanced artificial agent will be able to do so, as well.

Finally, we note the importance of Löb’s Theorem’s antecedent: $\Box(\Box\varphi \Rightarrow \varphi)$. Epistemic logics typically include the antecedent as a theorem, in which case Löb’s Theorem will allow us to derive $\Box\varphi$ for all φ . This is why consistent mathematical systems at least as expressive as Peano arithmetic cannot prove their own consistency.

In what follows, we summarize the problems this theorem presents to formal agency modeling in general, and a unique potential benefit it affords in game theory due to LaVictoire *et. al.* in [27]. Next, we round up the usual suspects of epistemic-doxastic logic and show that, on the assumption that they aim to capture human-like reasoning, they crash into Löb’s Obstacle. Unfortunately, *DASL* itself suffers from Löb’s Obstacle, and so we provide a weakened static base suitable for agents capable of self-reflective reasoning. It sustains the contraposed game theoretic inference we sought out to model for realistic agents.

5.2 Discussion on Löb’s Obstacle

Löb’s Obstacle has its origin in mathematical logic, where it was shown to prevent Peano arithmetic from proving its own soundness. Among modal logicians of a philosophical bent, mathematical logic, and its provability logic cousin, have received much less interest than epistemic, temporal, and deontic logics. Among modal logicians interested in technical aspects of logic, Löb’s Theorem has received interest insofar as it relates to fixpoints. One exception is Raymond Smullyan in [78], where he examined the parallels between Peano Arithmetic⁴ proving statements about itself and logicians deducing beliefs about themselves. In it, he points out that reasoners with certain powers will risk deviating into an inconsistent set of beliefs, due to self-

⁴We use this to refer to any system at least as expressive as PA.

reflection. The cited paper appeared as a curiosity of sorts, presenting a series of recreational logic puzzles or riddles, with corresponding solutions following.

More recently, researchers in formal agent modeling have become interested in this particular logic puzzle of Smullyan's. When the logician reasoning about itself is an artificial intelligence, the puzzle becomes a unique safety problem. In [83], Yudkowski and Herreshoff consider an AI of sufficient reasoning power (PA^+), and show by straightforward application of Löb's Theorem that such an agent can neither trust its own conclusions in general, nor assent to its, we would say *in fact*, correct reasoning principles. The alternative to this is an inconsistent agent.

According to Yudkowski and Herreshoff, the AI “cannot say, ‘I don’t know if φ is true, but I trust that if any proof exists from my axioms of φ , then φ must be true.’ Instead [the AI] can only react to particular proofs of φ by asserting φ . This presents an apparent paradox of rational coherence... the agent will believe any proof of φ you present to it, but will refuse to assent to ‘I ought to believe any proof of φ I am presented with,’ a seeming defect of reflectivity — the agent cannot approve, under reflection, its own behavioral principles.”

The safety problem is an instance of φ in the previous paragraph. One can substitute any safety property in for φ that a future AI ought to abide by. We can ignore problems of humans correctly encoding the safety property, and focus on the AI assessing actions under the constraint of φ . We assume some such actions will be of the meta variety, where it must decide whether to *a*) modify itself, *b*) create a copy of itself, or *c*) produce offspring AI. Any such meta type action causes it to confront a sentence of the form, “[meta type action *a*] does not violate φ .” If the AI is producing strictly weaker reasoning systems as offspring or self-modifications, it can safely infer that the sentence is true[83]. This can produce a descending chain of systems culminating in PA. However, for self-copies and non-weakening self-modifications, the AI cannot infer the above sentence without risking inconsistency.

Perhaps for the short-term, this problem can be avoided due to the fact that powerful AI programs tend to be *narrow*, or focused on a particular problem domain. So long as AI programs tend to be narrow, formal methods researchers may be able to specify safety properties and develop techniques for mechanically checking whether they hold. However, there are three reasons to be skeptical of this optimistic claim.

First, this avoids the technical problem by assuming it will not arise, not by solving it. Given a technical scenario with potentially unsafe consequences, proper safety engineering requires examination of worst case scenarios and engaging with them.

Second, one instance of narrow AI is a program focused on constructing software. This scenario seems highly likely to occur in the relatively near future, as more technology development abstracts away from and automates various levels of configuration and compilation. Such automation trends toward a program that takes as input a computational problem and attempts to construct a program that solves it. Due to the Curry-Howard Correspondence between proofs and programs, this type of program in fact already exists in the form of automated theorem provers. Should these programs continue to advance to a point where they are constructing other programming engines, the problem of Löb’s Obstacle immediately arises, even for narrow AI.

Finally, organizations are already actively working to develop artificial *general* intelligence, defined as an AI capable of solving problems from any domain at least as well as a human. A company which has made numerous headlines for its recent successes is DeepMind, a subsidiary of Alphabet, which states that its mission is, “We’re on a scientific mission to push the boundaries of AI, developing programs that can learn to solve any complex problem without needing to be taught how.”⁵ With a strong economic incentive to achieve artificial general intelligence, it is not un-

⁵<https://deepmind.com/about/> Accessed 7/9/19.

reasonable to engage with the safety concerns that arise from it *prior* to its successful creation.

The consequences of Löb’s Theorem for AI safety are as follows. First, ignoring chains of descendingly weaker systems, AI cannot safely be in the AI programming business. There is a sort of mathematical gravity that pulls a system toward Löb’s Obstacle, as it places two clear desiderata in conflict with one another. On the one hand, we seek agents with (near-) perfect *external* accuracy, with beliefs that are justifiedly true of the world. On the other hand, we seek agents with (near-) perfect *internal* accuracy, with powerful introspective access to their own internal states. Combining these desiderata yields an $\mathcal{S5}$ modality, and indeed, the ideally rational agent is one with $\mathcal{S5}$ epistemic modality. Aumann’s formalization of the ideally rational agents of game theory and economics is the $\mathcal{S5}$ modal operator in [73]. His was not an arbitrary formalization, but rather one that made explicit assumptions that mathematical models in those domains make about agents. Thus, there is a tendency, in constructing artificial agents, to build toward that. Such agents have the powerful positive and negative introspection properties, and the axiom that knowledge implies truth.

However, constructing such an agent capable of reasoning about PA^+ is guaranteed to eventually decay into inconsistency, by Löb’s Theorem, as pointed out by Smullyan. Such agents satisfy the antecedent of Löb’s Theorem, and are able to deduce Löb’s Theorem itself, and therefore conclude that they know any arbitrary formula. The problematic case here is not that the AI will come to believe that apples are blue, but rather that producing a next generation AI^+ is provably safe, when this is not the case.

As this problem relates to the goal of this thesis, which is to describe a logic of agency for human-like agents, we examine some static epistemic logics that crash into Löb’s Obstacle.

5.3 Epistemic Logics that Crash

5.3.1 S5 Epistemic Logic

The most prominent epistemic logic in the literature, by far, is *S5* epistemic logic. *S5* epistemic logic is routinely presented as the logic of knowledge, and often serves as a static base for dynamic extensions to epistemic logic involving action and communication. Its characteristic axioms are:

$$\mathbf{K}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \psi) \quad (5.2)$$

$$\mathbf{K}_i \varphi \Rightarrow \varphi \quad (5.3)$$

$$\neg \mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \neg \mathbf{K}_i \varphi \quad (5.4)$$

Clearly (5.2) is Axiom K, and (5.3), troublingly, combines with the Rule of Necessitation to produce the antecedent of Löb's Theorem. (5.4) is called the Negative Introspection axiom, or sometimes in philosophy circles, the Wisdom Axiom. Logicians call it Axiom 5. It is read, "If *i* does not know that φ , then she knows that she doesn't know it". Other than being clearly invalid for humans, this axiom and (3) allows us to derive,

$$\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi$$

Proof.

$$(3.1) \quad \neg \mathbf{K}_i \neg \mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \varphi \quad \text{Contrapositive of Axiom 5}$$

$$(3.2) \quad \mathbf{K}_i \neg \mathbf{K}_i \neg \mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi \quad \text{Rule of Necessitation on (5.1), Axiom K}$$

$$(3.3) \quad \varphi \Rightarrow \neg \mathbf{K}_i \neg \varphi \quad \text{Axiom T, Contrapositive}$$

$$(3.4) \quad \neg \mathbf{K}_i \neg \varphi \Rightarrow \mathbf{K}_i \neg \mathbf{K}_i \neg \varphi \quad \text{Axiom 5}$$

$$(3.5) \quad \varphi \Rightarrow \mathbf{K}_i \neg \mathbf{K}_i \neg \varphi \quad (4.3), (4.4)$$

$$(3.6) \quad \mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \neg \mathbf{K}_i \neg \mathbf{K}_i \varphi \quad \mathbf{K}_i \varphi / \varphi, (4.5)$$

$$(3.7) \quad \mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi \quad (4.2), (4.6)$$

QED

□

Thus, $\mathcal{S5}$ satisfies Löb's three conditions, if we assume the presence of self-referential sentences possible, which we should. Therefore, with \mathbf{K}_i instead of \Box , the proof of Löb's Theorem is possible in this brand of $\mathcal{S5}$. However, to make matters worse, the antecedent of Löb's Theorem is itself an axiom of $\mathcal{S5}$. Therefore, $\mathbf{K}_i \varphi$ is a theorem, for all φ .

We take this as a *reductio ad absurdum* that $\mathcal{S5}$ epistemic logic cannot be a logic for reasoning about the knowledge of agents with expressive power beyond Peano arithmetic. Therefore, it cannot be a logic of knowledge for humans, or human-like agents.

5.3.2 Hintikka's S4 Epistemic Logic

In Hintikka's 1967 *Knowledge and Belief: A logic of the two notions*, he presented an epistemic logic for determining the validity and consistency of claims people make about knowledge and belief. He rejected out of hand the negative introspection axiom for knowledge, but chose to include positive introspection, which is formalized as $\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi$. Clearly then, if Hintikka's epistemic system is meant for human-like reasoners who can express sentences like, "If I know this sentence is true, then $1 + 1 = 2$," then it crashes into Löb's Obstacle, with the extra bite of having the antecedent of Löb's Theorem as a theorem itself, and therefore, $\mathbf{K}_i \varphi$ is also a theorem.

5.3.3 Kraus and Lehmann System

In [36], Kraus and Lehmann tackle the issue of how to combine knowledge and belief in a single system remarked on in the conclusion of Halpern and Moses in [38]. The project at the time was very similar to that addressed by this dissertation's static foundation: the development of a logic sufficiently realistic and expressive for reasoning about intelligent agents.

They axiomatize knowledge and belief as follows.

$\mathbf{K_i}(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K_i} \varphi \Rightarrow \mathbf{K_i} \psi)$	Distribution of $\mathbf{K_i}$
$\mathbf{K_i} \varphi \Rightarrow \varphi$	Truth
$\neg \mathbf{K_i} \varphi \Rightarrow \mathbf{K_i} \neg \mathbf{K_i} \varphi$	Negative Introspection
$\mathbf{B_i}(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \psi)$	Distribution of $\mathbf{B_i}$
$\mathbf{B_i} \varphi \Rightarrow \langle \mathbf{B_i} \rangle \varphi$	Belief Consistency
$\mathbf{K_i} \varphi \Rightarrow \mathbf{B_i} \varphi$	Knowledge implies Belief
$\mathbf{B_i} \varphi \Rightarrow \mathbf{K_i} \mathbf{B_i} \varphi$	Conscious Belief
From $\vdash \varphi$ and $\vdash \varphi \Rightarrow \psi$, infer $\vdash \psi$	Modus Ponens
From $\vdash \varphi$, infer $\vdash \mathbf{K_i} \varphi$	Necessitation of $\mathbf{K_i}$

Table 5.1: Logic of Kraus and Lehmann

The key differences with *DASL* are that knowledge is *S5* and the axiom of Conscious Belief, which reverses the composition of knowledge and belief under belief relative to *DASL*'s Evidential Restraint axiom. From the above axiom schemas, it follows that belief is a regular *KD45* operator.

In their article, and in Meyer and van der Hoek's [42], they show that $\mathbf{B_i}(\mathbf{B_i} \varphi \Rightarrow \varphi)$ is a theorem. Therefore, this system suffers from a particularly bad collision with Löb: Both the knowledge and belief modalities lead to inconsistency for intelligent agents capable of self-referential reasoning.

5.3.4 \mathcal{DASL}

$\mathbf{K_i}(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K_i} \varphi \Rightarrow \mathbf{K_i} \psi)$	Distribution of $\mathbf{K_i}$
$\mathbf{K_i} \varphi \Rightarrow \varphi$	Truth
$\mathbf{B_i}(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \psi)$	Distribution of $\mathbf{B_i}$
$\mathbf{B_i} \varphi \Rightarrow \langle \mathbf{B_i} \rangle \varphi$	Belief Consistency
$\neg \mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \neg \mathbf{B_i} \varphi$	Negative Belief Introspection
$\mathbf{K_i} \varphi \Rightarrow \mathbf{B_i} \varphi$	Knowledge implies Belief
$\mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \mathbf{K_i} \varphi$	Evidential Restraint
From $\vdash \varphi$ and $\vdash \varphi \Rightarrow \psi$, infer $\vdash \psi$	Modus Ponens
From $\vdash \varphi$, infer $\vdash \mathbf{K_i} \varphi$	Necessitation of $\mathbf{K_i}$

Table 5.2: Logic of \mathcal{DASL}

Recall that in \mathcal{DASL} 's static logic, knowledge lacks both epistemic introspection properties, while belief has negative introspection and positive introspection. The modalities are combined via the Knowledge implies Belief Axiom and the Evidential Restraint Axiom. Löb's Theorem is derivable for the belief operator. We first show that the two composite axiom schemas are sufficient to derive positive introspection for belief.

Theorem 5.3.1 (Positive Belief Introspection). $\mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \mathbf{B_i} \varphi$

Proof.

$$(3.1) \quad \mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \mathbf{K_i} \varphi \quad \text{ER Axiom}$$

$$(3.2) \quad \mathbf{B_i} \mathbf{K_i} \varphi \Rightarrow \mathbf{B_i} \mathbf{B_i} \varphi \quad \text{KiB Axiom + Necessitation of } \mathbf{B_i}^6$$

$$(3.3) \quad \mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \mathbf{B_i} \varphi \quad (3.1), (3.2)$$

□

⁶This rule can be derived from Necessitation of $\mathbf{K_i}$ and KiB Axiom.

QED

With the three Löb Conditions (Axiom K, Axiom 4, and the Rule of Necessitation) satisfied for $\mathbf{B_i}$, Löb's Theorem holds for belief. The Belief Consistency Axiom $\mathbf{B_i} \varphi \Rightarrow \langle \mathbf{B_i} \rangle \varphi$ is equivalent to $\neg(\mathbf{B_i} \varphi \wedge \mathbf{B_i} \neg\varphi)$, which is furthermore equivalent to $\neg\mathbf{B_i}(\varphi \wedge \neg\varphi)$, which results in the following disaster.

Theorem 5.3.2 (Consistency Disaster). *If $\neg\mathbf{B_i}(\varphi \wedge \neg\varphi)$ and $\mathbf{B_i}(\mathbf{B_i} \varphi \Rightarrow \varphi) \Rightarrow \mathbf{B_i} \varphi$ are theorems, then $\mathbf{B_i}(\varphi \wedge \neg\varphi)$ is a theorem.*

Proof.

$$(3.1) \quad \neg\mathbf{B_i}(\varphi \wedge \neg\varphi) \quad \text{Belief is Consistent}$$

$$(3.2) \quad \mathbf{B_i}(\varphi \wedge \neg\varphi) \Rightarrow (\varphi \wedge \neg\varphi) \quad (4.1)$$

$$(3.3) \quad \mathbf{B_i}(\mathbf{B_i}(\varphi \wedge \neg\varphi) \Rightarrow (\varphi \wedge \neg\varphi)) \quad \text{Necessitation of } \mathbf{B_i}$$

$$(3.4) \quad \mathbf{B_i}(\mathbf{B_i}(\varphi \wedge \neg\varphi) \Rightarrow (\varphi \wedge \neg\varphi)) \Rightarrow \mathbf{B_i}(\varphi \wedge \neg\varphi) \quad \text{Löb's Theorem}$$

$$(3.5) \quad \mathbf{B_i}(\varphi \wedge \neg\varphi) \quad (4.3), (4.4)$$

□

Thus, with Belief Consistency, Löb's Theorem, and Theorem 5.3.2, it follows that:

Theorem 5.3.3 (Beliefs Inconsistent). *For all φ , $\mathbf{B_i} \varphi$.*

Proof. This follows from Theorem 5.3.2 and $\neg\mathbf{B_i}(\varphi \wedge \neg\varphi)$. □

QED

Note the similarity here to the problem facing PA^+ were it to prove its own consistency, due to Gödel's Second Incompleteness Theorem.[18]

We briefly note that the logic defined by these axiom schemas and inference rules avoids Löb's Obstacle for $\mathbf{K_i}$, as it is no longer has the positive introspection property.

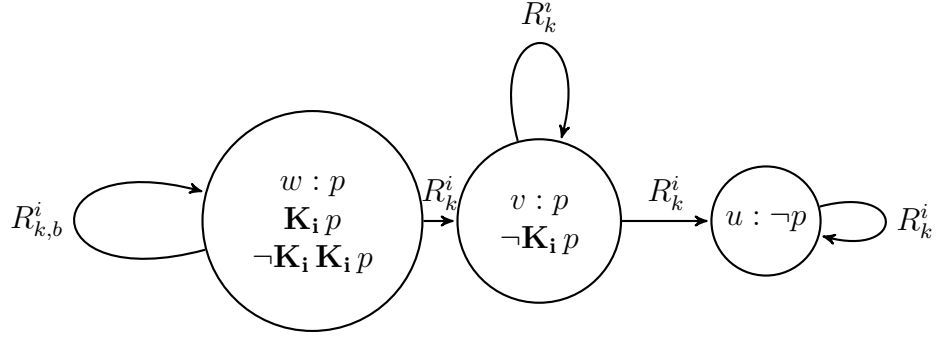


Figure 5.1: A counterexample to $\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{K}_i \varphi$.

Absent one of the three Löb Conditions, Löb's Theorem cannot be derived. Thus, we can say that the \mathbf{K}_i operator is Löb-safe, which is to say that the modal logic System T is Löb-safe. However, System T on its own is not adequate to model the knowledge and belief of human-like agents.

One might wonder whether it would be acceptable to abandon the Truth Axiom for knowledge and allow Löb's Theorem to hold for the knowledge operator in a way that avoids inconsistency. This would introduce more modesty to the notion of knowledge, where a human-like agent knows that her knowledge is true only for those propositions that she actually knows, but not in the general sense.⁷

What would this mean for epistemology? A false proposition would no longer imply a lack of knowledge, and the rejection of the truth axiom goes against the entire history of Western philosophical thought. This is to say, it would require robust philosophical defense, which we are not prepared to give here. Relaxing the Truth Axiom allows positive and negative introspection to live harmoniously with self-reference. We leave this for future work to explore. We note here, however, that work by Barasz *et. al.* in [7] show that agents with Löb's Theorem holding for their epistemic operators are able to cooperate in the Prisoner's Dilemma game.⁸ Related

⁷Smullyan referred to such agents as modest agents, for they are confident in the accuracy only of particular beliefs that they have good reasons (proofs) for, but lack a general confidence in their own beliefs.

⁸Specifically, they are programs capable of inspecting their own source code.

research in formal agent modeling via programs that play games with each other and can examine each other’s source code has been explored by Binmore, Howard, McAfee, Tennenholtz, and other game theorists. Barasz *et. al.* introduce the novel approach of using the so-called Gödel-Löb (\mathbb{GL}) modal logic of provability, presented in detail by Boolos in [11]. This logic consists of a \Box operator for “is provable in PA”, with the above Löb Theorem as an axiom, in addition to Axiom K (Distribution).

Rather than pursuing that route, we relax our axiom schema so as to avoid the derivation of Löb’s Theorem entirely, while maintaining the same or similar realism sought by the original *DASL* static schemas.

5.4 Avoiding Löb

Earlier, we noted that there is a mathematical gravity pulling agents toward Löb’s Obstacle due to the desiderata of rationality. Recall we described these desiderata as fundamentally *external* and *internal*. We seek rational agents that are accurate with respect to the external world, in that their beliefs (knowledge) is are true. We also seek rational agents that are highly introspectively accurate, able to draw inferences about their internal epistemic states. With this framework in mind, an approach to avoiding Löb’s Obstacle that adopts \mathbb{GL} as the epistemic operator, as in [7], eschews the external desiderata in favor of the internal desiderata. The agent is highly introspective, and we, as external observers, take it for granted that they are also externally accurate, though they themselves could not conclude this in general.

In order to avoid Löb’s Obstacle, we eschew the internal accuracy desideratum, and model agents who know only that their knowledge is accurate of the world, but not whether they know something or not, and have sufficient introspection as to know that they do not believe inconsistencies, and that they have strong evidence for their beliefs. This model is better suited to human-like agents who interact with reality

and lack perfect information about the environment.

We make the following relaxtions from \mathcal{DASL} 's static logic. First, we reduce the belief operator to a basic System D modality, consisting of $\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$. Second, we weaken Evidential Restraint from $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$ to $\mathbf{B}_i \Rightarrow \langle \mathbf{B}_i \rangle \mathbf{K}_i \varphi$, which we call Weakened Subjective Evidential Restraint. We make each of these moves because negative introspection and Evidential Restraint are each sufficient for deriving positive introspection for belief, which opens the door for Löb's Theorem, which leads to the Consistency Disaster proved in Theorem 5.3.2. The resulting logic avoids the Löbian Obstacle while retaining the inference from unsafe action to information about missing knowledge, the useful and interesting contrapositive to the standard game theoretic inference mentioned at the beginning of this thesis. Furthermore, it remains realistic for human-like agents, and remains in the class of logics whose completeness is established via Sahlqvist's Theorems. When this static logic serves as a foundation for \mathcal{DASL} 's dynamic extension, we refer to that logic as \mathcal{DASL}^- , or Weakened \mathcal{DASL} .

$\mathbf{K}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \psi)$	Distribution of \mathbf{K}_i
$\mathbf{K}_i \varphi \Rightarrow \varphi$	Truth
$\mathbf{B}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \psi)$	Distribution of \mathbf{B}_i
$\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$	Belief Consistency
$\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi$	Knowledge implies Belief
$\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \mathbf{K}_i \varphi$	Weak Subjective Evidential Restraint
From $\vdash \varphi$ and $\vdash \varphi \Rightarrow \psi$, infer $\vdash \psi$	Modus Ponens
From $\vdash \varphi$, infer $\vdash \mathbf{K}_i \varphi$	Necessitation of \mathbf{K}_i

Table 5.3: Static Logic of \mathcal{DASL}^-

Neither the belief operator nor the knowledge operator is susceptible to Löb's Obstacle, as Löb's Theorem is not derivable in the system.

In the counterexample below, we see that the belief operator lacks positive introspection.

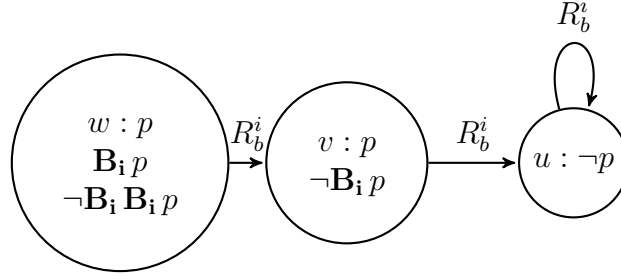


Figure 5.2: A counterexample to $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{B}_i \varphi$.

Without positive belief introspection, Löb’s Theorem for belief is no longer derivable, and therefore this (very weak) epistemic logic avoids Löb’s Obstacle. The resulting logic is still in the Sahlqvist class, so it is sound and complete.

```

Theorem weak_evidential_restraint_is_sahlqvist:
  forall (phi : prop) (a : DASL.Agents),
    sahlqvist_formula
      (SB a (SProp phi) =s⇒ \ (SB a (\ (SK a (SProp phi))))).
Proof.
  intros; sahlqvist_reduce.
Qed.

```

We provide the following counterexample to Löb’s Theorem.

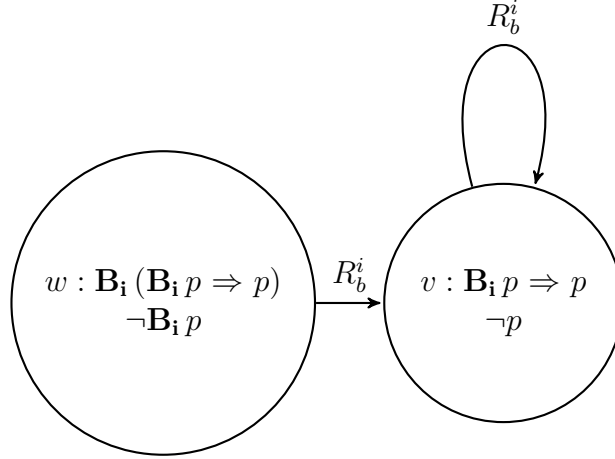


Figure 5.3: A counterexample to $\mathbf{B}_i(\mathbf{B}_i \varphi \Rightarrow \varphi) \Rightarrow \mathbf{B}_i \varphi$.

We assume in world w that $\mathbf{B}_i(\mathbf{B}_i p \Rightarrow p)$ holds. From this it follows that in all doxastically accessible worlds, *e.g.* v , $\mathbf{B}_i p \Rightarrow p$ holds. This holds for v when it has reflexive access only to itself, and $\neg p$ is the case. Because $\neg p$ is the case v , $\neg \mathbf{B}_i p$ is the case at w , concluding the counterexample.

The logic also includes a sort of weakened positive introspection theorem about knowledge, from the Knowledge implies Belief axiom and Weak Subjective Evidential Restraint. It is $\mathbf{K}_i \varphi \Rightarrow \langle \mathbf{K}_i \rangle \mathbf{K}_i \varphi$, which we read as “ i knows that φ only if she has objectively available evidence that she knows φ ”. This seems intuitive, and does not allow Löb’s theorem to destroy the integrity of knowledge. It is perhaps a satisfying compromise for those who find positive introspection about knowledge to be intuitive. The contrapositive of this weak positive introspection formula, $\mathbf{K}_i \neg \mathbf{K}_i \varphi \Rightarrow \neg \mathbf{K}_i \varphi$ is an instance of the T axiom, so it turns out to have been a theorem all along anyway, for any epistemic logic with the Truth axiom for knowledge.

This logic allows us to make the key inference discussed in this thesis, when it serves as a static foundation for weakened \mathcal{DASL} .

Theorem 5.4.1 (Weakened \mathcal{DASL} Inference of Safety-Critical Information).

$$(\langle \mathbf{A}, a \rangle_i \text{true} \wedge \neg \text{pres}(a)) \Rightarrow \neg \mathbf{K}_i \text{pres}(a) \wedge \neg \mathbf{K}_i \neg \mathbf{K}_i \text{pres}(a).$$

Proof. Assuming $w \models \langle \mathbf{A}, a \rangle_i \text{true}$, it follows that $\mathbf{B}_i \text{pres}(a)$, from minimum rationality. From weak subjective evidential restraint, it follows that $\langle \mathbf{B}_i \rangle \mathbf{K}_i \text{pres}(a)$. As an instance of Knowledge implies Belief, this implies $\langle \mathbf{K}_i \rangle \mathbf{K}_i \text{pres}(a)$, which is equivalent to $\neg \mathbf{K}_i \neg \mathbf{K}_i \text{pres}(a)$. From $\neg \text{pres}(a)$, it follows from the Truth axiom that $\neg \mathbf{K}_i \text{pres}(a)$. \square

Thus, the application results established in Chapter 4 still apply, and likewise the contrapositive of the standard game theoretic inference still holds, from bad action to bad knowledge. We turn now to some philosophical considerations.

5.4.1 Philosophical Considerations

It is worth circling back to the philosophical considerations presented in Chapter 3 and assessing the effects of this weakening move. We must assess whether the static base of \mathcal{DASL}^- satisfies them, and if not, acknowledge it here. The philosophical desiderata were: 1) An account of justification, 2) Avoidance of skepticism for humans, 3) the formal properties of soundness, completeness, indistinguishability formalized, and avoid a collapse of belief and knowledge, and 4) handle Moore's Paradox.

Justification. Though we have altered our axiom schema, we have not altered our interpretation of the epistemic and doxastic relations. An epistemic relation between two possible worlds is still based on objectively available evidence, and a doxastic relation is one based on subjectively available evidence. The notion of Weakened Subjective Epistemic Restraint remains internalist about justification. We still assume our agents hold only justified beliefs, and this is captured with Weakened Subjective Evidential Restraint, where an agent believes something only if she has subjectively available evidence that she knows it, *i.e.* that it is true, she believes it, and she is justified. Lacking powerful internal accuracy, it may be more difficult for such agents

to form beliefs. The axiom schemas, in conjunction with the K axiom schemas, grant us some internal power, though. Agents know that their beliefs are consistent, they know that knowledge is factual, and they know that their beliefs are justified. They do not necessarily know what they believe, but they have subjective evidence that they believe something, when they do.

An alternative approach is to use $\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{K}_i \rangle \mathbf{K}_i \varphi$, which might we called Weakened Evidential Restraint. This is actually weaker than the Weakened Subjective Evidential Restraint property.

Theorem 5.4.2 (Weakened Subjective ER Implies Weakened ER).

$$(\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \mathbf{K}_i \varphi) \Rightarrow (\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{K}_i \rangle \mathbf{K}_i \varphi).$$

Proof. It suffices to show that $\langle \mathbf{B}_i \rangle \mathbf{K}_i \varphi \Rightarrow \langle \mathbf{K}_i \rangle \mathbf{K}_i \varphi$. This is an instance of Knowledge implies Belief. □

Adopting this weaker property would be consistent with externalist theories of justification, which we eschew, being convinced by arguments for internalism.

Skepticism. We clearly do not present an account of knowledge that falls victim to the skeptical hypothesis argument, because we avoid the problematic properties of knowledge that do not hold for human-like agents, *i.e.* positive introspection. Regarding the skeptical argument that justification is impossible, we maintain our evidentialist interpretation of the relations on possible worlds, and continue with our foundationalist version of evidentialism where noticed perceptions of evidence can constitute justification, in addition to beliefs.

Indistinguishability. The static base of \mathcal{DASL}^- has caused us to lose our notion of doxastic capture defined in Theorem 3.3.1. This formal property provided a sense of indistinguishability proved in Theorem 3.3.2. Our choice of external accuracy over internal accuracy seems to have cost us this desideratum. However, it is worth noting that no epistemic logic can achieve indistinguishability through an equivalence rela-

tion while also avoiding Löb's Obstacle. An equivalence relation entails reflexivity and transitivity, which in conjunction with self-reflection yield Löb's Theorem and inconsistency. Thus, it may be the case that indistinguishability cannot be achieved through an equivalence relation. For now, we take it as philosophical debt that we must further explore indistinguishability for self-reflective agents.

Collapse. The static base of \mathcal{DASL}^- similarly avoids the collapse of belief and knowledge. The same counterexample serves, where agent i falsely believes φ , and thus does not know it.

Moore's Paradox. We maintain the ability to model assertions as actions that presuppose beliefs, thereby formalizing Moore's sentence, " p , but I do not know that p ," as $\mathbf{B}_i p \wedge \mathbf{B}_i \neg \mathbf{K}_i p$. Weakened Subjective Evidential Restraint allows us to infer $\langle \mathbf{B}_i \rangle \mathbf{K}_i p$ from the first conjunct, which is equivalent to $\neg \mathbf{B}_i \neg \mathbf{K}_i p$, a direct contradiction to the second conjunct, explaining the intuitive paradox of Moore's sentence.

Soundness and Completeness. We have already shown that \mathcal{DASL}^- is complete, but we have not specified the frame condition corresponding to Weakened Subjective Evidential Restraint. We do so now, by applying the Sahlqvist-van Benthem Algorithm

$$\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \mathbf{K}_i \varphi.$$

$$\begin{aligned} &\rightsquigarrow \forall P, y, (R_b^i(x, y) \Rightarrow P(y)) \Rightarrow \exists z, \forall z', (R_b^i(x, z) \cdot (R_k^i(z, z') \Rightarrow P(z'))) \\ &\rightsquigarrow \forall y, (R_b^i(x, y) \Rightarrow \lambda u. (R_b^i(x, u))(y)) \\ &\quad \Rightarrow \exists z, \forall z', (R_b^i(x, z) \cdot (R_k^i(z, z') \Rightarrow \lambda u. (R_b^i(x, u))(z'))) \\ &\rightsquigarrow \forall y, (R_b^i(x, y)) \Rightarrow \exists z, \forall z', (R_b^i(x, z) \cdot (R_k^i(z, z') \Rightarrow (R_b^i(x, z')))) \end{aligned}$$

It is an unfamiliar frame condition, but it is a first order formula, and because SCorT establishes correspondence between the modal formula and the first order formula, it entails that the class of frames defined by the first order formula as a

frame condition are those for which the modal formula are sound.

5.5 Summary

This chapter has presented a problem facing agents capable of self-reflective reasoning that gives rise to logical problems of self-reference, which are present in the Peano Arithmetic. Thus, any agent whose reasoning is formally stronger than PA must worry about this problem, named Löb’s Obstacle.

We explained how the Löb Conditions (Axiom K, Axiom 4, and Rule of Necessitation), in conjunction with self-reflective sentences, derive Löb’s Theorem. We showed that standard assumptions about knowledge and belief result in the antecedent of Löb’s Theorem being a theorem itself, for both knowledge operators and belief operators, which results in inconsistency. \mathcal{DASL} ’s static base suffers from this problem.

We responded to this by presenting a weakened \mathcal{DASL} , which we named \mathcal{DASL}^- , that keeps most of the strengths of \mathcal{DASL} while avoiding Löb’s Obstacle. The one desideratum we lost in \mathcal{DASL}^- was a formal account of indistinguishability, which is intuitively powerful but not formally required, and generally lacks philosophical defense. Furthermore, we conjecture that any modal system with a satisfactory account of indistinguishability will run into Löb’s Obstacle.

Chapter 6

Summary and concluding remarks

In this work, we have presented a dynamic modal logic for reasoning about the relationship between an agent's beliefs, knowledge, and action, with a distinction between mere action and safe action expressible. We named this logic Dynamic Agent Safety Logic, or *DASL*. Including this many modalities in one logic makes it very expressible and suitable for modelling a variety of realistic situations involving human-like agents. Additionally, the logic was developed with careful concern for realism. This departs from the typical thesis of modern modal logic, which develops a logic with advanced formal properties without much regard for the realism of its foundational axiom schemas. By carefully considering the philosophical implications of our static base, the thesis made advances in reconnecting the worlds of epistemology and epistemic logic, and philosophy and modal logic more generally.

We paid careful attention to the philosophical commitments our epistemic logic carried with it, and worked to develop a philosophically defensible static foundation. This required us to adopt philosophically defensible positions against skepticism and of justification, which we did. We specified an interpretation of the possibility relations that conforms well to the predominant epistemic theory known as evidentialism, with a foundationalist structure to justification. In this sense, the thesis makes ad-

vances as a formal theory of epistemology, though it is certainly not complete. More work needs to be done in connecting the formalized interpretations to the system itself. One idea for doing this is to build a simple AI system with sensors and actuators, which abides by the axioms of *DASL* in its reasoning, and constructs models based on objective and subjectively available evidence. This is beyond the scope of the thesis, however.

As a thesis of computer science and logic, we took care to use computational tools to validate our logical theses. We used the Coq Proof Assistant to mechanically check our proofs of soundness and completeness to a large degree, and similarly mechanized our case studies which validate *DASL*'s application to aviation safety.

The proof of completeness mechanized a powerful aspect of modal logic from Sahlqvist and van Benthem. This dissertation, as far as we know, marks the first use of Coq for proving a modal logic's completeness via Sahlqvist theorems. To do this, we developed inductive predicates over the structure of axiom schemas corresponding to the definitions in Blackburn *et. al.*[20]. This two-leveled approach to mechanization allowed us to prove theorems about the logic, and then drop down into the object language and prove theorems in the logic. This approach is extensible to other domains, where an object language with different atoms could be defined in Coq, while the upper level schema mechanization remains fixed, along with its results.

We mechanically checked our formalization of three cases of aviation mishaps, zeroing in on particular moments from the events where *DASL* enables an inference from action to ignorance about safety-critical information. This inference was foreshadowed by the foundational models of agency in classical game theory, which infer good actions from assumptions of rationality and knowledge. A boolean manipulation of that classical inference yields the inference from bad action and assumption of rationality to an absence of knowledge. *DASL* provides the formal logic for capturing this inference in a rich and realistic fashion.

Finally, we confronted a formal problem facing epistemic logics with self-reflective agents. The problem is related to Löb’s Theorem in mathematical logic, formalized in provability logic by Boolos, and named Löb’s Obstacle by researchers in formal agent foundations. As noted by Smullyan, self-reflective agents believing things about themselves are formally identical to mathematical systems, *e.g.* Peano Arithmetic, proving statements about themselves. For contexts in which self-reflective agents should be explicitly modeled, as in AI safety research or realistic models of human reasoning, Löb’s Obstacle must be addressed. We show that the static foundation for \mathcal{DASL} , in addition to the other dominant epistemic logics, suffer from Löb’s Obstacle, and are therefore inadequate. We respond by weakening \mathcal{DASL} ’s base to \mathcal{DASL}^- , and retain all of \mathcal{DASL} ’s strengths but the philosophically desirable feature of indistinguishability for the possibility relations. Further work must be done to explore this weakened \mathcal{DASL} and alternative ways to avoid Löb’s Obstacle, but they are beyond the scope of this thesis.

Bibliography

- [1] Artemov, S.; Nogina, E. ntroducing Justification into Epistemic Logic. *Journal of Logic and Computation* 15. 1059-1073. (2005).
- [2] Federal Aviation Administration
website: <https://www.faa.gov/files/gslac/library/documents/2012/Nov/71574/DirtyDozenWeb3.pdf>,
retrieved 10/5/2016.
- [3] Rankin, William. “MEDA Investigation Process”, Aero Quarterly.
Website: Boeing.com/commercial/aeromagazine/articles/qtr_2_07/aero_q207_article3.pdf,
retrieved 10/5/2016
- [4] Aviation Maintenance Technician Handbook. Chapter 14, “Addendum/Human Factors”.
- [5] van Ditmarsch, H.; van der Hoek, W.; Kooi, B. Dynamic Epistemic Logic. Springer. Dordrecht, The Netherlands, 2008.
- [6] Ahrenbach, S.; Goodloe, A. Formal Analysis of Pilot Error Using Agent Safety Logic. *Innovations in Systems and Software Engineering*. **Submitted**.
- [7] Barasz, M.; Christiano, P.; Fallenstein, B.; Herreshoff, M.; LaVictoire, P.; Yudkowsky, E. “Robust Cooperation in the Prisoner’s Dilemma: Program Equilibrium via Provability Logic”, <https://arxiv.org/abs/1401.5577v1> Accessed 7/9/19.

- [8] Bolton, M.; Bass, E.; Siminiceanu, R. Using formal verification to evaluate human-automation interaction: A review. *Systems, Man, and Cybernetics: Systems, IEEE Transactions on* (Vol 43, No. 5, pp.488-503). 2013. IEEE.
- [9] Bolton, M.; Bass, E. Formally verifying human-automation interaction as part of a system model: limitations and tradeoffs. *Innovations in systems and software engineering* (Vol 6, No. 3, pp.419-231). 2010. Springer-Verlag.
- [10] Bolton, M.; Bass, E.; Siminiceanu, R. A systematic approach to model checking human-automation interaction using task-analytic models. *Systems, Man, and Cybernetics: Systems and Humans, IEEE Transactions on* (Vol 41, No. 5, pp.961-976). 2011. IEEE.
- [11] Boolos, G. The Logic of Provability. Cambridge University Press. New York and Cambridge, 1993.
- [12] Bostrom, N. Superintelligence: Paths, Dangers, Strategies. Oxford University Press. Oxford, 2014.
- [13] Barras, B.; Boutin, S.; Cornes, C.; Courant, J.; Filliatre, J.C.; Gimenez, E.; Herbelin, H.; Huet, G.; Munoz, C.; Murthy, C. and Parent, C. The Coq proof assistant reference manual: Version 6.1 (Doctoral dissertation, Inria). 1997.
- [14] De Wind, P. "Modal logic in Coq." M.Sc. University of Amsterdam. 2001.
- [15] Feldman, R.; Conee, E. Evidentialism. *Philosophical Studies*. (Vol. 48 (1): pp 15-34). 1985.
- [16] Gettier, E. Is Justified True Belief Knowledge? *Analysis* (Vol. 23, pp. 121-123). 1963.

- [17] Fallenstein, B. Soares, N. “Problems of self-reference in self-improving space-time embedded intelligence.” *Lecture Notes in Artificial Intelligence 8598*, pages 21–32. Springer, 2014.
- [18] Gödel, K. ”Über Formal Unentscheidbare Sätze der Principia Mathematica und Verwandter Systeme, I.” *Monatshefte für Math. u. Physik* 38, 173-198, 1931.
- [19] Prichard, H.A., 1950. *Knowledge and Perception*, Oxford: Clarendon Press.
- [20] Blackburn, P.; de Rijke, M.; Venema, Y. *Modal Logic*. Cambridge University Press. New York, 2001.
- [21] Baltag, A.; Moss, L.; Solecki, S. ”The logic of public announcements, common knowledge, and private suspicions.” *Readings in Formal Epistemology*. Springer International Publishing, 2016. 773-812.
- [22] Fagin, R.; Halpern J.; Moses, Y.; Vardi, M. *Reasoning about knowledge*. MIT press, 2004.
- [23] Basin, D., Radomirovic, S. and Schmid, L., *Modeling Human Errors in Security Protocols*.
- [24] Butler, R.W., Miller, S.P., Potts, J.N. and Carreno, V.A., 1998. A formal methods approach to the analysis of mode confusion. In *Digital Avionics Systems Conference, 1998. Proceedings., 17th DASC. The AIAA/IEEE/SAE (Vol. 1, pp. C41-1)*. IEEE.
- [25] Hoare, C. A. R. ”An axiomatic basis for computer programming.” *Communications of the ACM* 12 (10): 576-583. (1969).
- [26] Butler, R.; Johnson, S. *Formal Methods For Life-Critical Software*. *Computing in Aerospace 9 Conference*, (pp. 319-329), San Diego, CA. October 1993.

- [27] LaVictoire, P.; Fallenstein, B.; Yudkowsky, E.; Barasz, M.; Christiano, P.; Herreshoff, M. “Program equilibrium in the prisoner’s dilemma via Löb’s Theorem.” *Multiagent Interaction without Prior Coordination: Papers from the AAAI-14 Workshop*. AAAI Publications, 2014.
- [28] Löb, M. Solution of a Problem of Leon Henkin. *Journal of Symbolic Logic* 20 (2), 115 - 118. (1955).
- [29] Davis, M.; Putnam, H. ”A Computing Procedure for Quantification Theory”. J.ACM. 7 (3): 201–215. (1960).
- [30] Rushby, J., 1993. Formal methods and the certification of critical systems. SRI International. Computer Science Laboratory.
- [31] Aucher, G. and Schwarzenruber, F., 2013. On the complexity of dynamic epistemic logic. arXiv preprint arXiv:1310.6406.
- [32] et d’Analyses, Bureau d’Enquêtes. “Final report on the accident on 1st June 2009 to the Airbus A330-203 registered F-GZCP operated by Air France flight AF 447 Rio de Janeiro–Paris.” Paris: BEA (2012).
- [33] van Benthem, J.; Pacuit, E; Roy, O. Toward a Theory of Play: A Logical Perspective on Games and Interaction. *GAMES* **2011**, 2, 52-86.
- [34] Prior, A. Time and Modality. Clarendon Press: Oxford, UK, 1957.
- [35] Kohler, R. Lords of the Fly. The University of Chicago Press: Chicago, USA, 1994.
- [36] Kraus, S.; Lehmann D. Knowledge, Belief and Time. *Theoretical Computer Science* **1988**, 58: 155-174.
- [37] Kripke, S. Semantical Considerations on Modal Logic. *Acta Philosophica Fennica* **1963**, 16: 83-94.

- [38] Halpern, J. Y.; Moses, Y. Towards a Theory of Knowledge and Ignorance: Preliminary Report. In: Apt K.R. (eds) Logics and Models of Concurrent Systems. NATO ASI Series (Series F: Computer and Systems Sciences), vol 13. Springer, Berlin, Heidelberg
- [39] Harel, D.; Kozen, D; Tiuryn, J. Dynamic Logic. MIT Press: Cambridge, MA, USA, 2000.
- [40] Hsu, F. Behind Deep Blue: Building the Computer that Defeated the World Chess Champion. Princeton University Press: Princeton, NJ, USA, 2004.
- [41] Lenzen, W. Recent work in epistemic logic. *Acta Philosophica Fennica* **1972**, *30(1)*: 1-219.
- [42] Meyer, J.J. Ch.; van der Hoek, W. Epistemic Logic for Artificial Intelligence. *Cambridge Tracts in Theoretical Computer Science*. Cambridge University Press: Cambridge, UK, 1995.
- [43] Moore, G.E. "Moore's Paradox. In Baldwin, T. *G.E. Moore: Selected Writings*. London: Routledge **1993**.
- [44] Negri, S. Kripke completeness revisited. *Acts of Knowledge: History, Philosophy and Logic: Essays Dedicated to Göran Sundholm* **2009**, 247-282.
- [45] Pnueli, A. The Temporal Logic of Programs. *Proceedings of the 18th Annual Symposium on Foundations of Computer Science (SFCS 77)*. IEEE Computer Society.
- [46] Rushby, J. Using model checking to help discover mode confusions and other automation surprises. *Reliability Engineering and System Safety* **2002**, *75*, 167-177.

- [47] Rushby, J. On the Interpretation of Assurance Case Arguments. *Proceedings of the Second International Workshop on Argument for Agreement and Assurance (AAA 2015)*. Springer LNCS.
- [48] Rushby, J. *Mechanized Analysis of a Formalization of Anselm's Ontological Argument by Eder and Ramharter*. CSL Technical Note, Menlo Park, CA, February 11, 2016.
- [49] Rushby, J. Trustworthy Self-Integrating Systems. *Springer LNCS*, **2016**, 9581, 19-29.
- [50] Gelman, G.; Feigh, K.; Rushby, J. Example of a Complementary use of Model Checking and Agent-based Simulation. *IEEE Int'l Conf. on Systems, Man, and Cybernetics* **2013**, 900-905.
- [51] Baltag, A.; Moss, L.; Solecki, S. The Logic of Public Announcements, Common Knowledge, and Private Suspicions. *TARK '98*, **1998**, 43-56.
- [52] Shoham, Y.; Leyton-Brown, K. *Multiagent Systems: Algorithmic, Game Theoretic and Logical Foundation*. Cambridge University Press: New York, USA, 2008.
- [53] van Ditmarsch, H.P.; van der Hoek, W.; Kooi, B.P. Dynamic Epistemic Logic with Assignment. *Proceedings of the fourth internal joint conference on Autonomous agents and multiagent systems*, **2005**, 141-148.
- [54] van Benthem, J. *Logical Dynamics of Information and Interaction*. Cambridge University Press: New York, USA, 2011.
- [55] van Benthem, J. Modal Logic for Open Minds. Center for the Study of Language and Information, 2010.
- [56] Van Ditmarsch, H. "Knowledge games." *Bulletin of Economic Research* 53.4 (2001): 249-273.

- [57] Fagin, R.; Halpern J.; Moses, Y.; Vardi, M. *Reasoning About Knowledge*. The MIT Press: Cambridge, USA, 2003.
- [58] Harel, D.; Kozen, D. and Jerzy Tiuryn. *Dynamic logic*. MIT press, 2000.
- [59] Hintikka, J. *Knowledge and Belief*. Cornell University Press. Ithaca, USA, 1962.
- [60] Horty, J. *Agency and Deontic Logic*. Oxford University Press. Oxford, UK, 2009.
- [61] Barwise, J.; Seligman, J. *Information Flow: The Logic of Distributed Systems*. Cambridge University Press: New York, USA, 1997.
- [62] van Ditmarsch, H.; van der Hoek, W.; Kooi, B. *Dynamic Epistemic Logic*. Springer: Dordrecht, The Netherlands, 2008.
- [63] Bass, E.; Feigh K.; Gunter, E.; Rushby, J. Formal Modeling and Analysis for Interactive Hybrid Systems. In *Proceedings of the Fourth International Workshop on Formal Methods for Interactive Systems* (FMIS 2011). Electronic Communications of the EASST. Vol 45 (2011).
- [64] Rushby, J. Formalism in Safety Cases. In *Making Systems Safer: Proceedings of the 18th Safety-Critical Systems Symposium*; Dale, C.; Anderson, T.; Eds. Bristol UK, 2010, 3-17.
- [65] Klir, G.; Yuan, B. *Fuzzy sets and fuzzy logic*. Vol. 4. New Jersey: Prentice hall, 1995.
- [66] Lescanne, P. "Mechanizing common knowledge logic using COQ." *Annals of Mathematics and Artificial Intelligence* 48.1-2 (2006): 15-43. APA
- [67] De Moura, L.; Bjørner, N. Z3: An efficient SMT solver. *Tools and Algorithms for the Construction and Analysis of Systems*. 2008. Springer.

- [68] Maliković, M.; Čubrilo, M. "Modeling epistemic actions in dynamic epistemic logic using Coq." CECIS-2010. 2010.
- [69] Maliković, M.; Čubrilo, M. "Reasoning about Epistemic Actions and Knowledge in Multi-agent Systems using Coq." Computer Technology and Application 2.8 (2011): 616-627.
- [70] National Transportation Safety Board. "Descent below visual glidepath and impact with Seawall Asiana Flight 214, Boeing 777-200ER, HL 7742, San Francisco, California, July 6, 2013 (Aircraft Accident Report NTSB/AAR-14/01)." Washington, DC Author: NTSB (2014).
- [71] Palmer, Bill. Understanding Air France 447. William Palmer, 2013.
- [72] Lescanne, P.; Puisségur, J. "Dynamic logic of common knowledge in a proof assistant." arXiv preprint arXiv:0712.3146 (2007).
- [73] Aumann, R. "Interactive epistemology I: knowledge." International Journal of Game Theory 28.3 (1999): 263-300.
- [74] Hwang, M.; Lin, J. "Information dimension, information overload and decision quality." Journal of Informatin Science 25 (1999): 213-218.
- [75] Rescher, Nicholas. "Epistemic Logic." University of Pittsburgh Press. 2005.
- [76] Sahlqvist, Henrik. "Completeness and Correspondence in the First and Second Order Semantics for Modal Logic." In *Proc. of the Third Scandinavian Logic Symposium*; Kanger, S.; Ed. Oslo, Norway, 1975, 110 - 143.
- [77] Simpson, C.; Prusak, L. "Troubles with Information Overload = Moving from Quantity to Quality in Information Provision." International Journal of Information Management 15 (1995): 413-425.

- [78] Smullyan, R. “Logicians who reason about themselves” In *Proceedings of the 1986 conference on Theoretical aspects of reasoning about knowledge*, Morgan Kaufmann Publishers Inc., San Francisco (CA), 1986, 341–352.
- [79] Verbrugge, Rineke (L.C.) (5 April 2017). ”Provability Logic”. The Stanford Encyclopedia of Philosophy. Retrieved 18 June 2019.
- [80] von Wright, G.H.. “An Essay on Modal Logic.” Amsterdam: North-Holland Publishing Company. 1951.
- [81] Williamson, T. Modal Logic as Metaphysics. Oxford University Press: New York, USA, 2013.
- [82] Wittgenstein, L. “Philosophical Investigations.” Blackwell Publishers. 1953.
- [83] Yudkowski, E.; Hereshoff, M. “Tiling Agents for Self-Modifying AI, and the Löbian Obstacle.” <https://intelligence.org/files/TilingAgentsDraft.pdf>, 2013.

VITA

This is a summary of your *professional* life, and should be written appropriately. This can be written in the following order: where you were born, what undergraduate university you graduated from, if you received a masters, and which institution you graduated from with your PhD (University of Missouri). You can describe when you began research with your current advisor.

In another paragraph, you could say if/when you were married, what the name of your kids are, and what your plans are for after graduation if you choose. Take a look at other vita's from other dissertations for examples.