

FORMAL METHODS WITH DYNAMIC AGENT SAFETY LOGIC

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. Guilherme DeSouza

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	12
2.2 Modal Logic	15
2.2.1 Modal Logic Syntax and Semantics	18
2.3 Epistemic Logic	21
2.3.1 Why S5 Is Appealing	23
2.4 Doxastic Logic	24
2.4.1 Dynamic Epistemic Logic	25
2.5 Formal Methods Tools	27
3 Dynamic Agent Safety Logic	28
3.1 Syntax and Semantics	29
3.1.1 Syntax	29
3.1.2 Metalanguage	30
3.1.3 Semantics	34
3.1.4 Hilbert System	35
3.1.5 Static Base	36

3.2	Soundness and Completeness	57
3.2.1	Soundness	58
3.2.2	Completeness	62
3.2.3	Completeness	69
4	Case Studies	76
4.1	Case Study and Mechanization	76
4.1.1	Air France 447	77
4.1.2	Mechanization in Coq	80
4.2	Additional Case Studies	90
4.2.1	Copa 201 and Asiana 214	90
4.3	Safety Properties	93
4.4	Decision Problem	94
5	Ensuring Delivery of Safety-Critical Information	97
5.1	SMT Solving and the Unsatisfiable Core	99
5.2	Encoding and Case Studies	100
5.2.1	Air France 447 Encoded	102
5.2.2	Copa Flight 201 Encoded	106
5.2.3	Asiana Flight 214 Encoded	109
5.3	The Unsatisfiable Core is the Safety-Critical Information	112
5.4	Application Design	114
5.5	Related Work	115
6	Summary and concluding remarks	116
	APPENDIX	117
A	Title of first appendix	117
A.1	Section title	117

A.1.1 Subsection title	117
BIBLIOGRAPHY	118
VITA	125

LIST OF TABLES

Table		Page
3.1	Hilbert System of Dynamic Agent Safety Logic	36
3.2	Static Base of Dynamic Agent Safety Logic	38
3.3	The reduction axioms of \mathcal{DASL}	57
3.4	Equivalent Static Base with Sahlqvist-friendly Negative Belief Intro- spection.	66

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 ”.	27
2.4	$M2$: The model after Alice announces “ φ ”.	27
3.1	Action Model for “Alice flips a coin”	31
3.2	Action Model for “Alice peeks”	31
3.3	Action Model for “Alice secretly peeks”	32

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 [5]. 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 [24, 27]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. This can be done at run-time, as demonstrated herein by an encoding of the problem space for use by the Z3 theorem prover, which computes the missing information. We do not propose that a future avionics monitor actually run the Z3 tool on the aircraft, but rather use it here as a proof of concept, with implementation details left for future work.

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 [22], they develop a modal logic for reasoning about the knowledge and decision-making of agents in games, and in [39], 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 [23], a graph theoretic semantics by Saul Kripke [25], and logics for belief and knowledge due to work by Rescher [60], von Wright [63], and Hintikka [44]. 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 [12], 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 $\Diamond_R p \Rightarrow p$. 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 [58], 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.

The final foundational school of modal logic comes 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* [31]. 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 [26].

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 [45]. For the

dynamic approach, see van Ditmarsch *et al.* [47].

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 present *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 uses the Z3 Theorem Prover’s Satisfiability Modulo Theory (SMT) solving capabilities to encode the aviation mishaps and demonstrate the inference of the missing safety-critical information.

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 Dynamic Epistemic Logic, upon which the logic presented in this thesis is based, with axioms corresponding to the more realistic model of agency we develop. Section 2.5 describes the formal tools we use to mechanically check that \mathcal{DASL} is sound and complete, that its application to aviation safety is without errors, and to encode the inference of safety-critical information. Coq is both a programming language and a tool for verifying proofs, based on the Calculus of Inductive Constructions. We explain the basics of each of these to the reader in this section, since Coq is the tool we use for proof and application checking. Z3 is a theorem prover and model checker, which interact with solely for its SMT solving abilities. We encode the actual instrument readings as a first order formula, and we encode the safety precondition of an action as formulas comprising the constraining theory. We then check whether the formula is satisfiable modulo that theory. This tells us whether the instrument readings are consistent with the safety preconditions for an action, and if not, where the conflict lies. There is some additional complexity here, but we describe this in more detail in Chapter 5.

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.

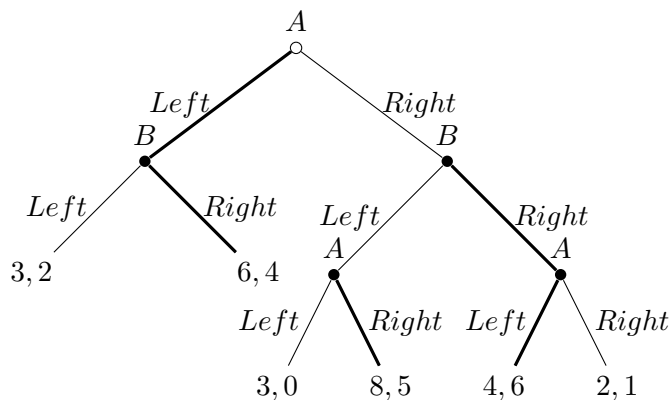


Figure 2.1: A game between players A and B

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* [22]. 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 [58].

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 [41]. 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 philo-

sophical interpretation, and expanding to epistemic and doxastic logic. Then, we introduce dynamic logic, and its expansion into Dynamic Epistemic Logic and Public Announcement Logic.

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 [40]. 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 [64]. 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 [23, 25]. 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 [25]. 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 [12]. 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} :

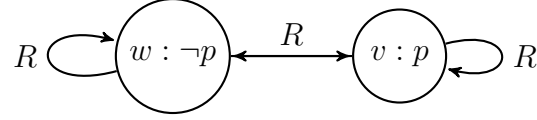


Figure 2.2: \mathcal{M} : A simple counterexample using possible world semantics.

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:

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

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

$$\forall x, y, z \ R(x, y) \wedge R(y, z) \Rightarrow R(x, z) \quad (2.3)$$

$$\Box\varphi \Rightarrow \Box\Box\varphi \quad (2.4)$$

- 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.*[42], we impose the frame conditions

of reflexivity and Euclidicity (Euclidean), or one that is transitive and reflexive, as in Hintikka [44]. 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.

2.3 Epistemic Logic

Epistemic logic began with philosophical concerns about knowledge [42, 44, 60]. 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 epistemic 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.

2.3.1 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.

2.4 Doxastic Logic

This section presents doxastic logic, the logic for reasoning about belief. Hintikka’s formalization in [44] 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.

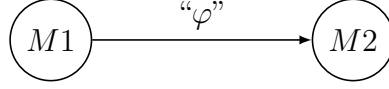
Knowledge involves a condition that is external to the agent's direct perception: the truth or falsity of the proposition. Belief, on the other hand, is entirely internal to the agent's perception. She can introspect on her evidence and deliberate on a proposition, and this is sufficient to see that 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.

2.4.1 Dynamic Epistemic Logic

Dynamic Epistemic Logic (DEL) formalizes situations in which agents' epistemic states change over time, due to announcements or other informational events. [47] 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 $\text{"}p\text{"}$ 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.

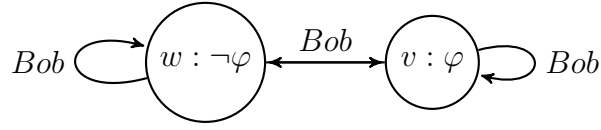


Figure 2.3: $M1$: The model before Alice announces $\text{"}\varphi\text{"}$.



Figure 2.4: $M2$: The model after Alice announces $\text{"}\varphi\text{"}$.

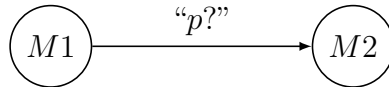
We are concerned with an additional element: the *safety* status of an action, and an agent’s knowledge and belief about that. To capture this, we extend DEL and call the new logic Dynamic Agent Safety Logic (\mathcal{DASL}), which we introduce in the next chapter. The next section lays out the state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

2.5 Formal Methods Tools

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 (DASL). The remainder of this section presents DASL’s syntax, semantics, and proves its soundness.

3.1 Syntax and Semantics

3.1.1 Syntax

The Dynamic Agent Safety Logic (DASL) used in this paper has the following syntax.

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

where $p \in AtProp$ is an atomic proposition, 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 $[\mathbf{A}, a]$ and $[\mathbf{A}, a]^S$ 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 $\langle \mathbf{K}_i \rangle$, $\langle \mathbf{B}_i \rangle$, $\langle \mathbf{A}, a \rangle$, and $\langle \mathbf{A}, a \rangle^S$ 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 relation

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 structures, called models, and action models, with functions for action pre- and post-conditions and the product update of Kripke and action models [47].

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 $\langle Actions, \chi_k^i, \chi_b^i, pre, pre_s, post, post_s \rangle$. It is a set of action tokens, sets of epistemic and doxastic relations on action tokens for agents, action pre- and post-conditions, and safe action pre- and post-conditions.

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- and post-conditions, and doxastic relations, for simplicity):

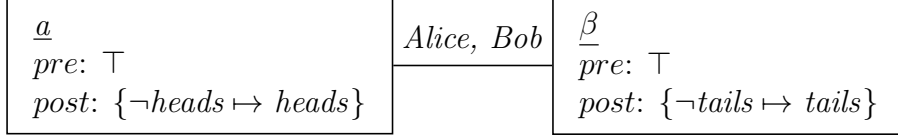


Figure 3.1: Action Model for “Alice flips a coin”

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:

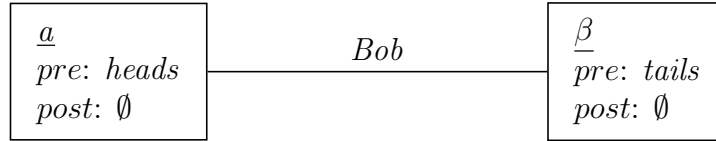


Figure 3.2: Action Model for “Alice peeks”

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:

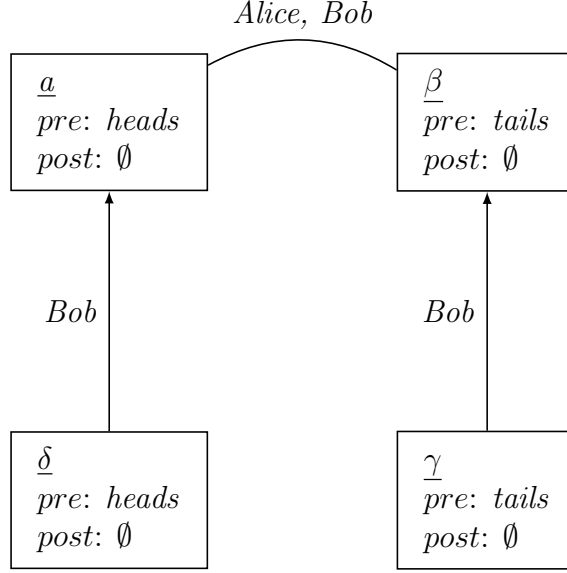


Figure 3.3: Action Model for “Alice secretly peeks”

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, $\llbracket (A, a) \rrbracket$ denotes a relation on Kripke model-world pairs. It represents the relation that holds between M, w and M', w' when 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.

Postcondition Function

The Postcondition function, $post :: A \times AtProp \mapsto AtProp$, takes an action structure and an atomic proposition, and maps to the corresponding atomic proposition after the action occurs.

$$post(A, p) = p \text{ if } update(M, A, w, a, i) \models p, \text{ else } \neg p.$$

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 executing the action, 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.

$$update(M, A, w, a, i) = (M', w') \text{ where :}$$

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. $V'(p) = post(A, p)$

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.

3.1.3 Semantics

The logic \mathcal{DASL} has the following Kripke semantics.

$$M, w \models p \Leftrightarrow w \in V(p)$$

$$M, w \models \neg\varphi \Leftrightarrow M, w \not\models \varphi$$

$$M, w \models \varphi \wedge \psi \Leftrightarrow M, w \models \varphi \text{ and } M, w \models \psi$$

$$M, w \models \mathbf{K}_i \varphi \Leftrightarrow \forall v, w R_k^i v \text{ implies } M, v \models \varphi$$

$$M, w \models \mathbf{B}_i \varphi \Leftrightarrow \forall v, w R_b^i v \text{ implies } M, v \models \varphi$$

$$M, w \models [\mathbf{A}, a]\varphi \Leftrightarrow \forall M', w', (M, w) \llbracket (A, a) \rrbracket (M', w')$$

$$\text{implies } M', w' \models \varphi$$

$$M, w \models [\mathbf{A}, a]^S \varphi \Leftrightarrow \forall M', w', (M, w) \llbracket (A, a) \rrbracket^S (M', w')$$

$$\text{implies } M', w' \models \varphi$$

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

$$(M, w) \llbracket (A, a) \rrbracket (M', w') \Leftrightarrow M, w \models \text{pre}(a)$$

$$\text{and } \text{update}(M, A, w, a, i) = (M', w')$$

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

$$\text{and } \text{update}(M, A, w, a, i) = (M', w')$$

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
$\mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \mathbf{B_i} \varphi$	Positive Belief Introspection
$\neg \mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \neg \mathbf{B_i} \varphi$	Negative Belief Introspection
$\mathbf{K_i} \varphi \Rightarrow \mathbf{B_i} \varphi$	Knowledge implies Belief
$\mathbf{B_i} \varphi \Rightarrow \mathbf{B_i} \mathbf{K_i} \varphi$	Confidence in Belief
$[\mathbf{A}, a]p \Leftrightarrow (post(a) \Rightarrow p)$	Atomic Consequence
$[\mathbf{A}, a]\neg\varphi \Leftrightarrow (pre(a) \Rightarrow \neg[\mathbf{A}, a]\varphi)$	Action Negation
$[\mathbf{A}, a](\varphi \wedge \psi) \Leftrightarrow ([\mathbf{A}, a]\varphi \wedge [\mathbf{A}, a]\psi)$	Action Conjunction
$[\mathbf{A}, a]\mathbf{K_i} \varphi \Leftrightarrow (pre(a) \Rightarrow \bigwedge_{a\chi_k^i\beta} \mathbf{K_i} [\mathbf{A}, \beta]\varphi)$	Action and Knowledge
$[\mathbf{A}, a][\mathbf{A}', a']\varphi \Leftrightarrow [(A, a); (A', a')]\varphi$	Action Composition
$[\mathbf{A}, a]^S p \Leftrightarrow (post_s(a) \Rightarrow p)$	Safe Atomic Consequence
$[\mathbf{A}, a]^S \neg\varphi \Leftrightarrow (pre_s(a) \Rightarrow \neg[\mathbf{A}, a]^S \varphi)$	Safe Action Negation
$[\mathbf{A}, a]^S (\varphi \wedge \psi) \Leftrightarrow ([\mathbf{A}, a]^S \varphi \wedge [\mathbf{A}, a]^S \psi)$	Safe Action Conjunction
$[\mathbf{A}, a]^S \mathbf{K_i} \varphi \Leftrightarrow (pre_s(a) \Rightarrow \bigwedge_{a\chi_k^i\beta} \mathbf{K_i} [\mathbf{A}, \beta]^S \varphi)$	Safe Action and Knowledge
$[\mathbf{A}, a]^S [\mathbf{A}', a']^S \varphi \Leftrightarrow [(A, a); (A', a')]\varphi$	Safe Action Composition
$[\mathbf{A}, a]\varphi \Rightarrow [\mathbf{A}, a]^S \varphi$	Inevitability
$\langle \mathbf{A}, a \rangle \varphi \Rightarrow \mathbf{B_i} \langle \mathbf{A}, a \rangle^S \varphi$	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 $\vdash [\mathbf{A}, a]\varphi$	Necessitation of $[\mathbf{A}, a]$

Table 3.1: Hilbert System of Dynamic Agent Safety Logic

Above are the axioms characterizing the logic. 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 S5 operator,

described earlier in section 2.3. This axiom schema for the knowledge operator was first formalized by Aumann in [58] as a logical foundation for the ideally rational agents of game theory and decision theory. Thus, we can consider that axiom schema to be the top tier of agent rationality axioms, where the belief operator is not even necessary.

Lower tiers attempt to formalize axioms for *boundedly rational* agents. An ideally rational agent has no limit to her ability to reason about higher order iterations of interactive epistemology. For example, common knowledge is formalized among agents i and j in the following way, where $i, j \in G$, the group:

$$\mathbb{C}_G \varphi \equiv (\varphi \Rightarrow (\mathbf{K}_i \varphi \wedge \mathbf{K}_j \varphi \wedge \mathbb{C}_G(\mathbf{K}_i \varphi \wedge \mathbf{K}_j \varphi))). \quad (3.1)$$

Thus, an ideally rational agent can carry this to any arbitrary n^{th} order of i knows that j knows that i knows ... Herein, I refer to the n levels of interactive epistemic knowledge as epistemic depth. Clearly, humans cannot actively carry out this depth of interactive reasoning, and this shows experimentally in games whose solutions require more than two or three levels of depth, *e.g.*, the centipede game. Agents with bounded rationality have their epistemic depths formally bounded, to better reflect the limitations of human reasoning.

Bounding agents in this way allows for formal models that capture how a human agent might realistically deviate from optimal Nash equilibria of games. However, bounding rationality in this way does not allow for agents who act based on false beliefs. A lower tier rationality model might impose bounds on the epistemic depth of its agents and also allow for the agents to make decisions based on false beliefs.

The static base of *DASL* logic does not cover epistemic depth at all, and instead attends to the formalization of agents who can make decisions based on false beliefs,

but whose beliefs are all justified. It is a very low tier of rationality.

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
$\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$	Confidence in 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 3.2: Static Base of Dynamic Agent Safety Logic

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 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.

Defense of 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 inferrable from what is known by an agent. Similarly for positive and negative introspection. We do not take this path in interpreting the knowledge operator. Instead, depart from the majority of formal epistemologists and interpret the knowledge operator as what is actively known by an agent at a moment.

We must justify this interpretation. While ours is a model that departs from those that appear in formal epistemology, we shall argue that it does a good job of capturing the prominent theories of knowledge espoused by non-formal epistemologists. Our

model is also situated well on the rationality tier for modeling realistic human agents whose beliefs are justified, which is a dramatic improvement, from the perspective of realism, over the higher tiers with S5 common knowledge or S5 bounded rationality.

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 neglected by formal efforts, because attention has been on properties of formal systems, see below, rather than primarily on their philosophical foundations.

Third, the model should be a sound and complete logical system to serve as a static foundation for an extension that treats actions. Just because attention must be paid to the philosophical foundations and motivations for a formal system does not mean that the formal system is off the hook for the standard desiderata of a logic.

The logic must avert a collapse of belief and knowledge. The logic must abide by the principle of indistinguishability, where an agent cannot distinguish between possible worlds it cannot rule out. If the logic allows for any introspection, it must avoid the Löbian Obstacle.

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 rational 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 S5 knowledge operator and limiting the agents' epistemic depth, 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 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. Thus, the achievement of realism in idealized games translates to little normative value in the real world.

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. 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, one level above 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. We call this agent a Grounded-Coherent Agent (GCA).

For completeness, we describe the lowest tier of model. It is a model of normal human reasoners studied by behavioral economists and psychologists. They believe things based on systemic biases and make no effort to make a belief system consistent or grounded in evidence. Modeling such an agent from a logical standpoint is challenging for obvious reasons, and certainly these models do not pretend to be normative beyond lessons about how reasoning and beliefs can go awry, for example through cognitive biases and fallacious inferences.

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 GCA model defended here in that its agents can make mistakes due to false beliefs.

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* [44], 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 can avoid Gettier-style counterexamples [11] 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 as formal theories of epistemology proper. The S5 with common knowledge epistemic logic does not strive for this, nor is it clear that theories of bounded rationality do, either. Recent work in epistemology has applied Bayesian reasoning to the justification problem, and probabilistic epistemic logics help formalize these Bayesian approaches. 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 the knowledge explicitly requires belief via the axiom *Knowledge implies belief*. Justification remains for us to elaborate. 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 is that, through reflection, one can know that she knows a proposition, and know that she believes but does not know a proposition [?]. 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 much make decisions regarding the above problems in a way that we can philosophically defend.

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.

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 *Confidence in Belief*. We must show that our system avoids $\mathbf{K} \mathbf{B}$ collapse.

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

Löbian Obstacle The Löbian Obstacle is a problem facing any sufficiently powerful (formally-defined) reasoner. It refers to Löb’s Theorem. The modal version of Löb’s Theorem is

$$\Box(\Box\varphi \Rightarrow \varphi) \Longrightarrow \Box\varphi, \tag{3.2}$$

where the standard interpretation is that \Box denotes provability in Peano Arithmetic. When 3.2 is included as a theorem of a normal modal logic, it can lead to undesirable results. For epistemic logics, or logics with reflective reasoners, the result is a reasoner who is guaranteed to conclude falsehoods, despite proving or knowing or believing (however one interprets the modal operator) that they do not conclude falsehoods.

The theorem grew out of Gödel’s incompleteness results about self-referential true but unprovable statements in systems at least as strong as Peano Arithmetic. Löb addressed the question of self-referential statements asserting their own provability. When quantified over statements, the question concerns a system’s ability to prove its own consistency. Löb’s Theorem states that a system can prove its own consistency only if it can prove every well formed formula in its language, including contradictions. For agents who can reflect on their own reasoning, this is a disaster. This is the Löbian Obstacle that systems must strive to avoid. Research at the foundations of artificial intelligence is actively engaged in this avoidance task, for if an artificial general intelligence, possessed of a reasoning system at least as powerful as Peano Arithmetic (which it likely shall possess), tries to establish that it can trust its own reasoning, it shall cease to reason consistently.

The standard modal provability logic is GL, for Gödel Löb. It is a normal modal logic with 3.2 as an axiom. The conditions that lead to the derivation of Löb’s Theorem for other modal logics are the following:

1. The modal logic contains the Rule of Necessitation: From $\vdash \varphi$ infer $\vdash \Box\varphi$.
2. The modal logic is normal: $\Box(\varphi \Rightarrow \psi) \Rightarrow (\Box\varphi \Rightarrow \Box\psi)$.
3. The modal operator's relation is transitive, *i.e.*, it is a theorem that $\Box\varphi \Rightarrow \Box\Box\varphi$.
4. The modal logic contains modal fixpoints: $\vdash \psi \Leftrightarrow \mathbf{F}(\Box\psi)$, for some formula \mathbf{F} containing ψ as a subformula.

If these four conditions are met, the modal logic will include Löb's Theorem. For logic's of knowledge and belief, conditions (1)-(3) are almost unavoidable for at least one of the operators. Condition (4) is likewise challenging to avoid, as common knowledge and common belief explicitly introduce modal fixpoints, but more fundamentally, there's no principled reason to say the human-like agents being modeled are incapable of expressing a sentence $\psi \stackrel{def}{=} \text{"If I know that } \psi, \text{ then } \varphi\text{"}$. Humans clearly *can* articulate self-referential sentences, and at the very least, they can reason on the level of Peano Arithmetic, which is sufficient for the problem to occur. Such agents will encounter the Löbian Obstacle, and come to know contradictions. This level of error is not desired in the model. Knowledge and belief serve as a foundation for action, so the potential to know and believe anything leads to the potential for any action to be executed in a way that seems rational to the agent.

OLD BELOW

This system warrants defending, as John Rushby has recently challenged the decisions. His concerns can be summarized as follows:

1. Belief is sufficient for reasoning about pilots' actions
2. The axioms of Knowledge and Belief are nonstandard and seem contrived

We will now address each concern and defend the axioms presented above. We note that the arguments below are our best attempts at charitably representing Rushby's arguments.

Solely Belief?

Rushby makes the following argument against including operators for knowledge and belief in a multi-modal system for reasoning about pilot error, and relying instead solely on an operator for belief.

1. A model treating only a pilot's beliefs is simpler than a model treating a pilot's beliefs about her knowledge.
2. If model A is simpler than model B, and A accurately models the phenomenon under investigation, then model A is preferable to model B.
3. Rushby's model treats only pilot belief and accurately models the phenomenon under investigation.
4. Therefore, the Rushby's model is preferable to the DASL model.

We argue that premise 3 is false. Specifically, we argue that Rushby's model does not accurately capture the relationship between the pilot's beliefs and unsafe action.

In Rushby's alternative model, he discusses a pilot's discernment of the presence of instrument conflict, which he assumes for simplicity to be always accurate. Is this not an assumption of knowledge for simplicity? Similarly, when the pilot is aware of the autopilot status through observing an indicator on a display or through inference from the aircraft's behavior, it is not clear that this is not a case of knowing the autopilot's status. Rushby's description of the true belief seems to make it fit with any epistemologist's theory of knowledge.

Knowledge and Belief Axioms

Knowledge is weaker here than in most epistemic logics [42, 47, 44]. Knowledge is almost universally formalized as an S5 modal operator, including axiom schema for positive and negative introspection. Rushby suspects that because the analyses

of aviation mishaps lays the blame on a failure of negative introspection in [5], the decision to exclude the negative introspection axiom from the schema appears contrived. Rushby proves that to include the negative introspection axiom would cause a collapse of knowledge and belief into equivalence, due to the axiom *Confidence in Belief*: $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$, which to our knowledge was first proved by Lenzen [28]. *Confidence in Belief*, additionally, is the converse of Hintikka’s choice for relating the two: $\mathbf{B}_i \mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi$. Axiom *Confidence in Belief* will be defended below. Rushby’s concern with the knowledge axioms, or lack thereof, can be summarized as follows.

1. Negative introspection is a standard axiom for epistemic logic.
2. *Confidence in Belief* is non-standard for multi-modal epistemic-doxastic logics.
3. Including negative introspection and excluding *Confidence in Belief* preserves an analysis based on false beliefs that is analytically sufficient.
4. Standard axioms are preferable to non-standard axioms when standard axioms are analytically sufficient.
5. Therefore, including negative introspection and excluding *Confidence in Belief* is preferable to excluding and including respectively.

Our representation of the argument here is based directly on Rushby’s claims in CITE for premises 1 through 3, but premise 4 is inferred from context and was not explicitly mentioned. However, in order for the argument to be deductively valid, the premise must be included. We take it to be an unstated assumption behind Rushby’s reasoning.

We shall argue that premise 4 is false. Rushby’s preference for the standard axioms, based on their preservation of an analysis based on false beliefs, fails to account for a feature of the present axioms which we shall endeavor to establish, namely that they more accurately describe real human knowledge and belief relations,

whereas the standard axioms do not. Thus, while it is true that the analytical strength of the standard axioms would be preserved based on a false belief analysis, it does not follow that these axioms are preferable, because they do not accurately model human agency.

First, we defend the exclusion of negative introspection for the knowledge operator. Negative introspection, recall, holds that $\neg K_i \varphi \Rightarrow K_i \neg K_i \varphi$: If an agent, i does not know φ , then i knows that she does not know it. When it is phrased this way, it can seem at worst a bit confusing, but not obviously false. However, consider one's response when it is phrased the following equivalent way: An agent does not know φ *only if* she knows that she does not know φ . This deformed version of the axiom makes it clear that knowledge of ignorance is being asserted as a *necessary condition* for ignorance, or lack of knowledge. One would be hard pressed to find an epistemologist willing to defend this view. To further elaborate on the point, consider the contrapositive: $\neg K_i \neg K_i \varphi \Rightarrow K_i \varphi$: If i does not know that she does not know φ , then she knows φ . The Pyrrhonian skeptics famously denied claims to all knowledge, even to the claim that they know nothing. Thus, for all formulas φ , $\neg K_{\text{Pyrrhonian}} \neg K_{\text{Pyrrhonian}} \varphi$. According to the negative introspection axiom, then, $K_{\text{Pyrrhonian}} \varphi$ for arbitrary φ , certainly a result they would be disturbed by.

Many find it intuitive to speak of the epistemic relation as an equivalence or indistinguishability relation, which requires it be symmetric. Symmetry, 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 3.1.0.3 (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 3.1.1 (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 . □

We also have the following.

Theorem 3.1.2 (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 3.1.3 (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 $wR_b^i w'$ for some arbitrary $w, w', v \in W$. It suffices to show that $w', v \in dc(w)$. Since $wR_b^i w'$, $w' \in dc(w)$. Since $w(R_k^i \circ$

$R_b^i)v$, the composed relation being a subset of R_b^i , it follows that $wR_b^i v$ and therefore $v \in dc(w)$. \square

Theorem 3.1.3 does not guarantee that the actual world w is in this bimodal composed equivalence relation, because it may not be the case that $wR_b^i w$. 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 3.1.4 (Strong BK Indistinguishability). *For all $w, v \in W$, if $wR_b^i w$ then for all $v \in W$, $wR_b^i v$ 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 $wR_b^i w$ and $wR_b^i v$. 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 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 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.

BELOW NEEDS WORK. TODO. Belief is standard [42]. They are related logically by EP(1-3), which hold that knowledge entails belief, belief entails that one believes that one knows, and belief entails that one knows that one believes. Finally, actions and safe actions are logically related by SP and PR, which hold that necessary consequences of *mere* action are also necessary consequences of *safe* actions, and that a pilot can execute an action only if he believes that he is executing a safe action.

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

$[A, a]p \Leftrightarrow (pre(a) \Rightarrow (post(A, p) \Rightarrow p))$	Aprop
$[A, a]\neg\varphi \Leftrightarrow (pre(a) \Rightarrow \neg[A, a]\varphi)$	AN
$[A, a](\varphi \wedge \psi) \Leftrightarrow ([A, a]\varphi \wedge [A, a]\psi)$	AC
$[A, a]K_i \varphi \Leftrightarrow (pre(a) \Rightarrow \bigwedge_{a\chi_k^i b} K_i [A, b]\varphi)$	AK
$[A, a]B_i \varphi \Leftrightarrow (pre(a) \Rightarrow \bigwedge_{a\chi_b^i b} B_i [A, b]\varphi)$	AB
<hr/>	
$[A, a]^S p \Leftrightarrow (pre_s(a) \Rightarrow (post(A, p) \Rightarrow p))$	Sprop
$[A, a]^S \neg\varphi \Leftrightarrow (pre_s(a) \Rightarrow \neg[A, a]^S \varphi)$	SN
$[A, a]^S (\varphi \wedge \psi) \Leftrightarrow ([A, a]^S \varphi \wedge [A, a]^S \psi)$	SC
$[A, a]^S K_i \varphi \Leftrightarrow (pre_s(a) \Rightarrow \bigwedge_{aR_k^i b} K_i [A, b]^S \varphi)$	SK
$[A, a]^S B_i \varphi \Leftrightarrow (pre_s(a) \Rightarrow \bigwedge_{aR_b^i b} B_i [A, b]^S \varphi)$	SB

Table 3.3: The reduction axioms of \mathcal{DASL} .

3.2 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 [61].

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.

3.2.1 Soundness

A logic has the property of soundness *if and only if* all of its derivable theorems are valid in its model theory. **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, subsets of R_k^i , and subsets 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 **DASL** is sound is to show that all of its theorems are satisfied by all such models.

Formally, this amounts to:

Theorem 3.2.1 (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 $(R_k^i \circ R_b^i) \subseteq R_b^i$, $(R_b^i \circ R_k^i) \subseteq R_k^i$, 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 [10] for the model theory mechanization and the works of Lescanne and Puisségur [51, 57] and Maliković and Čubrilo [53, 54] 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 $| \vdash$.

```

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

```

The proof can be found in the appendix.

3.2.2 Completeness

For completeness, our proof depends on the notion of a Sahlqvist formula, which itself depends on more foundational properties of modal formulas, which we must also define.

A *Sahlqvist formula* 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 $\mathbf{n} \geq \mathbf{0} \square$'s and **(b)** negative formulas, **(c)** using only \vee , \wedge , and \diamond to build the antecedent from **(a)** and **(b)** components.

To define these notions and reason about axiom schemas rather than actual well-formed formulas of the 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 formula : Type :=
  | FProp : prop → formula
  | FAnd : formula → formula → formula
  | FOr : formula → formula → formula
  | FImp : formula → formula → formula
  | FNeg : formula → formula
  | FK : DASL.Agents → formula → formula
  | FB : DASL.Agents → formula → formula.

```

We define a **negative** formula as follows.

```

Fixpoint negative_formula (phi : formula) : Prop :=
  match phi with
  | FProp p ⇒ False
  | FAnd phi1 phi2 ⇒ (negative_formula phi1)
    ∧ (negative_formula phi2)
  | FOr phi1 phi2 ⇒ (negative_formula phi1)
    ∧ (negative_formula phi2)
  | FImp phi1 phi2 ⇒ (negative_formula phi1)
    ∧ (negative_formula phi2)
  | FNeg phi' ⇒ not (negative_formula phi')
  | FK a phi' ⇒ negative_formula phi'
  | FB a phi' ⇒ negative_formula phi'
  end.

```

We define a **positive** formula as the opposite of a **negative** formula.


```

Fixpoint positive_formula (phi : formula) : Prop :=
  match phi with
  | FProp p  $\Rightarrow$  True
  | FAnd phi1 phi2  $\Rightarrow$  (positive_formula phi1)
     $\wedge$  (positive_formula phi2)
  | FOr phi1 phi2  $\Rightarrow$  (positive_formula phi1)
     $\wedge$  (positive_formula phi2)
  | FImp phi1 phi2  $\Rightarrow$  (positive_formula phi1)
     $\wedge$  (positive_formula phi2)
  | FNeg phi'  $\Rightarrow$  not (positive_formula phi')
  | FK a phi'  $\Rightarrow$  positive_formula phi'
  | FB 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 : formula) : Prop :=
  match phi with
  | FProp p  $\Rightarrow$  True
  | FAnd phi1 phi2  $\Rightarrow$  False
  | FOr phi1 phi2  $\Rightarrow$  False
  | FImp phi1 phi2  $\Rightarrow$  False
  | FNeg phi'  $\Rightarrow$  False
  | FK a phi'  $\Rightarrow$  boxed_formula phi'
  | FB 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 : formula) : Prop :=
  match phi with
  | FProp p  $\Rightarrow$  True
  | FAnd phi1 phi2  $\Rightarrow$  (s_a_component phi1)  $\wedge$  (s_a_component phi2)
  | FOr phi1 phi2  $\Rightarrow$  (s_a_component phi1)  $\wedge$  (s_a_component phi2)
  | FImp phi1 phi2  $\Rightarrow$  not (s_a_component phi1)  $\wedge$  (s_a_component phi2)
  | FNeg phi'  $\Rightarrow$  match phi' with
    | FProp p  $\Rightarrow$  True
    | FAnd p1 p2  $\Rightarrow$  positive_formula p1  $\wedge$  positive_formula p2
    | FOr p1 p2  $\Rightarrow$  positive_formula p1  $\wedge$  positive_formula p2
    | FImp p1 p2  $\Rightarrow$  negative_formula p1  $\wedge$  positive_formula p2
    | FNeg p'  $\Rightarrow$  s_a_component p'
    | FK a p'  $\Rightarrow$  not (s_a_component p')
    | FB a p'  $\Rightarrow$  not (s_a_component p')
  end
  | FK a phi'  $\Rightarrow$  boxed_formula phi'
  | FB a phi'  $\Rightarrow$  boxed_formula phi'
end.

```

Theorem 3.2.2 (Sahlqvist's Theorem). *Every Sahlqvist formula φ defines a complete modal logic when added as an axiom to the minimal modal logic defined by the axiom K , with respect to the first order frame relation corresponding to φ , $\mathbf{FO}(\varphi)$.*

For soundness and completeness, it suffices to show that each DASL axiom is a Sahlqvist formula. In order to make use of Theorem 3.2.2, we characterize the Negative Belief Introspection axiom as the equivalent contrapositive $\langle \mathbf{B}_i \rangle \mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \varphi$, as in below:

$\mathbf{K}_i (\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K}_i \varphi \Rightarrow \mathbf{K}_i \psi)$	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
$\langle \mathbf{B}_i \rangle \mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \varphi$	(Sahlqvist) Negative Belief Introspection
$\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi$	Knowledge implies Belief
$\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$	Confidence in 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 3.4: Equivalent Static Base with Sahlqvist-friendly Negative Belief Introspection.

OLD BELOW

Theorem 3.2.3 (Soundness). *Dynamic Agent Safety Logic is sound for Kripke structures with*

- (1) reflexive \mathbf{R}_k^i relations,
- (2) serial, transitive, Euclidean \mathbf{R}_b^i relations,
- (3) which are partially ordered $(\mathbf{R}_k^i \circ \mathbf{R}_b^i) \subseteq \mathbf{R}_b^i$, $(\mathbf{R}_b^i \circ \mathbf{R}_k^i) \subseteq \mathbf{R}_b^i$, and $\mathbf{R}_b^i \subseteq \mathbf{R}_k^i$.

Proof. (1) *and* (2) correspond to the axioms that \mathbf{K}_i is a T modality and \mathbf{B}_i is a KD45 modality in the usual way. (3) corresponds to EP1, EP2, and EP3. Axioms AP through SB are reduction axioms. This leaves (4), corresponding to SP, and (5) which corresponds to PR. Here we will prove (5). Let \mathbf{M} be a Kripke structure satisfying the five conditions above. Let \mathbf{A} be an Action structure with \mathbf{a} and \mathbf{i} as its actual action token and agent.

We prove (5) via the contrapositive of PR: $\langle \mathbf{B}_i \rangle [\mathbf{A}, \mathbf{a}]^{\mathbf{S}} \varphi \Rightarrow [\mathbf{A}, \mathbf{a}] \varphi$. Assume

$M, w \models \langle B_i \rangle [A, a]^S \varphi$. By the semantics of $\langle B_i \rangle$, there exists a v , such that $wR_b^i v$ and $v \models [A, a]^S \varphi$. From the semantics, it follows that for all M', v' , if $(M, v) \llbracket (A, a) \rrbracket^S (M', v')$ then $M', v' \models \varphi$. By slightly abusing the notation, and letting $(W, w)R_b^i(W, v)$ be equivalent to $wR_b^i v$, we can create the composed relation $(\llbracket (A, a) \rrbracket^S \circ R_b^i)$. It then holds, by condition (5), that $(M, w)(\llbracket (A, a) \rrbracket^S \circ R_b^i)(M', v')$ implies $(M, w) \llbracket (A, a) \rrbracket (M', v')$. So, for all M', v' , if $(M, w) \llbracket (A, a) \rrbracket (M', v')$, then $M', v' \models \varphi$. So, $M, w \models [A, a] \varphi$. \square

Aprop : \Rightarrow . Assume $M, w \models [A, a]p$. We must show that $M, w \models pre(a) \Rightarrow (post(A, p))$. By the semantics of $[A, a]$, for all (M', w') , if $M, w \models pre(a)$ and $update(M, A, w, a, i) = (M', w')$, then $M', w' \models p$. By definition of $post(A, p)$, if $update(M, A, w, a, i) = (M', w')$ and $M', w' \models p$, then $post(A, p) = p$. So, if $M, w \models pre(a)$, then $post(A, p) = p$, and thus $post(A, p) \Rightarrow p$.

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

AN : \Rightarrow . Assume $M, w \models [i, (A, a)] \neg \varphi$. It suffices to show that $M, w \models pre(a) \Rightarrow \langle i \rangle (A, a) \neg \varphi$. From the assumption and the semantics, for all (W', w') , if $(M, w) \llbracket (A, a) \rrbracket (M', w')$ then $M', w' \models \neg \varphi$. So, if $M, w \models pre(a)$ and $update(M, A, w, a, i) = (M', w')$, then $M', w' \models \neg \varphi$. Assume $M, w \models pre(a)$, and it follows that $update(M, A, w, a, i)$ is defined, so there exists a M', w' such that $(M, w) \llbracket (A, a) \rrbracket (M', w')$ and $update(M, A, w, a, i) = (M', w')$ and $M', w' \models \neg \varphi$. Therefore, $M, w \models pre(a) \Rightarrow \langle A, a \rangle \neg \varphi$.

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

then $M', w' \models \neg\varphi$, and therefore $M, w \models [A, a]\neg\varphi$.

AC is obvious.

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

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

\Leftarrow . Assume $M, w \models \text{pre}(a) \Rightarrow K_i [i, (A, a)]\varphi$. We must show $M, w \models [i, (A, a)]$

$K_i \varphi$. Thus, we must show $M, w \models \text{pre}(a)$ and $\text{update}(M, A, w, a, i) = (M', w')$ implies $M', w' \models K_i \varphi$. So it suffices to show that if $\text{update}(M, A, w, a, i) = (M', w')$ and $M, w \models K_i [i, (A, a)]\varphi$, then $M', w' \models K_i \varphi$. Assume $\text{update}(M, A, w, a, i) = (M', w')$ and $M, w \models K_i [i, (A, a)]\varphi$. Then for all v , if $wR_k^i v$, then $M, v \models [i, (A, a)]\varphi$. It follows that $M, v \models \text{pre}(a)$ and $\text{update}(M, A, v, a, i) = (M', v')$ 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, $M', w' \models K_i \varphi$.

Proofs for **AB** through **SB** follow the above proofs exactly analogously. \square

Assume $M, w \models \langle A, a \rangle \text{true}$. By the semantics of $\langle A, a \rangle$, $M, w \models \text{pre}(a)$ and $\text{update}(M, w, \chi) \models \text{true}$. Let $(M', w') = \text{update}(M, w, \chi)$. Then $M', w' \models \text{true}$. From (5) above, it holds that $R_b^i(w) \subseteq V^s(a)$. R_b is serial, so there is at least one such $v \in R_b^i$. Then $M, v \models \text{pre}_s(a)$. From (4), $M, v \models \text{pre}(a)$, so $\text{update}(M, v, \chi)$ is defined, call it (M'', v') . Because (M'', v') is defined, $M'', v' \models \text{true}$. So, $M, v \models \text{pre}_s(a)$ and $\text{update}(M, v, \chi) \models \text{true}$. This holds for all v , such that $wR_b v$. Thus, $M, w \models B_i \langle A, a \rangle^s \text{true}$. Therefore, $M, w \models \langle A, a \rangle \text{true} \Rightarrow B_i \langle A, a \rangle^s \text{true}$. \square

Next we turn to completeness.

3.2.3 Completeness

Completeness proofs in contemporary modal logic research follow the following format. First, a *canonical model* is defined such that it belongs to the Kripke frames for which the logic under investigation is sound. The model combines validity and deducibility so that the former entails the latter. The worlds of the canonical model are sets of formulas from the language. If a formula is a member of a such a set, then it is true at the world. We define this formally below. The objective is to show that all formulas valid for the relevant Kripke frames are likewise deducible in the logic. The key is to construct the sets in a particular way, to be discussed below.

The previous section proved soundness of DASL with respect to frames with a reflexive R_k^i relation, a serial, transitive, Euclidean R_b^i relation, which are partially ordered $(R_k^i \circ R_b^i) \subseteq R_b^i$, $(R_b^i \circ R_k^i) \subseteq R_b^i$, and $R_b^i \subseteq R_k^i$, and with action relations $\llbracket(A, a)\rrbracket \subseteq \llbracket(A, a)\rrbracket^S$, and $(\llbracket(A, a)\rrbracket^S \circ R_b^i) \subseteq \llbracket(A, a)\rrbracket$. The canonical model we proceed to define will belong to this frame. Because it belongs to this frame, the logical closure of a set of formulas will consist of valid DASL deductions. This allows us to infer deducibility from validity. First we must define the notion of a maximal consistent set.

Informally, a maximal consistent set of formulas is one that is a subset of the well-formed formulas of the language \mathcal{L}_{DASL} according to the syntax presented earlier, whose members are consistent with each other, such that the set has no consistent extension.

Definition. Maximal Consistent Set.

For a set of formulas $\Gamma \subseteq \mathcal{L}_{DASL}$, Γ is maximal consistent iff

1. Γ is consistent: $\Gamma \not\vdash \perp$.
2. Γ is maximal: there is no $\Gamma' \subseteq \mathcal{L}_{DASL}$ such that $\Gamma \subset \Gamma'$ and $\Gamma' \not\vdash \perp$.

A canonical model is a Kripke model of \mathcal{L}_{DASL} with a set of worlds W^C that collectively satisfy all formulas from a maximal consistent set. For our formalism here, we treat worlds as identical to sets of formulas.

Definition. Canonical Model. A canonical model $M^C = \langle W^C, R_{k,i}^C, R_{b,i}^C, w, V^C \rangle$ is defined:

1. $W^C = \{\Gamma \mid \Gamma \text{ is maximal consistent relative to } \mathcal{L}_{DASL}\}$
2. $\Gamma R_{k,i}^C \Delta$ iff $K_i \varphi \in \Gamma$ implies $\varphi \in \Delta$
3. $\Gamma R_{b,i}^C \Delta$ iff $B_i \varphi \in \Gamma$ implies $\varphi \in \Delta$
4. $V^C(p) = \{\Gamma \in W^C \mid p \in \Gamma\}$

Theorem 3.2.4 (Completeness). *The language of Dynamic Agent Safety Logic, \mathcal{L}_{DASL} , is complete for Kripke structures with*

- (1) reflexive R_k^i relations,
- (2) serial, transitive, Euclidean R_b^i relations,
- (3) which are partially ordered $(R_k^i \circ R_b^i) \subseteq R_b^i$, $(R_b^i \circ R_k^i) \subseteq R_b^i$, and $R_b^i \subseteq R_k^i$,
- (4) $\llbracket (A, a) \rrbracket \subseteq \llbracket (A, a) \rrbracket^S$ and
- (5) $(\llbracket (A, a) \rrbracket^S \circ R_b^i) \subseteq \llbracket (A, a) \rrbracket$.

Proof. Completeness states that if a formula φ is valid, then it is deducible. The sketch for this proceeds by contraposition. So we must show that for every formula φ in the language \mathcal{L}_{DASL} , $\not\models \varphi$ implies $\not\vdash \varphi$.

The proof appeals to the following lemmas, proven also in [47, 30]:

Lemma 3.2.5 (Lindenbaum). *Every consistent set of formulas is a subset of a maximal consistent set of formulas.*

Proof. If we begin with a consistent set of \mathcal{L}_{DASL} formulas β , we construct a maximal

consistent set Γ as follows:

1. enumerate the formulas of \mathcal{L}_{DASL} $\varphi_0, \dots, \varphi_n, \dots$
2. $\Gamma_0 \equiv \beta$
3. $\Gamma_{n+1} \equiv \{\varphi_n\}$ if $\Gamma_n \vdash \varphi_n$
4. $\Gamma_{n+1} \equiv \{\neg\varphi_n\}$ otherwise
5. $\Gamma \equiv \bigcup_{n \geq 0} \Gamma_n$

We now show that Γ is consistent. Γ is consistent if and only if $\forall \Gamma_n \subset \Gamma$, $\Gamma_n \not\vdash \perp$. The base set Γ_0 is consistent because it is equivalent to β , which is consistent by assumption. The induction hypothesis is if Γ_{n+1} is consistent then so is Γ_n . There are two ways to construct Γ_{n+1} : via step (3), or via step (4). If via (3), then the added formula φ_n logically follows from Γ_n , and \mathcal{L}_{DASL} is sound, so Γ_{n+1} is consistent. Otherwise, if via (4), then φ_n does not logically follow from Γ_n . Because the formulas are enumerated and examined in order, φ_n is not in Γ_n . Thus, $\Gamma_n \cup \{\neg\varphi_n\} \not\vdash \perp$. Therefore, Γ is consistent.

We show Γ is maximal by contradiction, assuming there is a $\Gamma' \subset \mathcal{L}_{DASL}$ such that $\Gamma \subset \Gamma'$ and $\Gamma' \not\vdash \perp$. Let $\varphi_n \in \Gamma'$ and $\varphi_n \notin \Gamma$ be an arbitrary extension of Γ . We show $\Gamma' \cup \{\varphi_n\} \vdash \perp$. When constructing Γ , the procedure examined φ_n because it is a formula of \mathcal{L}_{DASL} . Since $\varphi_n \notin \Gamma$, it must be the case that $\Gamma_n \not\vdash \varphi_n$. In step (4) we set $\Gamma_{n+1} \equiv \Gamma_n \cup \neg\varphi_n$. Since $\neg\varphi_n \in \Gamma$, we would have $\Gamma \cup \varphi_n \vdash \perp$, and so $\Gamma' \vdash \perp$, contradicting the assumption. Therefore, Γ is maximal. \square

Lemma 3.2.6 (Properties).

If Γ and Δ are maximal consistent sets and β a consistent set, then:

1. Γ and Δ are deductively closed.
2. $\varphi \in \Gamma$ iff $\neg\varphi \notin \Gamma$.

3. If $\beta \not\vdash \varphi$, then there exists a maximal consistent set Γ such that $\beta \subset \Gamma$ but $\varphi \notin \Gamma$

Proof. (1) clearly follows from the construction of the maximal consistent sets as seen in the proof of Lindenbaum, specifically step (3) of the process. (2) holds because each formula of \mathcal{L}_{DASL} was enumerated and considered, and either it or its negation is added to the growing set at steps (3) and (4). (3) follows from step (4) of the construction process of maximal consistent sets in the Lindenbaum proof. \square

Lemma 3.2.7 (Truth). *For every $\varphi \in \mathcal{L}_{DASL}$, and every maximal consistent set Γ :*

$$\varphi \in \Gamma \text{ iff } (M^C, \Gamma) \models \varphi \quad (3.3)$$

Proof. By induction on φ , all the non-modal formulas follow straight-forwardly from the definition of $V^C(p) = \{\Gamma \in W^C \mid p \in \Gamma\}$. The modal formulas require some explanation. For the left to right direction, assume we have a formula $K_i \varphi \in \Gamma$. We must show that $(M^C, \Gamma) \models K_i \varphi$. The definition of $\Gamma R_{k,i}^C \Delta$ is that $K_i \varphi \in \Gamma$ implies $\varphi \in \Delta$. Thus, for all $\Delta \in W^C$, if $\Gamma R_{k,i}^C \Delta$, then $\varphi \in \Delta$. Therefore, by the semantics of K_i , $(M^C, \Gamma) \models K_i \varphi$.

We prove the right to left direction by its contrapositive, and assume $K_i \varphi \notin \Gamma$. We show $(M^C, \Gamma) \not\models K_i \varphi$. Since $K_i \varphi \notin \Gamma$, from property (2) of Lemma 3.2.6 it follows that $\neg K_i \varphi \in \Gamma$. From the semantics of K_i , there is some $\Delta \in W^C$ such that $\Gamma R_{k,i}^C \Delta$ and $(M^C, \Delta) \models \neg \varphi$. Therefore, $(M^C, \Gamma) \models \neg K_i \varphi$, and equivalently, $(M^C, \Gamma) \not\models K_i \varphi$.

A similar proof shows the $B_i \varphi \in \Gamma$ case, and the dynamic modalities follow from the reduction axioms. \square

Lemma 3.2.8 (Canonicity). *The canonical model satisfies the frame conditions of the \mathcal{L}_{DASL} semantics.*

Proof. We must show that

1. $R_{k,i}^C$ relation is reflexive,
2. $R_{b,i}^C$ relation is serial, transitive, and Euclidean,
3. $R_{b,i}^C \subseteq R_{k,i}^C$,
4. $(R_k^i \circ R_b^i) \subseteq R_b^i$,
5. $(R_b^i \circ R_k^i) \subseteq R_b^i$,
6. $\llbracket (A, a) \rrbracket^S \subseteq \llbracket (A, a) \rrbracket$, and
7. $(\llbracket (A, a) \rrbracket^S \circ R_b^i) \subseteq \llbracket (A, a) \rrbracket$.

We prove these by appealing to the fact that the worlds in the canonical model \mathbf{M}^C are maximal consistent sets that each contains all of the theorems of \mathcal{L}_{DASL} . Therefore, each of the axiom schemas of the DASL Hilbert system is true at each world. We use this to prove that the doxastic, belief, and dynamic relations on \mathbf{M}^C satisfy the above properties, thereby placing \mathbf{M}^C in the appropriate class of Kripke frames.

For (1), we appeal to the fact that all instances of the axiom schema $\mathbf{K}_i \varphi \Rightarrow \varphi$ are elements of each world in \mathbf{W}^C , the set of maximal consistent sets in \mathbf{M}^C . Assume an arbitrary world $(\mathbf{M}^C, \Gamma) \models \mathbf{K}_i \varphi$. By the axiom schema just noted, it follows that $(\mathbf{M}^C, \Gamma) \models \varphi$. Thus, $(\mathbf{M}^C, \Gamma) \models \mathbf{K}_i \varphi$ implies $(\mathbf{M}^C, \Gamma) \models \varphi$. From the definition of the canonical model above, it follows that $\Gamma R_{k,i}^C \Gamma$. This holds for all Γ , and so $R_{k,i}^C$ is reflexive.

Condition (2) depends on the axiom schemas of the doxastic operator.

The schema $\mathbf{B}_i \varphi \Rightarrow \langle \mathbf{B}_i \rangle \varphi$ establishes that $R_{b,i}^C$ is serial. We assume an arbitrary world $(\mathbf{M}^C, \Gamma) \models \langle \mathbf{B}_i \rangle \varphi$. By the schema, it follows that $(\mathbf{M}^C, \Gamma) \models \mathbf{B}_i \varphi$.

By the semantics of $\langle \mathbf{B}_i \rangle$, there exists a world $(M^C, \Delta) \models \varphi$ such that $\Gamma R_{b,i}^C \Delta$. Thus, for all $\Gamma \in W^C$, $\exists \Delta \in W^C$ such that $\Gamma R_{b,i}^C \Delta$, which is what it means to be serial.

The schema $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{B}_i \varphi$ establishes transitivity of $R_{b,i}^C$.

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

Condition (3) follows from the axiom schema $\mathbf{K}_i \varphi \Rightarrow \mathbf{B}_i \varphi$.

Recall the definition of two relations S and R *composing*:

$$S \circ R \equiv \{(x, z) \in X \times Z \mid \exists y \in Y : (x, y) \in R \wedge (y, z) \in S\} \quad (3.4)$$

Condition (4) follows from the axiom schema $\mathbf{B}_i \varphi \Rightarrow \mathbf{B}_i \mathbf{K}_i \varphi$.

Condition (5) follows from the axiom schema $\mathbf{B}_i \varphi \Rightarrow \mathbf{K}_i \mathbf{B}_i \varphi$.

Condition (6) follows from the axiom schema $[\mathbf{A}, a] \varphi \Rightarrow [\mathbf{A}, a]^S \varphi$.

Condition (7) follows from the axiom schema

$$\begin{aligned} \llbracket (\mathbf{A}, a) \rrbracket^{S \circ R_{b,i}^C} &\equiv \{(\Gamma, \Delta') \in W^C \times W^C \mid \exists \Delta \in W^C : (\Gamma, \Delta) \in R_{b,i}^C \wedge (\Delta, \Delta') \in \\ &\llbracket (\mathbf{A}, a) \rrbracket^S\} \end{aligned}$$

Must show: $\{(\Gamma, \Delta')\} \subset \llbracket (\mathbf{A}, a) \rrbracket$.

$\langle \mathbf{iA}, a \rangle \varphi \Rightarrow \mathbf{B}_i \langle \mathbf{A}, a \rangle^S \varphi$. Assume an arbitrary world $(M^C, \Gamma) \models \langle \mathbf{A}, a \rangle \varphi$. By the axiom schema, it follows that $(M^C, \Gamma) \models \mathbf{B}_i \langle \mathbf{A}, a \rangle^S \varphi$. By the semantics of \mathbf{B}_i , it follows that $\forall \Delta \in W^C$, if $\Gamma R_{b,i}^C \Delta$, then $(M^C, \Delta) \models \langle \mathbf{A}, a \rangle^S \varphi$. We use the reduction law SN : $[\mathbf{A}, a]^S \neg \varphi \Leftrightarrow (pre_s(a) \Rightarrow \neg [\mathbf{A}, a]^S \varphi)$ to deduce the equivalent if $\Gamma R_{b,i}^C \Delta$, then $(M^C, \Delta) \models pre_s(a)$ and $(M^C, \Delta) \llbracket (\mathbf{A}, a) \rrbracket^S (M^C, \Delta')$ implies $(M^C, \Delta') \models \varphi$, which implies $update(M^C, \mathbf{A}, \Delta, a, i) = (M^C, \Delta')$.

The semantics of $[\mathbf{A}, a]^S$ yield $update(M^C, \mathbf{A}, \Delta, a, i) = (M'^C, \Delta')$. From the assumption that $(M^C, \Gamma) \models \langle \mathbf{A}, a \rangle \varphi$ and the semantics of $\langle \mathbf{A}, a \rangle$, it is true that $(M^C, \Gamma) \llbracket (\mathbf{A}, a) \rrbracket (M'^C, \Gamma')$ and $(M'^C, \Gamma') \models \varphi$. \square

With these lemmas, we assume $\not\models \varphi$, and show that $\not\models \varphi$. Since $\not\models \varphi$, the set

$\{\neg\varphi\}$ is consistent. From Lindenbaum, $\{\neg\varphi\}$ is part of a maximal consistent set, call it Γ . From the Truth Lemma, $(M^C, \Gamma) \models \neg\varphi$. Therefore, $\not\models \varphi$. \square

Because the static logic is complete and we have translation axioms that convert the dynamic formulas to equivalent static ones, we can conclude that the entirety of **DASL** is complete. In the next chapter we examine case studies and a mechanization of the logic.

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 [53, 54], in which they mechanize an analysis of the game of Cluedo using Dynamic Epistemic Logic, based on van Ditmarsch’s formalization of the game [41]. 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[21]. 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[50].

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 [56] 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). $B_i \varphi \Rightarrow \neg K_i \neg \varphi$.

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

$$B_i \varphi \Rightarrow \langle B_i \rangle \varphi,$$

which is equivalent to

$$B_i \varphi \Rightarrow \neg B_i \neg \varphi.$$

Due to axiom EP1, it follows that

$$B_i \varphi \Rightarrow \neg K_i \neg \varphi. \quad \square$$

We now formalize the critical moment.

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

Before beginning, we note that our use of sets in the following Coq code requires the following argument passed to coqtop before executing: `-impredicative-set`. In CoqIDE, this can be done by selecting the ‘Tools’ dropdown, then ‘Coqtop arguments’. Type in *-impredicative-set*.

We first formalize the set of agents.

```
Inductive Agents: Set := 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.

```
Inductive Inputs : Set :=
  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 : Set := 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 reading *VertUp1* may represent a nose up reading between 0° and 10° , while *VertUp2* represents a reading between 11° and 20° .

```

Inductive Readings (s : Side) : Set :=
  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.

```

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

```

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

```

Next we follow Maliković and Čubrilo [53, 54] in defining a set *prop* of propositions in predicate calculus, distinct from Coq’s built in type *Prop*. 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. Interestingly, modal logic cannot be directly represented in Coq’s framework [51]. We first define propositions in first-order logic, which we then use to define DASL. This appears to be the standard technique for mechanizing modal logics in Coq.

```

Inductive prop : Set :=
| atm : Atoms → prop
| imp : prop → prop → prop
| Forall : forall (A : Set), (A → prop) → prop
| K : Agents → prop → prop
| B : Agents → prop → prop

```

We use the following notation for implication and universal quantification.

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

```
Notation "\- / p" := (Forall _ p) (at level 70, right associativity).
```

We likewise follow Maliković and Čubrilo [53, 54] by defining an inductive type *theorem* representing a theorem of DASL. 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 " $| \text{--- } p$ " := (theorem p) (at level 80).

We encode actions as records in Coq, recording 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, $\neg[i, (A, a)]\neg True$, equivalently $\langle i \rangle(A, a)True$, 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))  $\implies$ 
NOT (atm (InstrumentsG (Global ms rls rms rrs)))  $\implies$ 
(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.

1. $(\mathbf{LeftAIHorLeft2}) \wedge \neg(\mathbf{MiddleAIHorLeft2}) \dots$ — configuration.
2. $\langle \mathbf{pilot} \rangle \mathbf{hardwingrighttrue}$ — pilot input.
3. $\mathbf{B}_{pilot}(\mathbf{MiddleAIHorLeft2})$ — from axiom PR, \mathbf{pre}_s .
4. $\neg \mathbf{K}_{pilot}(\mathbf{MiddleAIHorLeft2})$ — from axiom K-Reflexive.
5. $\mathbf{B}_{pilot} \mathbf{K}_{pilot}(\mathbf{MiddleAIHorLeft2})$ — from (3), axiom EP2.
6. $\neg \mathbf{K}_{pilot} \neg \mathbf{K}_{pilot}(\mathbf{MiddleAIHorLeft2})$ — from (5), Lemma 4.1.1.
7. $\neg \mathbf{K}_{pilot}(\mathbf{MiddleAIHorLeft2}) \wedge \neg \mathbf{K}_{pilot} \neg \mathbf{K}_{pilot}(\mathbf{MiddleAIHorLeft2})$ — from (4), (6).
8. $\neg(\neg \mathbf{K}_{pilot}(\mathbf{MiddleAIHorLeft2}) \Rightarrow \mathbf{K}_{pilot} \neg \mathbf{K}_{pilot}(\mathbf{MiddleAIHorLeft2}))$ — from (7).

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 [55]. 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.

1. $(A/T = \text{Off}) \wedge (\text{AirspeedSlow3}) \dots$ — configuration.
2. $\langle \text{pilot} \rangle \text{hardthrustminustrue}$ — pilot input.
3. $B_{\text{pilot}}(A/T = \text{On})$ — from axiom PR, pre_s .
4. $\neg K_{\text{pilot}}(A/T = \text{On})$ — from axiom K-Reflexive.
5. $B_{\text{pilot}} K_{\text{pilot}}(A/T = \text{On})$ — from (3), axiom EP2.
6. $\neg K_{\text{pilot}} \neg K_{\text{pilot}}(A/T = \text{On})$ — from (5), Lemma 4.1.1.
7. $\neg K_{\text{pilot}}(A/T = \text{On}) \wedge \neg K_{\text{pilot}} \neg K_{\text{pilot}}(A/T = \text{On})$ — from (4), (6).
8. $\neg(\neg K_{\text{pilot}}(A/T = \text{On}) \Rightarrow K_{\text{pilot}} \neg K_{\text{pilot}}(A/T = \text{On}))$ — from (7).

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.

$$\neg pre_s(a) \Rightarrow K_i \neg K_i pre_s(a)$$

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

$$\neg pre_s(a) \Rightarrow K_i \neg pre_s(a)$$

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 $K_i \neg \varphi \Rightarrow K_i \neg K_i \varphi$. Assume $K_i \neg \varphi$ holds. From EP1, it follows that $B_i \neg \varphi$, and because knowledge is a normal modality, it follows that $K_i B_i \neg \varphi$ holds. From lemma 1, and again the fact that knowledge is normal modality, it follows that $K_i \neg K_i \neg \neg \varphi$, or equivalently, that $K_i \neg K_i \varphi$. \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: $\{\text{False}, \text{True}\}$. Let the epistemic relation include (u, v) , (v, u) , and the reflexive relations. Then at world u φ is false, and $K_i \neg K_i \varphi$ is true, but $K_i \neg \varphi$ is false.

The formalizations show that from the pilot's unsafe action, it follows that he

lacks negative introspection of the safety precondition.

$$\langle i \rangle a_{true} \wedge \neg pre_s(a) \Rightarrow \neg K_i \neg K_i pre_s(a) \quad (4.1)$$

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.

$$\langle i \rangle a_{true} \wedge \neg pre_s(a) \Rightarrow \neg K_i \neg pre_s(a) \quad (4.2)$$

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.

$$K_i \neg pre_s(a) \Rightarrow \neg \langle i \rangle a_{true} \vee pre_s(a) \quad (4.3)$$

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

$$K_i \neg pre_s(a) \Rightarrow \neg \langle i \rangle a_{true}. \quad (4.4)$$

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: $a : \textit{Action}$, $C : \textit{Configuration}$, $pre_s : \textit{Action} \mapsto \textit{Configuration}$, $k : \textit{Int}$.

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 $(c_1 \wedge c_2 \wedge \dots \wedge c_n) \vee (c'_1 \wedge c'_2 \wedge \dots \wedge c'_m)$. The negation of this has the form $(\neg c_1 \vee \neg c_2 \vee \dots \vee \neg c_n) \wedge (\neg c'_1 \vee \neg c'_2 \vee \dots \vee \neg c'_m) \wedge \dots$, 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

It turns out that the Z3 Theorem Prover can solve this problem, which we turn to in the next chapter.

Chapter 5

Ensuring Delivery of Safety-Critical Information

Theorem 5.0.1 (Awareness of Safety-Critical Information). *If a rational pilot \mathbf{i} has awareness about the safety precondition of action a and the safety precondition is false, then \mathbf{i} does not execute a .*

Proof.

1. $\neg K_{\mathbf{i}} pre_s(a) \Rightarrow K_{\mathbf{i}} \neg K_{\mathbf{i}} pre_s(a)$ Awareness Assumption
2. $\neg pre_s(a)$ Safety Precondition Assumption
3. $\neg pre_s(a) \Rightarrow \langle \mathbf{i} \rangle a True \Rightarrow \neg K_{\mathbf{i}} \neg K_{\mathbf{i}} pre_s(a)$ Theorem 1
4. $\langle \mathbf{i} \rangle a True \Rightarrow \neg K_{\mathbf{i}} \neg K_{\mathbf{i}} pre_s(a)$ 2, 3 Modus Ponens
5. $\neg K_{\mathbf{i}} pre_s(a)$ 2, K reflexive
6. $K_{\mathbf{i}} \neg K_{\mathbf{i}} pre_s(a)$ 1, 5 Modus Ponens
7. $\neg \langle \mathbf{i} \rangle a True$ 4, 6 Modus Tollens
8. $[i, a] False$ Def. $\langle \mathbf{i} \rangle a$ \square

□

Thus, if the countermeasure succeeds in establishing the pilot’s awareness of the safety-critical information’s being false, it follows by the logic of DASL that the agent discontinues the unsafe action. The goal is to convince her that she does not know what she thinks she knows, like the airspeed, or the horizontal attitude.

The insight gained from the formal analysis allows us to increase aviation safety, but we require a means for computing the safety-critical information that the pilot is unaware of. We do this by encoding the actual configuration of the flight deck and the safety precondition of an input action as a set of constraints and a formula to be examined by a SMT solver, and then check whether it is satisfiable. If it is not, we are interested in retrieving the specific instrument readings that conflict with the safety precondition. In SMT solving, this collection of clauses is called the unsatisfiable core. In the next section, we describe the formal properties of the unsatisfiable core and other ways it is used in formal methods research.

Before providing background on the unsatisfiable core’s role in other applications of formal methods, we will provide a bit more discussion on the safety precondition. For this paper, we take it as a given that the safety precondition of an action can be determined by formally representing the contents of flight safety manuals. The specificity of instructions one finds in aviation manuals is one of the principle reasons we have this confidence. Aviation safety frequently discusses the notion of a “flight envelop” as a range of conditions in which safe actions can be executed. In a sense, we are formalizing the flight envelop as constraints for a SMT solver. However, we do not represent the safety preconditions for actions as they would appear in a higher fidelity system, but rather approximate them with a little bit of common sense, in order to achieve a balance between illustration of the approach and brevity. Were this approach to be engineered for real-world applications, the safety preconditions, and indeed the instrument configurations, would no doubt be more complex.

5.1 SMT Solving and the Unsatisfiable Core

Satisfiability Modulo Theories (SMT) is a decision problem from theoretical computer science where a formula from first order logic, constrained by theories in the form of other first order formulas, is assessed for satisfiability. Tools exist to solve engineering problems amenable to formalization in first order logic with constraining theories.

Consider the following example. The first order logic formula is

$$(x < y) \wedge (y = 10) \wedge (x > 20) \wedge (z = 40).$$

The constraining theory is the standard set of axioms of arithmetic. The formula is unsatisfiable.

When a formula constrained by theories is unsatisfiable, the unsatisfiable core refers to the set of subformulas responsible for a clausal formula's being unsatisfiable. That is, the unsatisfiable core of the formula is the portion of it that conflicts with respect either to internal consistency or with the constraining theories. In the example above, the unsatisfiable core is the first three propositions. The unsatisfiable core is successfully applied as a formal method for engineering in the domains of hardware and software verification. When engineers design a chip, before fabricating it they need to verify that the design meets the specifications. One approach to this is to formalize the chip components as variables in a first order formula, and to formalize the specification as a set of constraints on the formula. Then, running the design and specification through a SMT solver, one can identify which aspects of the design are in conflict with the specification. Software is similarly formalized and analyzed during its specification stage.

Some recent work has applied SMT solving to what are known as cyber physical

systems, or hybrid systems [48]. In particular, some hybrid systems are those that involve human machine interaction. Efforts to apply SMT solving to human machine interaction have succeeded in identifying execution paths the system can go down that could cause the human user to become confused. Furthermore, this approach has been applied to the domain of aviation safety. However, as yet no efforts have been undergone to apply SMT solving for the run-time diagnosis and error correction of human error in hybrid systems.

This work does so by leveraging the logical connection between a pilot’s action, her knowledge and beliefs about the system, and the features of the system that make the action safe or unsafe. We model the conditions of the action’s safe execution as a constraining theory of instrument readings, and the actual instrument readings as a first order formula to be tested for satisfiability modulo that theory. If the formula modeling the actual instrument readings is not satisfiable, we extract the unsatisfiable core. Based on the logical connection established by DASL, we identify the unsatisfiable core with the safety-critical information that the pilot is unaware of when she executes an unsafe action.

5.2 Encoding and Case Studies

We develop an encoding of the instruments and the safety preconditions of actions so that they can be fed into the Z3 SMT solver and, if the actual instrument readings conflict with the constraining safety precondition, their unsatisfiable core may be extracted [52]. The case studies we present as examples in this section are meant to illustrate the encoding and the usefulness of the approach. With length limitations, we include three case studies in order to show the versatility of the approach, but sacrifice the complexity of the representations. The actual values of the case studies have

been simplified for presentation, and the encodings themselves have been abridged to include only enough of the instrumentation and safety preconditions so as to illustrate the approach.

The value of representing the safety precondition as a first order formula, rather than as a propositional formula, is that it is both more succinct and a more accurate representation of the way aviation safety experts reason about the relationship between instrument readings and actions. Because instruments have values drawn from the real number space, it is natural to represent them as variables over numbers, as is possible in first order logic. The previous approach using propositional logic that “chunks” the instruments into distinct predicates of values, like *airspeedlow*, *airspeedmedium*, and *airspeedhigh*, is a clunky way of achieving a less accurate goal.

The propositional approach unnecessarily imposes a trade-off between the space complexity of the encoding and its precision. For example, if each instrument is “chunked” into three predicates, that means there are three propositions for each instrument, a subformula specifying that their truth values are mutually exclusive (because a single instrument cannot have two different readings at the same time), and the safety precondition must explicitly encode each distinct combination of instrument readings that allow the action to be safely executed. Imagine a very precise “chunking” of the instrument space, say each natural number value gets its own proposition. Then there would have to be an explosion of clauses enumerating every permitting configuration of instrument readings. The files would be huge, and the computation would be less efficient.

In this section we describe the first order logic encoding of actions and their safety preconditions, and present case studies as examples. We begin with Air France 447¹.

¹The source code for the case studies can be found here: <https://github.com/sethahrenbach/UnsatCore>

5.2.1 Air France 447 Encoded

Air France flight 447 from Rio de Janeiro, Brazil, to Paris, France, occurred 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 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 [56] 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. The official Bureau d'Enquêtes et d'Analyses pour la Sécurité de l'Aviation Civile (BEA) investigation found that a procedure for safely dealing with unknown airspeeds was known to the pilots but not engaged in, and this may have saved the flight as well.

For this case study, we encode the segment of the mishap chain concerning the inconsistent airspeed indicator readings and the hard nose up pitch command that the pilot continuously executed. We begin by declaring the variables of the scenario.

```
;; Air France 447 Case
(set-option :produce-unsat-cores true)

;; right side instruments
(declare-const right-airspeed Int)
(declare-const right-pitch Int)
(declare-const right-bank Int)

;; middle instruments
(declare-const middle-airspeed Int)
(declare-const middle-pitch Int)
(declare-const middle-bank Int)

;; left side instruments
(declare-const left-airspeed Int)
(declare-const left-pitch Int)
(declare-const left-bank Int)
```

Next we encode the actual instrument readings during the flight. We give these clauses names, so that the unsatisfiable core extraction can refer to them if they are responsible for a logical conflict.


```
;; actual instrument readings
(assert (! (= left-airspeed 135) :named lair))
(assert (! (= right-airspeed 100) :named rair))
(assert (! (= middle-airspeed 100) :named mair))
(assert (! (= right-pitch 15) :named rpitch))
(assert (! (= middle-pitch 15) :named mpitch))
(assert (! (= left-pitch 15) :named lpitch))
(assert (! (= right-bank 15) :named rbank))
(assert (! (= middle-bank 15) :named mbank))
(assert (! (= left-bank 15) :named lbank))
```

Next, we encode the safety precondition of the hard nose up input action. We capture the requirement that a hard nose up should not be engaged in if the indicators for the same aspect of flight, like airspeed, are in disagreement. We also capture some upper and lower bounds between which the indicators must read in order for the input to be considered safe. Normally, these bounds are thought of as the flight envelope.

```

;; safety precondition of action: Hard Nose Up
(assert (= left-airspeed right-airspeed))
(assert (= left-airspeed middle-airspeed))
(assert (= middle-airspeed right-airspeed))
(assert (= right-pitch left-pitch))
(assert (= right-pitch middle-pitch))
(assert (= middle-pitch left-pitch))
(assert (< right-airspeed 700))
(assert (> right-airspeed 99))
(assert (< middle-airspeed 700))
(assert (> middle-airspeed 99))
(assert (< left-airspeed 700))
(assert (> left-airspeed 99))

(check-sat)
(get-unsat-core)

```

The final commands of the .smt2 file tell Z3 to check whether the formula is satisfiable, and if not, to retrieve and print the unsatisfiable core. The command line input and output for running Z3 on the file appear below.

```

input  : z3 af447.smt2
output: unsat
output: (lair rair)

```

After downloading the Z3 executable² and creating the af447.smt2 file, we run Z3

²<https://github.com/Z3Prover/z3/releases>

on the file in a single command line input, and receive two output lines. One line tells us that the formula is unsatisfiable, which means that there is some conflict between the actual instrument readings and the safety precondition for the action. The second line tells us which instrument readings are responsible for the conflict. In this case, the instrument readings *lair* and *rair* are determined to be the conflict. If we look up at the file contents, we see that the clauses with those labels capture the left airspeed indicator, reading 135, and the right airspeed indicator, reading 100. According to the safety precondition, these values should be equal. Otherwise, a hard nose up input is unsafe, because the airspeed is not known.

The unsatisfiable core of the formula, as extracted by the Z3 SMT solver, infers the safety-critical information of the case, because it is that pair of instrument readings that render the safety precondition false, by contradicting the condition that they be in agreement.

At this point, one might fairly wonder why *mair* is not likewise included in the set. This is because it is not necessary for identifying the source of the contradiction. Having found that *lair* and *rair* conflict, these suffice to serve as the unsatisfiable core. And indeed, it is irrelevant whether the core is identified as *lair* and *rair* or *lair* and *mair*, because either will adequately serve as the safety critical information. Both sufficiently demonstrate a conflict in the relevant indicators.

5.2.2 Copa Flight 201 Encoded

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 rolling dive. At some point one of the pilots switched the input to the backup Attitude Indicator to be that of the pilot’s faulty one, furthering confusion. 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 encode the variables the same as before, and leave out the code to save space. A simplification of the actual instrument readings is encoded as follows.

```
;; actual instrument readings
(assert (!=(= left-airspeed 135) :named lair))
(assert (!=(= right-airspeed 135) :named rair))
(assert (!=(= middle-airspeed 135) :named mair))
(assert (!=(= right-pitch 15) :named rpitch))
(assert (!=(= middle-pitch 15) :named mpitch))
(assert (!=(= left-pitch 15) :named lpitch))
(assert (!=(= right-bank -15) :named rbank))
(assert (!=(= middle-bank 0) :named mbank))
(assert (!=(= left-bank 0) :named lbank))
```

Then the safety precondition for a hard right bank would include the condition that the instrument readings agree on the plane's actual bank position.

```

;; safety precondition of action: Hard Right Bank
(assert (= left-airspeed right-airspeed))
(assert (= left-airspeed middle-airspeed))
(assert (= middle-airspeed right-airspeed))
(assert (= right-bank left-bank))
(assert (= right-bank middle-bank))
(assert (= middle-bank left-bank))
(assert (< right-airspeed 700))
(assert (> right-airspeed 99))
(assert (< middle-airspeed 700))
(assert (> middle-airspeed 99))
(assert (< left-airspeed 700))
(assert (> left-airspeed 99))

(check-sat)
(get-unsat-core)

```

When we save the file as `copa201.smt2` and run Z3 on it, we see the following in the terminal.

```

input  : z3 copa201.smt2
output: unsat
output: (rbank lbank)

```

Again, as before, the Z3 SMT solver allows us to infer which safety-critical information in particular the pilot is unaware of, because it is the information that conflicts with the safety precondition. What the pilot's action reveals is that he mistakenly believed,

at least at first, that he was in a steep left bank, because he was relying on the faulty instrument in front of him. Had he grasped the safety-critical information, that the attitude indicators were in stark disagreement, he would have realized that he did not know the plane's bank position. In this case, safety procedure is to maintain a level horizontal input, absent other convincing evidence that a steep left bank is occurring, which there was not in this case. Before executing such an extreme action, proper cross-checking of other instruments must occur.

5.2.3 Asiana Flight 214 Encoded

Asiana flight 214 from South Korea to San Francisco departed in the evening of July 6, 2013 and was scheduled to land just before noon that morning [55]. 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. When the pilot realized the plane's glide path was too high, he accidentally disconnected the autopilot's monitoring of the airspeed. This caused an autothrottle (A/T) protection to turn off, so when he pitched the nose down, the plane descended faster than was safe, causing it to collide with the embankment in front of the runway. We will formalize the moment at which the pilot pitches the nose down, unaware that the autothrottle is turned off.

In encoding this state of affairs, we add a new global state variable to capture the autothrottle status.

```

;; Asiana 214 Case

(set-option :produce-unsat-cores true)

;; global state

(declare-const autothrottle Bool)

```

We likewise include right, middle, and left instruments for altitude in the variable declarations, omitted here for space. The encoding of the actual instrument readings is as follows.

```

;; actual instrument readings

(assert (!(= autothrottle false) :named at-status))

(assert (!(= right-alt 5000) :named ralt))

(assert (!(= middle-alt 5000) :named malt))

(assert (!(= left-alt 5000) :named lalt))

(assert (!(= left-airspeed 200) :named lair))

(assert (!(= right-airspeed 200) :named rair))

(assert (!(= middle-airspeed 200) :named mair))

```

Finally, for the action of a Hard Nose Down, we show how to encode a conditional constraint, that if the altitude is less than 10000 feet, then the autothrottle must be on in order for a Hard Nose Down to safely be executed. Logically, this conditional is equivalent to the proposition that either the altitude is greater than (or equal to, but we omit this in the formalization) 10000 feet, or the autothrottle is on.

```

;; safety precondition of action: Hard Nose Down
(assert (= left-alt right-alt))
(assert (= left-alt middle-alt))
(assert (= middle-alt right-alt))
(assert (= left-airspeed right-airspeed))
(assert (= left-airspeed middle-airspeed))
(assert (= middle-airspeed right-airspeed))
(assert (or (> right-alt 10000) (= autothrottle true)))

(check-sat)
(get-unsat-core)

```

When we save the file as `asiana.smt2` and run `Z3` on it, we see the following in the terminal.

```

input  : z3 asiana.smt2
output: unsat
output: (at-status ralt)

```

This indicates that the safety-critical information is the combination of the altitude, here `Z3` returns only the right side altimeter reading because it is sufficient, and the fact that the autothrottle is off, just as the NTSB report concluded.

5.3 The Unsatisfiable Core is the Safety-Critical Information

This section establishes the identification of the unsatisfiable core extracted from Z3 with the safety-critical information that a pilot is unaware of when she executes an unsafe action. To do this, we examine a case study formalized in DASL, and compare the results with the same case study encoded as an SMT problem that was run through Z3.

Consider the following formalization of Copa flight 201 in DASL. Line 1 represents a conjunction over the actual instrument readings, encoded as atomic propositions. The relevant propositions included in the text state that the right attitude indicator (*RightBank*) reads that the aircraft is pitching left considerably (*SteepLeft*), but that the left attitude indicator (*LeftBank*) reads level flight. We ignore the middle attitude indicator to try and keep it readable, but if it were included, it would be present in the conclusion as well. Line 2 expresses in DASL that the pilot executes a hard right bank action. Line 3 concludes with the theorem’s application, which states that the pilot lacks awareness either the left indicator or the right indicator: she is unaware of the falsity of one of them. This analysis establishes that the attitude indicators’ values together are a piece of safety-critical information, because it shows that the pilot engages in the unsafe action *only if* she remains unaware of them jointly.

1. $(\mathbf{RightBankSteepLeft}) \wedge (\mathbf{LeftBankLevel}) \dots\dots\dots$ configuration.
tion.
2. $\langle \mathbf{pilot} \rangle \mathbf{hardrightbanktrue} \dots\dots\dots$
pilot input.
3. $\neg K_{\mathbf{pilot}} (\mathbf{LeftBankSteepLeft} \vee \mathbf{RightBankLevel}) \wedge$
 $\neg K_{\mathbf{pilot}} \neg K_{\mathbf{pilot}} (\mathbf{LeftBankSteepLeft} \vee \mathbf{RightBankLevel}) \dots\dots\dots$ 1,2, The
1.

Compare this with the results of running the SMT encoding through Z3. The instrument readings construed as a first order formula are unsatisfiable modulo the constraints imposed by the action's safety precondition. The unsatisfiable core identified by Z3, that is, the reason it is unsatisfiable, is because the right attitude indicator (***rbank***) and the left attitude indicator (***lbank***) are in disagreement. This is consistent with the DASL analysis, which identifies the left attitude indicator's value as the one the pilot is critically unaware of.

The way these two formalisms work together is that the DASL theorem establishes the connection between the action, pilot knowledge, and safety conditions, and the SMT solver performs the safety-critical inference in an efficient and expressive way. They are different in an interesting way, though. DASL works by identifying a proposition from the safety precondition of the action that is *false*, not currently present on the instrument panel. For example, a hard right bank action's safe execution requires the attitude indicators to be in agreement. If the attitude indicators disagree, and the pilot executes the action anyway, then DASL allows us to infer unawareness of the false propositions from the safety precondition. The SMT formalization operates by identifying which of the *true* propositions from the instrument panel conflict with the safety precondition of the action. Thus, the result of DASL: the pilot is unaware

regarding the left indicator’s being in a steep left bank or the right indicator’s being level; Z3 infers that the actual instrument readings of the right airspeed indicator and the left airspeed indicator are the ones that the pilot needs to see. We built it this way so that the instruments themselves, not the safety precondition, are the output of the Z3 engine, for reasons to be discussed in the next section.

5.4 Application Design

This section briefly describes how this encoding of aviation activity and the Z3 theorem prover might be used in a runtime monitor application.

The application would execute the following algorithm in a loop.

Algorithm 1 Monitor Algorithm

```

while flying do
  smt2_actual_readings  $\leftarrow$  instrument readings
  smt2_safety_preconditions  $\leftarrow$  lookup_safety_precondition(flight control input)
  smt2_file  $\leftarrow$  concat(smt2_actual_readings, smt2_safety_preconditions)
  unsat_core  $\leftarrow$  Z3 smt2_file
  if unsat_core is empty then undim all
  else if unsat_core is not empty then
    map(dim, [x | x  $\notin$  unsat_core])

```

This algorithm, though simple, represents a dramatic departure from contemporary safety practices, in that it aims to draw the pilot’s attention to the safety-critical information by *dimming* the information that is not, at the moment, safety-critical. Dimming, of course, is not a technical term, and has not been defined. Its definition and implementation depends on the context and the technical details of the cockpit.

In most emergency situations, a number of alarms are going off at once, and the pilots suffer from information overload, which decreases their situational awareness. The safety-critical information is right in front of them, and perhaps even has an alarm and blinking light accompanying it. And yet they still crash. This algorithm

addresses the problem by temporarily *removing* the alarms and blinking lights that are not related to the safety-critical information identified by Z3's returned unsatisfiable core, based on the encoding of the configuration and action. If the only alarm or blinking lights are those highlighting the fact that, say, the airspeed indicators are in disagreement, immediately following the pilot's action to dramatically reduce thrust, then the chances that he or she notices this problem must increase.

This runtime monitor will increase the likelihood that the pilot notices the safety-critical information. Noticing it, he or she will cease the unsafe action.

5.5 Related Work

Over the past decade related work in applying formal methods to the human component of aviation safety have advanced the state of the art and contributed to aviation safety [48, 32, 19, 16]. However, previous work applies formal methods to the specification stage of systems, typically using model checking or SAT/SMT solving to discover possible states of the system in which a human might become disoriented or confused. This work is important. The present work advances a different but related line of effort, applying formal methods to the run-time stage of systems, with the goal of diagnosing and correcting unsafe system behaviors as they occur, especially those caused by human confusion. To do this, we have taken a formal model of the pilot's reasoning in the form of DASL and used the resulting analysis to develop a tool for delivering safety-critical information to the pilot when it is needed.

Chapter 6

Summary and concluding remarks

Congratulations on completing your dissertation.

Appendix A

Title of first appendix

A.1 Section title

Here is some additional information which would have detracted from the point being made in the main article.

A.1.1 Subsection title

This section even has subtitles

Bibliography

- [1] Federal Aviation Administration
website: <https://www.faa.gov/files/gslac/library/documents/2012/Nov/71574/DirtyDozenWeb3.pdf>,
retrieved 10/5/2016.
- [2] Rankin, William. “MEDA Investigation Process”, Aero Quarterly.
Website: Boeing.com/commercial/aeromagazine/articles/qtr_2_07/aero_q207_article3.pdf,
retrieved 10/5/2016
- [3] Aviation Maintenance Technician Handbook. Chapter 14, “Addendum/Human Factors”.
- [4] van Ditmarsch, H.; van der Hoek, W.; Kooi, B. Dynamic Epistemic Logic. Springer. Dordrecht, The Netherlands, 2008.
- [5] Ahrenbach, S.; Goodloe, A. Formal Analysis of Pilot Error Using Agent Safety Logic. *Innovations in Systems and Software Engineering*. **Submitted**.
- [6] 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.

- [7] 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.
- [8] 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.
- [9] 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.
- [10] De Wind, P. “Modal logic in Coq.” M.Sc. University of Amsterdam. 2001.
- [11] Gettier, E. Is Justified True Belief Knowledge? *Analysis* (Vol. 23, pp. 121-123). 1963.
- [12] Blackburn, P.; de Rijke, M.; Venema, Y. Modal Logic. Cambridge University Press. New York, 2001.
- [13] 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.
- [14] Fagin, R.; Halpern J.; Moses, Y.; Vardi, M. Reasoning about knowledge. MIT press, 2004.
- [15] Basin, D., Radomirovic, S. and Schmid, L., Modeling Human Errors in Security Protocols.
- [16] 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.
- [17] Butler, R.; Johnson, S. Formal Methods For Life-Critical Software. *Computing in Aerospace 9 Conference*, (pp. 319-329), San Diego, CA. October 1993.
 - [18] Davis, M.; Putnam, H. "A Computing Procedure for Quantification Theory". J.ACM. 7 (3): 201–215. (1960).
 - [19] Rushby, J., 1993. Formal methods and the certification of critical systems. SRI International. Computer Science Laboratory.
 - [20] Aucher, G. and Schwarzenruber, F., 2013. On the complexity of dynamic epistemic logic. arXiv preprint arXiv:1310.6406.
 - [21] 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).
 - [22] van Benthem, J.; Pacuit, E; Roy, O. Toward a Theory of Play: A Logical Perspective on Games and Interaction. *GAMES* **2011**, 2, 52-86.
 - [23] Prior, A. Time and Modality. Clarendon Press: Oxford, UK, 1957.
 - [24] Kohler, R. Lords of the Fly. The University of Chicago Press: Chicago, USA, 1994.
 - [25] Kripke, S. Semantical Considerations on Modal Logic. *Acta Philosophica Fennica* **1963**, 16: 83-94.
 - [26] Harel, D.; Kozen, D; Tiuryn, J. Dynamic Logic. MIT Press: Cambridge, MA, USA, 2000.

- [27] Hsu, F. Behind Deep Blue: Building the Computer that Defeated the World Chess Champion. Princeton University Press: Princeton, NJ, USA, 2004.
- [28] Lenzen, W. Recent work in epistemic logic. *Acta Philosophica Fennica* **1972**, 30(1): 1-219.
- [29] Moore, G.E. "Moore's Paradox. In Baldwin, T. *G.E. Moore: Selected Writings*. London: Routledge **1993**.
- [30] Negri, S. Kripke completeness revisited. *Acts of Knowledge: History, Philosophy and Logic: Essays Dedicated to Göran Sundholm* **2009**, 247-282.
- [31] Pnueli, A. The Temporal Logic of Programs. *Proceedings of the 18th Annual Symposium on Foundations of Computer Science (SFCS 77)*. IEEE Computer Society.
- [32] Rushby, J. Using model checking to help discover mode confusions and other automation surprises. *Reliability Engineering and System Safety* **2002**, 75, 167-177.
- [33] 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.
- [34] Rushby, J. Trustworthy Self-Integrating Systems. *Springer LNCS*, **2016**, 9581, 19-29.
- [35] 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.
- [36] Baltag, A.; Moss, L.; Solecki, S. The Logic of Public Announcements, Common Knowledge, and Private Suspicions. *TARK '98*, **1998**, 43-56.

- [37] Shoham, Y.; Leyton-Brown, K. *Multiagent Systems: Algorithmic, Game Theoretic and Logical Foundation*. Cambridge University Press: New York, USA, 2008.
- [38] 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.
- [39] van Benthem, J. *Logical Dynamics of Information and Interaction*. Cambridge University Press: New York, USA, 2011.
- [40] van Benthem, J. Modal Logic for Open Minds. Center for the Study of Language and Information, 2010.
- [41] Van Ditmarsch, H. "Knowledge games." *Bulletin of Economic Research* 53.4 (2001): 249-273.
- [42] Fagin, R.; Halpern J.; Moses, Y.; Vardi, M. *Reasoning About Knowledge*. The MIT Press: Cambridge, USA, 2003.
- [43] Harel, D.; Kozen, D. and Jerzy Tiuryn. *Dynamic logic*. MIT press, 2000.
- [44] Hintikka, J. *Knowledge and Belief*. Cornell University Press. Ithaca, USA, 1962.
- [45] Horty, J. *Agency and Deontic Logic*. Oxford University Press. Oxford, UK, 2009.
- [46] Barwise, J.; Seligman, J. *Information Flow: The Logic of Distributed Systems*. Cambridge University Press: New York, USA, 1997.
- [47] van Ditmarsch, H.; van der Hoek, W.; Kooi, B. *Dynamic Epistemic Logic*. Springer: Dordrecht, The Netherlands, 2008.
- [48] 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).
- [49] 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.
 - [50] Klir, G.; Yuan, B. Fuzzy sets and fuzzy logic. Vol. 4. New Jersey: Prentice hall, 1995.
 - [51] Lescanne, P. "Mechanizing common knowledge logic using COQ." *Annals of Mathematics and Artificial Intelligence* 48.1-2 (2006): 15-43. APA
 - [52] De Moura, L.; Bjørner, N. Z3: An efficient SMT solver. *Tools and Algorithms for the Construction and Analysis of Systems*. 2008. Springer.
 - [53] Maliković, M.; Čubrilo, M. "Modeling epistemic actions in dynamic epistemic logic using Coq." *CECHS-2010*. 2010.
 - [54] 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.
 - [55] 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).
 - [56] Palmer, Bill. Understanding Air France 447. William Palmer, 2013.
 - [57] Lescanne, P.; Puisségur, J. "Dynamic logic of common knowledge in a proof assistant." *arXiv preprint arXiv:0712.3146* (2007).

- [58] Aumann, R. "Interactive epistemology I: knowledge." *International Journal of Game Theory* 28.3 (1999): 263-300.
- [59] Hwang, M.; Lin, J. "Information dimension, information overload and decision quality." *Journal of Informatin Science* 25 (1999): 213-218.
- [60] Rescher, Nicholas. "Epistemic Logic." University of Pittsburgh Press. 2005.
- [61] 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.
- [62] 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.
- [63] von Wright, G.H.. "An Essay on Modal Logic." Amsterdam: North-Holland Publishing Company. 1951.
- [64] Williamson, T. *Modal Logic as Metaphysics*. Oxford Universty Press: New York, USA, 2013.
- [65] Wittgeinstein, L. "Philosophical Investigations." Blackwell Publishers. 1953.

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.