# The adequacy criteria for epistemic logics

**Do you have a subtitle?**
**If so, write it here**

**Seth Ahrenbach · Second Author**

**Abstract** Most epistemic logics model agents with an S5 knowledge operator, which precludes unknown unknowns and is defined semantically by an equivalence relation. The S5 knowledge operator is appealing because the equivalence relation matches our intuition that epistemically possible states should be indistinguishable from an agent's perspective, which we call the indistinguishability criterion for ignorance (ICI). However, the S5 knowledge operator is inappropriate for human-like agents, who have unknown unknowns. Thus, the epistemic logician interested in modeling human-like knowledge face a dilemma: deny the existence of unknown unknowns, or satisfy the ICI. This paper defines a multi-modal logic of knowledge and belief that is appropriate for human-like agents, allowing for unknown unknowns, mistaken but justified beliefs, with what we call a subjective equivalence relation.

## 1 Introduction

Recent advances in epistemic logic tend to focus on layering formal machinery of dynamic logics on top of a static base logic that is epistemic. This allows for the development of formal systems that combine agent knowledge with actions, as in Dynamic Epistemic Logic (DEL), and public announcements, as in Public Announcement Logic (PAL). This 'dynamic turn' as it is called has led to a plethora of highly technical logics that involve aspects of

F. Author
first address
Tel.: +123-45-678910
Fax: +123-45-678910
E-mail: fauthor@example.com

S. Author
second address

agency, including knowledge. The static epistemic base almost always consists of an S5 epistemic modal operator. An S5 modal operator is powerful and has very nice technical properties for logicians to work with when proving soundness, completeness, and decidability. However, the benefits of these various dynamic logics of agency are undercut by the fact that an S5 epistemic operator is inappropriate for human-like knowledge, despite its formal niceties. This paper argues for a static epistemic base of epistemic and doxastic operators that can serve as a foundation for agency-relevant dynamic extensions. The static base we propose allows for the sort of ignorance humans tend to have called 'unknown unknowns', maintains a reasonable level of normativity that can support game-theoretic reasoning, and satisfies what we call the indistinguishability criterion for epistemic logic.

This paper proceeds as follows. In Section 2 we provide the necessary formal background a reader might need in order to understand the context of the thesis.

## 2 Background

Epistemic logic grew out of efforts in the 1950's and 1960's to formalize important philosophical concepts. Epistemology is the philosophical study of knowledge, and an epistemic logic is a formal logic for reasoning about knowledge. Most epistemic logics are a type of modal logic, which similarly shares its roots in philosophy, developed for the study of necessity and possibility. A modal logic in its most basic form modifies propositional logic by adding modal operators. The usual modal operators are $\Box$ for "it is necessary that" and $\Diamond$ for "it is possible that". The operators are defined by a relational semantics, sometimes called Kripke semantics, with a set of possible worlds and a binary possibility relation $R$ defined over worlds. A proposition is true or false at each world, and we use the symbols $w \models \varphi$ to say that world $w$ satisfies $\varphi$, or equivalently, that $\varphi$ is true at $w$. A proposition is necessarily true at a world $w$ just in case it is true at all possible worlds relative to $w$: $w \models \Box\varphi$ just in case for all $v$, $wRv$ implies $v \models \varphi$. A proposition is possibly true at $w$ just in case it is true at at least one possible world relative to $w$: $w \models \Diamond\varphi$ just in case for some $v$, $wRv$ and $v \models \varphi$.

Epistemic logic interprets the possibility relation as one of epistemic possibility. This is typically explained as a state of affairs that is possible given what evidence an agent $i$ has. The knowledge operator for $i$ is $\mathbf{K_i}$, and the corresponding possibility operator is $\langle \mathbf{K}_i \rangle$, which is typically somewhat difficult to interpret other than something like "$i$ considers $\varphi$ possible, given her evidence" for the formula $\langle \mathbf{K}_i \rangle \varphi$.

The theorems a modal logic proves depends heavily on the properties of the binary possibility relation which defines the operators. Certain algebraic properties correspond to axioms of the logic. A reflexive possibility relation is one such that

**Definition 1 (Reflexivity)**

$$\textit{For all } w, \; wRw. \tag{1}$$

A transitive possibility relation is one such that

**Definition 2 (Transitivity)**

$$\textit{For all } w, v, u, \; (wRv \wedge vRu \Rightarrow wRu). \tag{2}$$

A Euclidean possibility relation is one such that

**Definition 3 (Euclidean)**

$$\textit{For all } w, v, u, \; (wRv \wedge wRv \Rightarrow vRu). \tag{3}$$

The reflexive and euclidean properties together entail a symmetric property, which is

**Definition 4 (Symmetry)**

$$\textit{For all } w, v \; (wRv \Rightarrow vRw). \tag{4}$$

*Proof* We assume 1 and 3 properties hold of relation $R$. Then we have for all $w, v$, $wRw$ by reflexivity, and $wRw \wedge wRv \Rightarrow vRw$, by 3. Therefore, for all $w, v$, $wRv \Rightarrow vRw$. QED

When a knowledge operator is defined with a possibility relation that is *reflexive* and *Euclidean*, it becomes an equivalence relation, defined as a relation that is reflexive, transitive, and symmetric. We briefly lay out the syntax of epistemic logic and provide English characterizations of the symbols.

$$\varphi := p \mid \varphi \wedge \varphi \mid \neg\varphi \mid \mathbf{K_i}\,\varphi \mid \langle \mathbf{K}_i \rangle\,\varphi.$$

The remaining Boolean operators are defined and interpreted as usual. $\mathbf{K_i}\,\varphi$ reads in English that "agent $i$ knows that $\varphi$ is the case." This should be unproblematic. $\langle \mathbf{K}_i \rangle$, however, is not easily interpretable into English. This can result in rather verbose translations of formulas it appears in. For example, $\langle \mathbf{K}_i \rangle\,\varphi$ might be translated as "$i$ cannot rule out that $\varphi$ might be the case, given her evidence". Imagine this translation for a formula of any real complexity! One might translate it as "$i$ considers $\varphi$ possible, given her evidence", but we resist this translation for reasons to be clear later. Another common translation is, "for all $i$ knows, it may be that $\varphi$." We interpret this translation and the evidence-based translation as equivalent, but we use the following shorthand to capture them: "*i may have* that $\varphi$."

Later in the paper we consider logics with modal operators for belief. These operators are $\mathbf{B_i}$ and $\langle \mathbf{B}_i \rangle$. $\mathbf{B_i}$ has an intuitive interpretation as "$i$ believes that $\varphi$. We reserve "$i$ considers it possible that $\varphi$" for $\langle \mathbf{B}_i \rangle\,\varphi$, and we will typically shorten this to "$i$ considers that $\varphi$."

We choose these interpretations because there is something more objective about when something may be the case for $i$, and something more subjective about when $i$ considers something possible, and we think this divides nicely along epistemic and doxastic (belief) lines.

We proceed with our criteria for assessing an epistemic logic.

## 3 Criterion 1: Descriptive Accuracy

A formal system first and foremust must accurately formalize some target phenomenon. Propositional logic formalizes and accurately describes valid inferences among propositions. First order logic formalizes and accurately describes valid inferences among quantified formulas with objects and predicates. Probability theory formalizes and accurately describes reasoning about the relations among probabilistic events. Can the same be said about epistemic logic for reasoning about knowledge?

As with modal logic in general, there are infinitely many ways to define an epistemic logic, because there are infinitely many first order conditions one an impose on the epistemic operator. The most common epistemic logic uses an S5 operator defined by a possibility relation that is Euclidean and transitive, and therefore is an equivalence relation. This corresponds to the following axioms schema.

| | |
|---|---:|
| $\mathbf{K_i}\,(\varphi \Rightarrow \psi) \Rightarrow (\mathbf{K_i}\,\varphi \Rightarrow \mathbf{K_i}\,\psi)$ | Distribution of $\mathbf{K_i}$ |
| $\mathbf{K_i}\,\varphi \Rightarrow \varphi$ | Truth Axiom |
| $\neg\mathbf{K_i}\,\varphi \Rightarrow \mathbf{K_i}\,\neg\mathbf{K_i}\,\varphi$ | Negative Introspection |
| 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 1** S5 Axiom Schema

Distribution of $\mathbf{K_i}$ is an axiom schema that every normal modal logic contains, which is a logic that uses possible world semantics with standard quantification over worlds.

The Truth Axiom requires that known propositions be true, which is considered a distinguishing feature of knowledge.

Negative introspection presents a problem. Hintikka considered this formula to be clearly unacceptable and did not include it in his system of epistemic logic. It states that an agent $i$ does not know something only if she knows that she does not know it. Presented in this form, it places a strong knowledge condition on ignorance. One can get a good sense of a potential axiom by considering its equivalent forms. Consider its contrapositive.

$$\neg \mathbf{K_i} \neg \mathbf{K_i} \varphi \Rightarrow \mathbf{K_i} \varphi. \tag{5}$$

That is to say, if $i$ does not know that she does not know $\varphi$, then she knows $\varphi$. By the definition of $\langle \mathbf{K}_i \rangle$, this is equivalent to

$$\langle \mathbf{K}_i \rangle \mathbf{K_i} \varphi \Rightarrow \mathbf{K_i} \varphi. \tag{6}$$

This states that if $i$ has MAYBES that she knows $\varphi$, then she knows it.

A theorem of S5 is

$$\mathbf{K_i} \varphi \Rightarrow \mathbf{K_i} \mathbf{K_i} \varphi, \tag{7}$$

Which holds that knowledge implies positive introspective knowledge. (7) is in fact referred to as positive introspection. By focusing on the fact that positive introspection here is a necessary condition for what we might call first order knowledge, we can see that if we are ever in doubt about whether we know something, then it follows that we do not know it. Many philosophers find this property to be counterintuitive, and hold that such a necessary condition on knowledge would have the effect of incorrectly ruling out cases of genuine knowledge.

Yet another counterintuive theorem of S5 is the so-called B Theorem, after L.E.J. Brouwer,

$$\varphi \Rightarrow \mathbf{K_i} \langle \mathbf{K}_i \rangle \varphi. \tag{8}$$

If the reader considers what this mean, she should recognize that S5 is a striking mischaracterization of human-like knowledge. We consider two ways in which S5 inaccurately describes human-like knowledge, other than its inclusion of (8). First, that it denies that unknown unknowns exist. Second, that it denies that unknown knowns exist.

### 3.0.1 Unknown Unknowns

The negative introspection axiom denies that unknown unknowns exist for the agents it models. We argue that such agents are not human-like, and that this counts against S5 epistemic logic on grounds of accuracy.

The reason we say S5 denies that unknown unknowns exists is because the negative introspection axiom states that "if $i$ does not know that $\varphi$, then $i$ knows that she does not know that $\varphi$." In this sentence, $\varphi$ is the unknown, and by the axiom this implies that she knows that it is unknown to her.

This is inaccruate for human-like knowledge, and it should be clear that this is so to the reader, but we shall make the case anyway. First, we need not show that all unknowns are unknown unknowns in order to invalidate the axiom; we must show only that sometimes human-like knowledge involves unknown unknowns. Therefore, a simple counterexample will do. Suppose Alice believes that there is milk in the fridge, but that this is false. Bob finished off the milk earlier. Thus, the proposition being false, it is not the case that Alice knows it. But, since she believes it to be the case, she clearly does not know that she does not know it. She is totally unaware of her ignorance regarding the milk. This suffices to show that the axiom does not describe many common occurences for human-like knowledge.

Another case concerns a true belief that is unjustified. Suppose Alice is watching a mystery film, and Bob walks in near the beginning, and asks Alice what she is watching. She tells him it is a who-dunnit. He looks at the screen and predicts that the character $M$ is the murderer, without knowing anything about the story. He simply sees $M$ on the screen, and forms the confident belief that $M$ is the murderer. It turns out he is correct, but he had no evidence to base his prediction on. He'd have guessed whichever character was on the screen when he walked in. When the film reveals the murderer is $M$, Bob proclaims that he knew it, which is false. He did not know it, and he did not know that he didn't know it.

A final case is perhaps most interesting. It exploits the human condition of not being aware of all the propositions in the world. Alice and Bob are firefighters, and their fire station receives an emergency call. Approaching the burning building, there are many things that they know they do not know. They know that they don't know if there are any people trapped in there. They know that they don't know what caused the fire. But they don't know that the construction company violated regulations when erecting the structure, and that these decisions led to a foundation that is more flammable than the remaining structure. They do not know that they don't know this, because it does not even occur to them; it is a highly irregular situation. The proposition that "the construction company built the foundation out of a material more flammable than the remaining structure" is not one that they think to assess their evidence for. If you mention it to them as a possibility, they would rightly dismiss it as highly unlikely. This proposition is an unknown unknown.

### 3.0.2 Unknown Knowns

Positive introspection imposes a necessary condition on knowledge such that the agent must know that she knows a proposition. We present an argument against the principle of positive introspection based on taking the higher order character of the principle seriously.

Before presenting our argument, we consider a regular case of knowledge. Suppose Alice knows that the tile in front of her is a rectangle. In order to know this, Alice must know the correct defnition of a rectangle, and she must have evidence that the tile she is looking at satisfies the definition, namely that it is a four sided polygon with 90° vertices. This should be uncontroversial.

However, if we apply the same standard to the higher order question, whether Alice knows that she knows that the tile she is looking at is a rectangle, she must now know the correct definition of knowledge, and have evidence that her mental attitude toward the proposition satisfies that definition. With this in mind, we present the argument that achieving higher order knowledge about knowledge is not possible for most human-like agents, whether in reality or described by a formal system.

For the case of reality, we have Alice and the rectangle. She knows that the tile is a rectangle. Unfortunately, her notion of what constitutes knowledge is that it is a justified true belief; she has not read Gettier. Furthermore, when she seeks evidence that the proposition "the tile I am looking at is a rectangle" is true, she realizes she must now identify the correct notion of truth, and identify evidence that that the proposition satisfies those conditions. Maybe truth is just as Tarski said, and the proposition is true just in case the tile is a rectangle. But maybe not. Whichever definitions of knowledge and truth are correct, we suppose Alice subscribes to a competing view, and therefore she does not have higher order knowledge.

For formal agents defined with knowledge that includes the positive introspection property, things get more interesting. As with Alice, the agent, call it *AlphaKnow*, must have the correct definition of a rectangle, and evidence demonstrating that the tile satisfies that definition. When deliberating on its own knowledge status, *AlphaKnow* might crack open its own hood and see that it knows that the tile is a rectangle just in case "the tile is a rectangle" is true in all possible worlds consistent with the evidence, which is the definition of its own knowledge. It sees that a sentence is true at a possible world just in case the world satisfies the sentence according to the inductively defined satisfaction relation. It sees further that knowledge entails truth, because it sees that the its possibility relation is reflexive, guaranteeing that the actual world is in the equivalence relation. It computes that "the tile is a rectangle" is true in all possible worlds consistent with its evidence. Therefore, it knows that the tile is a rectangle, and having reasoned to this conclusion soundly, it knows that it knows it.

In what follows, we present *AlphaKnow*'s reasoning formally, with the outermost $\mathbf{K_{ak}}$ referring to its object level reasoning, and everything within its scope is the content of its object level reasoning. Thus, *AlphaKnow* is capable of self-reference. Let $\varphi$ be the proposition that "the tile is rectangle". We represent object level propositional encoding of the metalanguage with $\ulcorner$these$\urcorner$. *AlphaKnow* reasons:

1. $\mathbf{K_{ak}}\left(\mathbf{K_{ak}}\,\varphi \Leftrightarrow \ulcorner(\forall v, wR_k^{ak}v \text{ implies } v \models \varphi) \text{ and } wR_k^{ak}w\urcorner\right)$
2. $\mathbf{K_{ak}}\left(\ulcorner wR_k^{ak}w\urcorner\right)$

3. $\mathbf{K_{ak}}\,(\ulcorner \forall v, w R_k^{ak} v \text{ implies } v \models \varphi \urcorner)$
4. $\mathbf{K_{ak}}\,(\mathbf{K_{ak}}\,\varphi)$
5. $\mathbf{K_{ak}}\,(\mathbf{K_{ak}}\,\mathbf{K_{ak}}\,\varphi)$

So, *AlphaKnow* reasons that if it knows that knowing the definition of knowledge and that the tile is a rectangle in all possible worlds entails that the tile is a rectangle, then it knows that the tile is a rectangle. Therefore, it knows that the tile is a rectangle, and because it deduced this by valid inference rules, it knows that it knows this.

Unfortunately, if *AlphaKnow* is capable of such self-referential reasoning, it can deduce the following theorem,

$$\mathbf{K_{AlphaKnow}}\,(\mathbf{K_{AlphaKnow}}\,\varphi \Rightarrow \varphi) \Rightarrow \mathbf{K_{AlphaKnow}}\,\varphi, \qquad (9)$$

which is better known as Löb's Theorem. An agent with Löb's Theorem in its head will always end up knowing false propositions. Agent foundations researchers refer to this as the *Löbian Obstacle*. Smullyan identified it as a problem for agents modeled by doxastic logic which we address in a future section. The obstacle presents itself to agents who (1) can conceive of self-referential sentences like, "if I know that this sentence is true, then $\varphi$," (2) are normal modal reasoners, (3) have a reflexive and transitive modal possibility relation. (2) and (3) obviously hold for S5 epistemic agents, and we make the assumption that by being human-like reasoners they have a sufficient expressive language buried beneath our propositional abstraction.

Thus, they face the Löbian Obstacle, which is, for an arbitrary $\varphi$:

| | | |
|---|---|---|
| (3.1) | $\mathbf{K_{ak}}\,(\psi \Leftrightarrow (\mathbf{K_{ak}}\,\psi \Rightarrow \varphi))$ | Löb Sentence[1] |
| (3.2) | $\mathbf{K_{ak}}\,(\mathbf{K_{ak}}\,\psi \Leftrightarrow \mathbf{K_{ak}}\,(\mathbf{K_{ak}}\,\psi \Rightarrow \varphi))$ | Distribution |
| (3.3) | $\mathbf{K_{ak}}\,(\mathbf{K_{ak}}\,\psi \Rightarrow (\mathbf{K_{ak}}\,\mathbf{K_{ak}}\,\psi \Rightarrow \mathbf{K_{ak}}\,\varphi))$ | Taut, Distribution |
| (3.4) | $\mathbf{K_{ak}}\,(\mathbf{K_{ak}}\,\psi \Rightarrow \mathbf{K_{ak}}\,\mathbf{K_{ak}}\,\psi)$ | Positive Introspection |
| (3.5) | $\mathbf{K_{ak}}\,(\mathbf{K_{ak}}\,\psi \Rightarrow \mathbf{K_{ak}}\,\varphi)$ | Taut (2.4) |
| (3.6) | $\mathbf{K_{ak}}\,(\mathbf{K_{ak}}\,\varphi \Rightarrow \varphi)$ | Truth Axiom |
| (3.7) | $\mathbf{K_{ak}}\,(\mathbf{K_{ak}}\,\psi \Rightarrow \varphi)$ | Taut (2.5) |
| (3.8) | $\mathbf{K_{ak}}\,\psi$ | Taut (2.1) |
| (3.9) | $\mathbf{K_{ak}}\,\varphi$ | Taut (2.8), (2.5) |

$$\mathcal{QED}$$

Thus, if a formal agent *really* has positive introspection, and can look inside itself to see whether the conditions for knowledge are satisfied, and reason in a human-like way, then it will reason itself into falsehoods. If formal agents lack positive introspection of the above sort, but their knowledge is governed by the positive introspection property, then they cannot know any proposition. Therefore, the positive introspection axiom is inappropriate for representing human-like knowledge in both humans and human-like machines.

---

[1] The literature sometimes refers to this as a Curry Sentence, after the logician Haskell Curry, who used it to derive similar paradoxes in the lambda calculus.

This section has shown that negative and positive introspection are inappropriate properties for human-like knowledge. Since S5 epistemic logic includes these properties, we conclude that it does not perform well by the criterion of accuracy as a formal system, at least regarding human-like knowlege.

The next section presents the criterion of normativity.

## 4 Criterion 2: Normativity

Just as propositional logic accurately describes the valid rules of inference for reasoning about propositions, it also has a significant normative aspect. Philosophers, lawyers, mathematicians, and anyone who takes reasoning and argumentation seriously, learns the basic validities of propositional logic, as well as perhaps first order logic and the classical syllogisms. They use these rules to guide their own thinking and to evaluate the arguments presented in defense of positions. Human thinking does not naturally abide by all the validities of propositional, first order, and syllogistic logic, so knowing these formal rules can serve as a way to improve one's own thinking through conscious effort.

A similar normativity is present in probabilistic reasoning. The axioms and theorems of probability theory describe the valid probabilistic inferences relating to conjunction, disjunction, negation, and conditionals. The normative aspect of probability theory evaluates human probabilistic inferences in terms of how well they abide by these rules, and a human can consciously avoid fallacious reasoning by keeping in mind the theorems of probability theory. One does better to make inferences in accordance with probability theory, both in terms of expected utility and in terms of forming accurate degrees of belief. This thesis does not address degrees of belief, so we abstract away from probability theory, but the principle stands that correct formalisms for describing a method of reasoning tend to have a normative character.

Epistemic logic, and modal logics generally, do not tend to have the same normative character. The normative aspects of S5 epistemic logic might be due to the fact that more knowledge allows for better decisions, as is the case with ideally rational agents in game theory. If you can approximate the knowledge conditions of these agents, and use an optimal decision procedure to choose which action to take, then you can approximate the optimal outcomes of ideally rational agents.

Consider a case in which Bob does not know there is a car in his blind spot on the highway. Bob is driving at a reasonable speed, and checked his blind spot a reasonable time ago. But there is a car there now because it is moving unreasonably fast. Bob is about to change lanes. He reasons, "I checked my blind spot a short time ago and saw no car, so I have recent evidence that no car is there, but I should check again to make sure." Has Bob reasoned invalidly? The contrapositive of negative introspection states that if Bob *may*

*have*[2] that he knows that the lane next to him is clear, then he knows it. But Bob identified that he may have that he knows, but does not know. According to S5 epistemic logic, his reasoning is invalid. But he clearly is better off by reasoning this way and checking his mirror again before changing lanes.

Thus, in this case, if Bob were to reason by S5 epistemic logic, he would be strictly worse off, because he would infer that he knows the lane next to him is empty from his maybe having knowledge of it. This is a mark against S5 epistemic logic as a normative standard to strive for.

Next we present the criterion of indistinguishability.


## 5 Criterion 3: 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 capture the idea of indistinguishability because each world in the equivalence class is indistinguishable from each other relative to the relation defining the class. 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.

This property of S5 epistemic logic is a strong point in its favor, despite the shortcomings detailed in the previous section. Any formal system seeking to model human-like knowledge must surely represent the idea that an agent does not know *whether* a proposition $\varphi$ is true or false just in case he cannot distinguish the world he is in from the possible worlds in which either is the case. For epistemic logic, this condition, which we call the *indistinguishability criterion*, applies to the underlying epistemic relation. S5, having an equivalence relation, clearly satisfies the criterion.

Hintikka's system described in the seminal work on epistemic logic did not meet this criterion. In rejecting the negative introspection axiom, Hintikka formalized knowledge as an S4 operator, which is weaker than S5, with a reflexive and transitive epistemic relation. This is known as a preorder relation, an instance of which is the familiar $\leq$ relation. For this relation, it is not the case that an agent cannot distinguish related worlds from each other. However, this relation does nicely formalize the notion that implicit knowledge is available to someone in her own head as a sort of *reachability* relation. The reachability relation of directed acyclic graphs is a preorder. One could consider the process of trying to remember a piece of information as stepping down the epistemic possibility relation in a directed acyclic graph in one's head, making inferences and recollections and transitioning from a state of

---

[2] Recall this is our shorthand for $(\langle \mathbf{K}_{Bob} \rangle)$

not knowing, ruling out possibilities along the way, until reaching a state of knowledge. However, this representation has a temporal quality not shared by the indistinguishability criterion, nor represented by the formal system itself, which is static. Furthermore, it does not represent the idea that the agent at present has all the information that allows her to rule out the possibilities inconsistent with the proposition she implicitly knows. In the world at the root of the preorder chain, she does not know the proposition, by assumption.

Therefore, S5 epistemic logic, specifically the negative introspection axiom, allows for this very appealing notion of epistemic indistinguishability. Anything weaker than an equivalence relation will not capture this notion. It is a benchmark against which candidate epistemic logics should be measured, along with the ability to accurately describe human-like knowledge and offer normative criteria for evaluating epistemic states.

## 6 Adding Belief

One aspect of human-like knowledge that we have not mentioned formally, but have referenced informally, is belief. To many philosophers, an agent's believing a proposition is a necessary condition for her knowing it. S5 epistemic logic on its own makes no reference to belief. Agents know things or not, and they know whether they know them or not. If one wishes to formally model an agent's knowledge and belief, one must devise an epistemic logic that synthesizes the two notions. The modal logic for reasoning about beliefs is called *doxastic logic*. We refered to synthesized logics for reasoning about knowledge and belief as epistemic logic, though, because these logics are still logics for reasoning about knowledge, of which belief is a component.

In formalizing knowledge and belief, a logician identifies a set of axioms for the knowledge operator, a set of axioms for the belief operator, and a set of axioms logically relating the two operators. This creates an epistemic logic that synthesizes knowledge and belief. Human-like agents sometimes have false beliefs, so an epistemic logic that can formalize inferences about knowlede, beliefs, and particularly false beliefs, increases its accuracy in describing human-like knowledge. If it can do so while maintaining a normative character, it maintains value as a normative formal system. We examine various attempts at synthesizing knowledge and belief, evaluating them according to the above criteria, and then present our own attempt.

### 6.1 Species of Logics for Knowledge and Belief

In examining the standard epistemic logic presented in the current literature, which we referred to as S5 epistemic logic due to the properties of its modal operator, we identified the following desirable properties of an epistemic logic. First, the epistemic logic should accurately describe human-like knowledge. Second, the epistemic logic should have a normative character, such that

agents who violate its theorems are properly criticizable in some way, and it should offer guidance to agents in how to improve their epistemic states by abiding by the theorems. Finally, the epistemic logic must capture the intuition that epistemically possible worlds should be indistinguishable from each other, which requires an equivalence relation. This may seem to be a set of criteria that are jointly unachievable. Before drawing that conclusion, we shall consider some other epistemic logics and asses them by these criteria.

## 6.2 Hintikka

## 6.3 KL

## 6.4 Autoepistemic

## 6.5 GC Logic

*Paragraph headings* Use paragraph headings as needed.

$$a^2 + b^2 = c^2 \tag{10}$$

## References

1. Author, Article title, Journal, Volume, page numbers (year)
2. Author, Book title, page numbers. Publisher, place (year)