

Evidentialist Logic

Matthew P. Wampler-Doty

Contents

1	Philosophy	4
1.1	Forward	4
1.2	Thermometers	5
1.3	Explicit Justification	7
1.4	Sketch	8
1.5	The Human Condition	11
1.6	Soundness	12
1.7	Descartes	13
1.8	Contradictions	13
1.9	Irrationality	15
1.10	Quine	16
1.11	Closing Remarks	18
2	Introduction to EvIL	19
2.1	Elementary EvIL	19
2.1.1	Grammar & Semantics	19
2.1.2	Intuitions	22
2.1.3	Validities	24
2.2	Basic EvIL	26
2.2.1	Elimination	26
2.2.2	Multiple Agents	29
2.2.3	Kripke Structures	31
2.3	EvIL Completeness	35
2.3.1	Axiom Systems	36
2.3.2	Subformula Model Construction	38
2.3.3	Bisimulation	41
2.3.4	Abstract Completeness	41
2.3.5	Translation	41

2.3.6	Completeness	42
2.3.7	Conservativity, Decidability & Complexity	43
3	Applications	45
3.1	Collapse	45
3.2	Epistemic Plurality	45
3.2.1	Different Kinds of Knowledge	45
3.2.2	Moore's Paradox	45
3.2.3	Fitch's Paradox	45
3.3	Intuitionistic Logic	45
3.3.1	The Gödel Tarski McKinsensy Embedding	45
3.3.2	Knowledge	45
3.3.3	Imagination	45
3.3.4	van Benthem $S4$	45
3.3.5	ImK_{\Box}	45
4	Epilogue	45
4.1	Comparison to Other Approaches	45
4.2	Failures	45
A	Grammars	46
B	Alternate Semantics	46
C	An Application of Pure Model Theory to $EviL$ Semantics	50
	References	52

1 Philosophy

1.1 Forward

The idea of applying modal logic to the study of knowledge more or less began with [Hin69], *Knowledge and Belief*, by. In this it is suggested that one can use the possible world framework of modal logic to model ideal logical agents, and reason about concepts like knowledge and belief as modal boxes. In Hintikka’s text, some philosophical emphasis is put on the ideas of *introspection*, which have two formulations:

- Positive: $\Box\phi \rightarrow \Box\Box\phi$ - “If the agent knows a fact, then she knows that she knows this.”
- Negative: $\Diamond\phi \rightarrow \Box\Diamond\phi$ - “If the agent does not know a fact, she knows that she does not know this.”

Intuitively, the second idea seems like something one ought to reject outright. Many will recall the somewhat famous piece of sophistry put forward by former US secretary of defense Donald Rumsfeld [Rum02]:

Reports that say that something hasn’t happened are always interesting to me, because as we know, there are known knowns; there are things we know we know. We also know there are known unknowns; that is to say we know there are some things we do not know. But there are also unknown unknowns – the ones we don’t know we don’t know. And if one looks throughout the history of our country and other free countries, it is the latter category that tend to be the difficult ones.

This quote was ultimately part of a rather hideous justification for criminal military action. Still, it is undeniably at variance with negative introspection; and despite its malicious intent, it is compelling. Like Rumsfeld, Hintikka also rejects negative introspection. Furthermore, it would seem from the above quote that Rumsfeld does not reject positive introspection. Hintikka does not either, and goes one step further to endorse it explicitly.

Despite philosophical objections, the received view in modern epistemic logic embraces both negative and positive introspection. In addition, the following axiom is also embraced:

- Reflection: $\Box\phi \rightarrow \phi$ - “If the agent knows a some statement, that statement is true”

These three axioms together, along with the axioms and rules of elementary modal logic, form C.I. Lewis’ system *S5* [LL51]. Under correspondence theory, these axioms express that the underlying modal accessibility relation is an equivalence relation. That is, they express that the ideal agent under investigation has partitioned her state space into *information states*. It is well known that game theory shares an equivalent notion of information states (see, for instance, [Hal99] and [Rub98], chapter 3).

Although this view of knowledge has been the focus of exhaustive research, even finding industrial application such as in [AHV02] and [HMdV05, H MV04], it is natural for a practitioner to hold lingering philosophical concerns. The purpose of this thesis is to present a framework which tries to address some of these issues. It shall conform to the following structure:

- §1 First, we shall elaborate the philosophical issues that will be addressed, and sketch an epistemic logic framework which tackles them
- §2 Next, we give formal details of the system we will develop. This section climaxes in the exposition of a completeness theorem.
- §3 Third, we shall look at philosophical applications of the framework developed. It shall be revealed that *intuitionistic logic* is profoundly linked to the perspective on epistemic logic we shall investigate here.
- §4 Finally, the framework developed shall be compared to other approaches.

That being said, we shall now turn to investigating the philosophical issues with epistemic logic, and suggest a manner of addressing them.

1.2 Thermometers

Imagine a 1 m^3 box with a thermometer sealed hermetically inside, as in Fig. 1. Further, pretend that the thermometer reads 290 Kelvin. How many moles of gas are in the chamber?

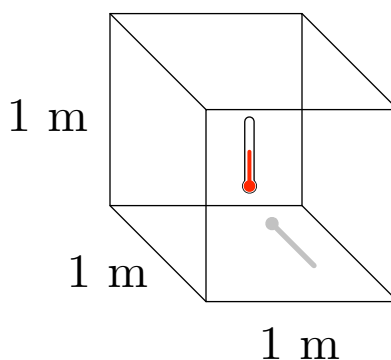


Figure 1: A thermometer in a box

The answer is indeterminate. Recall the *ideal gas law*, originally discovered by Émile Clapeyron [Cla34], which in modern parlance it reads:

$$PV = nRT$$

Where:

- P is the pressure in Pascals
- V is the volume in cubic meters
- n is the number of moles of gas
- T is the temperature in Kelvins

- R is the *ideal gas constant*, $\approx 8.3 \text{ J}\cdot\text{K}^{-1}\cdot\text{mol}^{-1}$

Applying the ideal gas law, one can observe that the thermometer is effectively in something analogous to an epistemic space. To be explicit, consider a basic modal language with the following grammar $\mathcal{L}_{\text{therm}}$:

$$\phi ::= x \text{ Pascals} \mid y \text{ moles} \mid z \text{ Kelvin} \mid \phi \rightarrow \psi \mid \perp \mid \Box \phi$$

Interpreting this language appropriately, the thermometer is evidently an epistemic agent in an $S5$ model. One model for the thermometer-in-a-box is the triple $\langle W, V, \sim \rangle$, where:

- W is pairs (P, n) where P is some positive pressure in Pascals and n is some positive number of moles.
- V is defined as follows:
 - $(P, n) \in V(x \text{ Pascals})$ if $P = x$
 - $(P, n) \in V(y \text{ moles})$ if $n = y$
 - $(P, n) \in V(z \text{ Kelvin})$ if $z = \frac{P}{n \cdot R}$
- Finally, $(P, n) \sim (P', n')$ holds if and only if $P \cdot n' = P' \cdot n$

2 provides a visual representation of the information states in the above model, which form rays emanating from the origin.

The view in philosophy of mind that thermometers are epistemic agents originates Daniel Dennett's *The Intentional Stance* [Den98], with the proviso that Dennett's original discussion originates around thermostats rather than thermometers as we have argued. Dennett writes:

It is not that we attribute (or should attribute) beliefs and desires only to things in which we find internal representations, but rather that when we discover some object for which the intentional strategy works, we endeavor to interpret some of its internal states or processes as internal representations. What makes some internal feature of a thing a representation could only be its role in regulating the behavior of an intentional system.

Now the reason for stressing our kinship with the thermostat should be clear. There is no magic moment in the transition from a simple thermostat to a system that really has an internal representation of the world around it. The thermostat has a minimally demanding representation of the world, fancier thermostats have more demanding representations of the world, fancier robots for helping around the house would have still more demanding representations of the world. Finally you reach us.

The aim of epistemic logic is to model agents modeling the world; and in doing so its development mirrors the increasing levels of complexity that Dennett illustrates. To give an exemplar of the modern level of sophistication achieved by epistemic logic, consider probabilistic dynamic epistemic logic as developed in [vBGK09].

Moreover, just as thinking about thermostats serves as a vehicle for philosophy of mind for Daniel Dennett, thinking about thermometers can elucidate intuitions behind basic epistemic logic, and

how it might be extended. Imagine we were to go up to one of the agents modeled in basic epistemic logic, and ask her why she knows some proposition ϕ . What would she possibly say? She would say she feels ϕ with every fiber of her being, that it is true in every world she can conceive. The reason that ϕ occurs to the agent is because it is what her sensory instruments (modeled as her accessibility relations) tell her. In this respect, the analogy of the thermometer seems pretty apt.

Some natural philosophical features of knowledge cannot be captured by this basic approach. If we were to ask a person on the street why she believes a proposition ϕ , an appeal that she cannot imagine or conceive of the contrary as possible would not slake our curiosity. We typically want some kind of *explanation*, especially if ϕ were a piece of mathematics, for instance. It would certainly make the enterprise of mathematics far simpler if proving theorems amounted to exhibiting that their negation is not imaginable. This gives rise to the following philosophical observation:

Anti-Thermometer Principle: *Traditional epistemic agents, like thermometers, don't always have knowledge, since one must sometimes have reasons for the things they believe.*

How might epistemic logic be saved by the objection that the above principle proposes? To find out, we turn to diagnosing the underpinning of the above principle.

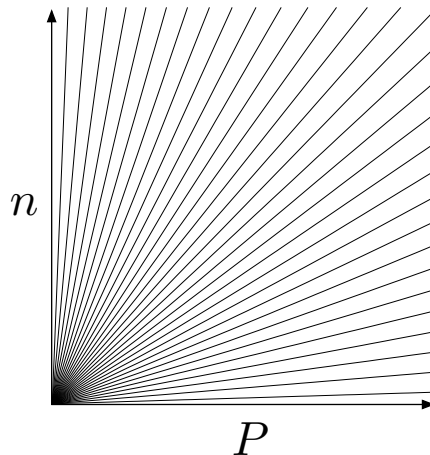


Figure 2: Thermometer information states

1.3 Explicit Justification

I hold that the *Anti-Thermometer Principle* is an expression of the following more fundamental idea:

Justification Principle: *In order to know something, one must sometimes demonstrate inferential justification.*

Demanding the *Justification Principle* is tantamount to demanding a sort of *explicit justification* for beliefs in epistemic logic. This particular desiderata been suggested previously. The hunt for logics of explicit justification was initiated in [vB91]¹. One framework which has been proposed to achieve this is *Justification Logic* [AN05, Art07, Fit04, Fit05]. Alternative frameworks for reasoning about implicit/explicit information have also been proposed in [vBV09] and [Vel09].

In this text we shall investigate a novel framework in this line of research, in the hopes of offering further response that practitioners of epistemic logic may exhort in answer to the problems raised by the aforementioned *justification principle*. To model beliefs with justifications, we shall modify the semantics of modal logic to incorporate certain *basic beliefs*, which we should interpret as non-inferentially justified. These basic beliefs then inferentially generate the rest of what the agent believes.

This perspective amounts to what is called *classical foundationalism* in the philosophical literature. Richard Fumerton describes the view as follows:

[A] foundationalist is someone who claims that there are noninferentially justified beliefs and that all justified beliefs owe their justification, ultimately, in part, to the existence of noninferentially justified beliefs². A belief is noninferentially justified if its justification is not constituted by the having of other justified beliefs. [DeP01, pg. 3]

The view presented would not deny that thermometers or traditional epistemic agents have knowledge, no more than it would deny that one has knowledge of that her right hand has five fingers attached to the end. Rather, our aim is to try to modify the semantics of epistemic logic, more specifically *doxastic logic*, with the ingredients for a foundationalist analysis of knowledge. As will be demonstrated, this can be done without modifying the basic modal logic syntax.

1.4 Sketch

In this section we shall see a very informal presentation of the basic elements which shall compose the forthcoming analysis. A formal development of the ideas sketched in this section shall be given in §2.1.1. With this proviso, consider the basic modal language $\mathcal{L}_K(\Phi)$:

$$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp \mid \Box \phi$$

Further, let $\mathfrak{M} \subseteq \wp\Phi \times \wp\mathcal{L}_K(\Phi)$, that is, let \mathfrak{M} be pairs of sets of letters and formulae. Define the following truth predicate \models recursively:

¹While this paper is considered seminal, it should be remarked that research into this subject began prior to it. Specifically, the phrase “explicit belief” appears to have its origins in [Lev84].

²In the proceeding discussion, we shall refer to *noninferentially justified beliefs* as *basic beliefs* or premises.

Definition 1.4.1.

$$\mathfrak{M}, (a, A) \models p \iff p \in a$$

$$\mathfrak{M}, (a, A) \models \phi \rightarrow \psi \iff \mathfrak{M}, (a, A) \models \phi \text{ implies } \mathfrak{M}, (a, A) \models \psi$$

$$\mathfrak{M}, (a, A) \models \perp \iff \text{False}$$

$$\mathfrak{M}, (a, A) \models \Box \phi \iff \text{for all } (b, B) \in \mathfrak{M}, \mathfrak{M}, (b, B) \models A \text{ implies } \mathfrak{M}, (b, B) \models \phi$$

Since the semantics like the above shall be the principle objects of study, we will give how to read them philosophically. In these semantics, instead of thinking of every world individually, we think of every world as containing facts and a part of the agent's mind. This part of the agent's mind is represented by what we shall refer to as propositions which she ascents to. We shall refer to these interchangeably as *premises*, *assumptions*, *basic beliefs*, *experiences*, or *evidence*. These sets of propositions also represent, in a way, the agent's *frame of mind*; we shall return to developing this perspective in §1.9. For now we will stick to the former readings.

To understand \Box , we think of the agent as producing derivations in a logical calculus on the basis of her evidence. So we alternately read $\Box \phi$ as *the agent believes ϕ* , *the agent can deduce ϕ* or *can compose an argument for ϕ* . We shall make further reference to this reading explicitly again in §.

Like the original formulation of epistemic logic in [Hin69], we assume that agents are *doxastically omniscient* - that is they believe all of the consequences of their beliefs. We shall prefer to think of this particular idealization as thinking about what an agent might conclude *eventually*.

By focusing on basic items of evidence as forming the foundation for belief, we consider this approach to be roughly in line with the *evidentialist* view on epistemology, which may be described as follows:

[E]videntialism is a supervenience thesis according to which facts about whether or not a person is justified in believing a proposition supervene on facts describing the evidence that the person has. [CF04], pg. 5

However, while our sympathies are with this perspective on epistemology, they differ foundationally - while evidentialism develops intuitions using analytical philosophy, our approach shall be founded in formal semantics like the one above.

As alluded to in §1.2, we shall read $\Diamond \phi$ as *the agent can conceive that ϕ* or *the agent can imagine ϕ being possible*. The former is the standard reading in epistemic logic (see, for instance, [MvdH95]). In general, we shall prefer to read \Diamond as in terms of *imagination*.

That is, if a proper foundation can be provided at all. Admittedly, the above formulation of truth immediately runs into a *paradox* - for instance, let

- $a := \emptyset$
- $A := \{\Box \perp\}$
- and $\mathfrak{M} := \{(a, A)\}$

Under this assignment, $\mathfrak{M}, (a, A) \models \Box \perp$ has no determinate truth value. So let $\mathcal{L}_0(\Phi)$ be the

propositional fragment of \mathcal{L}_K , with the following grammar:

$$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp$$

We shall restrict the truth value to $\mathfrak{M} \subseteq \Phi \times \wp\mathcal{L}_0(\Phi)$. This suffices to make every truth value of this logic determinate.

We may observe that the logic of these semantics is familiar:

Proposition 1.4.2. *Assuming that the set of proposition letters Φ is infinite*

$$\vdash_K \phi \text{ if and only if } \mathfrak{M}, (a, A) \models \phi \text{ for all finite } \mathfrak{M} \text{ for all } (a, A) \in \mathfrak{M}$$

Here K is basic modal logic.

Proof. Left to right is trivial, so we shall focus on right to left. Assume that $\not\models_K \phi$, then we know from completeness and the finite model property that there's some finite model $\mathbb{M} = \langle W, V, R \rangle$ and world $w \in W$ such that $\mathbb{M}, w \not\models \phi$ (see [BRV01], chapters 2 & 4 for details of these facts).

Now let $PropVars_\phi$ be the proposition letters that occur as subformulae of ϕ , and let $q : W \hookrightarrow \Phi \setminus PropVars_\phi$ be an injection. In other words q assigns *fresh letters* to worlds in the model³. Define $\theta : W^{\mathbb{M}} \rightarrow \wp\Phi \times \wp\mathcal{L}_0(\Phi)$ as follows

$$\theta(x) := (\{p \in \Phi \mid \mathbb{M}, w \models p\} \cup \{q_w\}, \{ \bigvee_{v \in R[w]} q_v \})$$

Now let $\Theta := \theta[W]$. An induction on the complexity of subformulae ψ of ϕ shows that $\mathbb{M}, w \models \psi \iff \Theta, \theta(w) \models \psi$ for all $w \in W^{\mathbb{M}}$. Since $\mathbb{M}, w \not\models \phi$ then we know that $\Theta, \theta(w) \not\models \phi$, which completes the proof. QED

Intuitively, the idea behind the above construction is to label all of the worlds with fresh letters, and then construct a special formula from these fresh letters for each world. The extension of each of these formulae is, in every case, exactly the worlds the agent could have accessed with the accessibility relation. A far more elaborate construction, based on the similar principles, shall be presented in §2.3.5.

Armed with this, we can see that these semantics are adequate for modeling agents according to our declared intentions. Recall the following definitions from basic logic and modal model theory⁴:

Definition 1.4.3. (1) *For any model \mathfrak{M} , define $Th(\mathfrak{M})$:*

$$Th(\mathfrak{M}) := \{\phi \in \mathcal{L}_K(\Phi) \mid \mathfrak{M}, (a, A) \models \phi \text{ for all } (a, A) \in \mathfrak{M}\}$$

*$Th(\mathfrak{M})$ is called **the theory of \mathfrak{M}** .*

(2) *$A \subseteq_\omega B$ means that A is a finite subset of B*

³ In this vein, we shall abbreviate $q(w)$ as q_w . Note that because $PropVars_\phi$ are *finite* and Φ is assumed to be *infinite*, such an inject always exists. This is a consequence of *The Axiom of Choice*.

⁴This notation consciously imitates the notation employed in [BRV01].

(3) Define $\Gamma \vdash_K \phi$ to mean:

$$\vdash_K \bigwedge \Delta \rightarrow \phi \text{ for some } \Delta \subseteq_\omega \Gamma$$

If $\Gamma \vdash \phi$, we say that ϕ is **derivable from** Γ .

The following theorem equates belief in at a world in a model with possession of a derivation:

Proposition 1.4.4. *For all $A \subseteq_\omega \mathcal{L}_0(\Phi)$, then $\mathfrak{M}, (a, A) \models \Box\phi$ if and only if $Th(\mathfrak{M}) \cup A \vdash_K \phi$.*

Proof. The proof of the above hinges on two basic facts. The first is the *deduction theorem* (provided that B is finite):

$$A \cup B \vdash_K \phi \iff A \vdash_K \bigwedge B \rightarrow \phi \quad (1.4.1)$$

The above follows from Definition 1.4.3 part (3), and is one of the standard results in modal logic.

The next observation is also rather basic:

$$Th(\mathfrak{M}) \vdash_K \phi \iff \phi \in Th(\mathfrak{M}) \quad (1.4.2)$$

The proof of this follows from the fact that if $\vdash_K \phi$ then $\phi \in Th(\mathfrak{M})$, and $Th(\mathfrak{M})$ can be observed to be closed under modus ponens.

So assume that $A \subseteq_\omega \mathcal{L}_0$. With the above key facts we have the following chain of reasoning:

$$\begin{aligned} Th(\mathfrak{M}) \cup A \vdash_K \phi &\iff Th(\mathfrak{M}) \vdash \bigwedge A \rightarrow \phi && \text{by (1.4.1)} \\ &\iff \bigwedge A \rightarrow \phi \in Th(\mathfrak{M}) && \text{by (1.4.2)} \\ &\iff \mathfrak{M}, (b, B) \models \bigwedge A \rightarrow \phi \text{ for all } (b, B) \in Th(\mathfrak{M}) && \text{by Def. 1.4.3 part (1)} \\ &\iff \mathfrak{M}, (a, A) \models \Box\phi \text{ for any } a \text{ where } (a, A) \in \mathfrak{M} && \text{by Def. 1.4.1} \end{aligned}$$

These equivalences suffice to prove the result. QED

A natural way to read $Th(\mathfrak{M})$ is the background knowledge the agent has about the universe she lives in. This approach presents an analysis of modal logic whereby an idealized agent is modeled as closed under deduction; this is the *doxastic omniscience* mentioned previously. Under this view, evidently the agent's beliefs correspond to those things for which she has proofs. This shall be the basis of our future investigations.

1.5 The Human Condition

To supplement to this basic framework, it shall be illustrated how further inspiration can be drawn from a philosophical perspective. This is in stark contrast to the received view in epistemic logic [Len78, pg. 34]:

The search for the correct analysis of knowledge, while certainly of extreme importance and interest to epistemology, seems not significantly to affect the object of epistemic logic, the question of the validity of certain epistemic-logical principles.

Quite to the contrary, we urge that epistemic logic should not turn its back on philosophy. Philosophy critically provides guidance for the intuitions behind how knowledge should be correctly modeled. It also provides a solid grounding in a proper treatment of knowledge. However, engaging with philosophy is evidently not the thrust of mainstream epistemic logic.

Most mainstream epistemic logic, the object of study is really the nature of information, not human knowledge. It applies equally well to robots, *homo economicus*, or thermometers as suggested in 1.2. Its inspiration is not really in what it is like to be a living person; it is more naturally based in artificial intelligence, information theory, automata theory, algebra, topology, and other abstract disciplines.

In contrast, we shall propose the following principle:

The Human Condition: *The analysis of knowledge should strive for a basis in human experience*

The above principle indeed underpins the justification principle provided in §1.3. This is because we feel that the belief in a proposition can be thought of human only if the agent has a reason associated with it. Otherwise, it seems that in the absence of reason, no account can be given for how the belief came about other than through instrumentation, which is the thermometer view.

By embracing the human condition, we shall now turn to the development of the logical perspective on knowledge extended here from its philosophical origins.

1.6 Soundness

So to give a shallow example of a basic application of a philosophical idea, it is natural to insist that if knowledge is based on beliefs generated via deduction from some set of premises, then those premises have to be *sound*. This can be done by introducing a new operator \odot with the following semantics:

$$\mathfrak{M}, (a, A) \models \odot \iff \mathfrak{M}, (a, A) \models A$$

Armed with these semantics, a first guess at what constitutes knowledge suggests it might be nothing more than possession of a belief based on a sound set of premises. So a first approximation of knowledge might be equated with the formula:

$$\odot \wedge \Box \phi.$$

But is this anything like an adequate analysis of knowledge?

No. To illustrate why, let us consider a thought experiment. Imagine that Charlotte suspects, correctly, that if John has tried to murder on Alex, then Alex has survived. She further learns, correctly, that John has indeed tried to murder Alex. But later, she “learns” some erroneous information asserting Vietnam is south of Malaysia. If we codify all of this as a set C , and let the real world be denoted c and the universe \mathfrak{M} , evidently we have $\mathfrak{M}, (c, C) \not\models \odot$, so this previous definition of knowledge fails. But should it? This is doubtful; Charlotte’s knowledge about John’s unspeakable betrayal of Alex is correct, as well as her inference that Alex is tough as nails. Just

because she has been deluded regarding irrelevant facts about geography shouldn't have any bearing on her knowledge about Alex.

1.7 Descartes

In reflection on the previous section, it should be remarked that philosophers have historically been concerned with defeasible experiential data, going back at least as early as Plato's *The Republic VII* [Pla98]. In answer to the problem faced by the above analysis of knowledge, guidance can be found in Descartes' *Meditations* [VMTD05]. In *Meditations I*, Descartes suggests that he might be in an enlightenment era version of *The Matrix* created by an all powerful demon. In *Meditations II*, he famously suggests how one might escape this trap:

The Meditation of yesterday has filled my mind with so many doubts, that it is no longer in my power to forget them. Nor do I see, meanwhile, any principle on which they can be resolved; and, just as if I had fallen all of a sudden into very deep water, I am so greatly disconcerted as to be unable either to plant my feet firmly on the bottom or sustain myself by swimming on the surface. I will, nevertheless, make an effort, and try anew the same path on which I had entered yesterday, that is, proceed by casting aside all that admits of the slightest doubt, not less than if I had discovered it to be absolutely false; and I will continue always in this track until I shall find something that is certain, or at least, if I can do nothing more, until I shall know with certainty that there is nothing certain.[VMTD05, *Meditations II*]

This tactic proposes a natural solution to the problem the previous thought experiment: *Charlotte can know that Alex survives if she argues **only** from her experience involving Alex and John.* If like Descartes she can forget some of what she has come to believe that's a little suspicious, she might be able to compose an argument with a sound basis that Alex is alive. Taking Descartes as inspiration, we might think of a novel semantic operation:

$$\mathfrak{M}, (a, A) \models \exists \phi \iff \text{for all } (b, B) \in \mathfrak{M} \text{ such that } a = b \text{ and } B \subseteq A \text{ then } \mathfrak{M}, (b, B) \models \phi$$

This mechanism lets Charlotte access subsets of her beliefs, which would then form the basis for various arguments she might compose. Provided that $(c, C') \in \mathfrak{M}$, where C' is the same as C but doesn't mention erroneous beliefs about geographical data, it might serve as a basis for Charlotte's knowledge that Alex survives. This suggests that the following equation might reasonably express a more adequate notion of knowledge:

$$\Diamond(\bigcirc \wedge \Box \phi)$$

1.8 Contradictions

There's hidden virtue in the previous analysis. To see what it is, inspiration can be found in the 19th century philosopher Ralph Waldo Emerson, who writes in his essay *Self-Reliance* [Eme08]:

Why drag about this corpse of your memory, lest you contradict somewhat you have stated in this or that public place? Suppose you should contradict yourself; what then?

It seems to be a rule of wisdom never to rely on your memory alone, scarcely even in acts of pure memory, but to bring the past for judgment into the thousand-eyed present, and live ever in a new day. . . .

A foolish consistency is the hobgoblin of little minds, adored by little statesmen and philosophers and divines. With consistency a great soul has simply nothing to do. He may as well concern himself with his shadow on the wall. Speak what you think now in hard words and to-morrow speak what to-morrow thinks in hard words again, though it contradict every thing you said to-day. – ‘Ah, so you shall be sure to be misunderstood.’ – Is it so bad then to be misunderstood?

A healthy lack of consistency is just part of what makes up the day to day life of any living, sane person. Isn’t error-prone reasoning a hallmark of human thought? And if a love sick epistemic agent \exists is getting mixed signals from another epistemic agent \forall , why can’t she draw inconsistent conclusions about \forall ’s feelings on the one hand, but still have basic knowledge that $734 \times 12 = 8808$ and other such irrelevant facts? There is no compelling reason why not. Under these considerations, the following is compelling:

Emerson’s Principle: *One can be inconsistent and still have knowledge*

We may observe how the framework so far developed accommodates this. We may draw further inspiration by a friend and contemporary of Emerson’s, the poet Walt Whitman. In *Leaves of Grass* [Whi08], he writes:

Do I contradict myself?
Very well then I contradict myself,
(I am large, I contain multitudes.)

So consider the model \mathfrak{M} in Fig. 3; this is intended to be a toy model of how one might interpret Walt Whitman in the above stanza. This figure should be read as follows:

- if one point (a, A) is above another point (b, B) and connected by a densely dotted line \cdots , this means that $a = b$ and $B \subset A$.
- if one point (a, A) is connected to another point (b, B) by a line with an arrow \longrightarrow , this means that $\mathfrak{M}, (b, B) \models A$

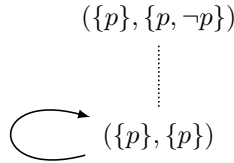


Figure 3: Inconsistent, yet still has knowledge

Observe that $\mathfrak{M}, (\{p\}, \{p, \neg p\}) \models \Box \perp$; it's obvious that in this state Walt is being inconsistent since he clearly believes contradictory things. Simultaneously, we have that $\mathfrak{M}, (\{p\}, \{p, \neg p\}) \models \Diamond(\Diamond \wedge \Box p)$; so we figure that Walt has a sound argument that p . Walt might be inconsistent, but it would appear that *at least one* of his arguments makes sense. And this is naturally because Walt contains a multiplicity of inner selves, just like he says, which the Ξ modality gives access to.

1.9 Irrationality

Embracing contradiction runs contrary to the received view on epistemic logic. For instance, [HS06] write:

Epistemic logic does carry epistemological significance but in an inevitably idealized sort of way: One restricts attention to a class of rational agents where rationality is defined by certain postulates. Thus, agents have to satisfy at least some minimal conditions to simply qualify as rational. This is by and large what Lemmon originally suggests [LH59].

Furthermore, it is conventional to think that rational agents do not hold contradictions.⁵ For instance, in [KL86], $\neg \Box \perp$ is taken as an axiom (it is A9 in their numbering).

This is similar to the thermometer concept of knowledge we provided in §1.3, since like the thermometer view, it is incompatible with a human perspective. Hence we shall extend the following:

Irrationality Principle: *Since humans are not rational, views on epistemic logic that postulate this should be rejected*

I should mention that while this perspective is not typically embraced in epistemic logic⁶, it finds sympathy in other logical traditions, namely in *relevance logic* and *paraconsistent logic*, as already noted [see GG02, chapters 1 & 4].

Apart from inconsistency, we do not really accommodate very much irrationality; frameworks like [Ran82] and [Lev84] employing *impossible world* semantics are far more accommodating to irrationality than the semantics we are investigating. However, these frameworks do not provide compelling explanation for the mechanism of irrationality, contrary to the perspective presented here. Regardless, allowing for an agents beliefs to be generated from inconsistent premises is already orthogonal to the assumption that agents are rational.

⁵It should be remarked that [Pri06] explicitly rejects this perspective on rationality. Priest points out that in times of scientific revolution, rational people naturally hold contradictory views. He suggests that a paraconsistent logic framework could account for a rational agent holding contradictory beliefs. While our hearts are indeed sympathetic to Priest's perspective, we are nonetheless confident that this does not represent the received view which we are attacking.

⁶Noted exceptions to this are [Ran82] and [Lev84].

1.10 Quine

To recap, so far we have suggested adding a novel modality \boxplus which corresponds to taking subsets of an agent’s set of beliefs. In the context of conventional modal logic, this means a shift in perspective - instead of thinking of each world as a situation where the agent can imagine other situations, now each world corresponds to a network of beliefs ordered by inclusion. These networks of beliefs form a poset, or partially ordered set. Thus the choice to visually represent them as *Hasse diagrams*, as we have done in Fig. 3, follows the standard practice in lattice theory.

Furthermore, consider the following phenomenon - as higher nodes in a belief network are considered, the agent is employing more premises for the arguments they are composing, and using less pure logic to come to conclusions. This suggests that as we consider levels higher and higher in the poset of an agent’s beliefs, this corresponds to embracing an agent’s experience and interpretation of their sensory data. Arguments that rest on more premises are *prima facie* more fallible than arguments that rely on fewer assumptions.

A similar perspective has been presented before, however in a different setting, in *Two Dogmas of Empiricism*:

Certain statements, though about physical objects and not sense experience, seem peculiarly germane to sense experience – and in a selective way: some statements to some experiences, others to others. Such statements, especially germane to particular experiences, I picture as near the periphery. But in this relation of “germaneness” I envisage nothing more than a loose association reflecting the relative likelihood, in practice, of our choosing one statement rather than another for revision in the event of recalcitrant experience. For example, we can imagine recalcitrant experiences to which we would surely be inclined to accommodate our system by re-evaluating just the statement that there are brick houses on Elm Street, together with related statements on the same topic. We can imagine other recalcitrant experiences to which we would be inclined to accommodate our system by re-evaluating just the statement that there are no centaurs, along with kindred statements. A recalcitrant experience can, I have already urged, be accommodated by any of various alternative re-evaluations in various alternative quarters of the total system; but, in the cases which we are now imagining, our natural tendency to disturb the total system as little as possible would lead us to focus our revisions upon these specific statements concerning brick houses or centaurs. These statements are felt, therefore, to have a sharper empirical reference than highly theoretical statements of physics or logic or ontology. *The latter statements may be thought of as relatively centrally located within the total network, meaning merely that little preferential connection with any particular sense data obtrudes itself.* [Qui51]

The emphasis on the last sentence is our addition. The above paragraph importantly anticipates ideas in belief revision theory (such as in [AGM85] and subsequent studies), as well as recent trends in probabilistic dynamic epistemic logic [such as in vB03, vBGK09, BS08, Koo03, etc.]. However, in the framework has been developing so far, what Quine refers to as the “periphery” of his web of belief corresponds to a higher node in a belief poset, while what Quine refers to as the “center” reflects something like a lower node. This is visually depicted in Fig. 4. Beliefs that are members

of lower nodes, and the ideas that follow from them, can be thought of as belonging to the agent's world-view.

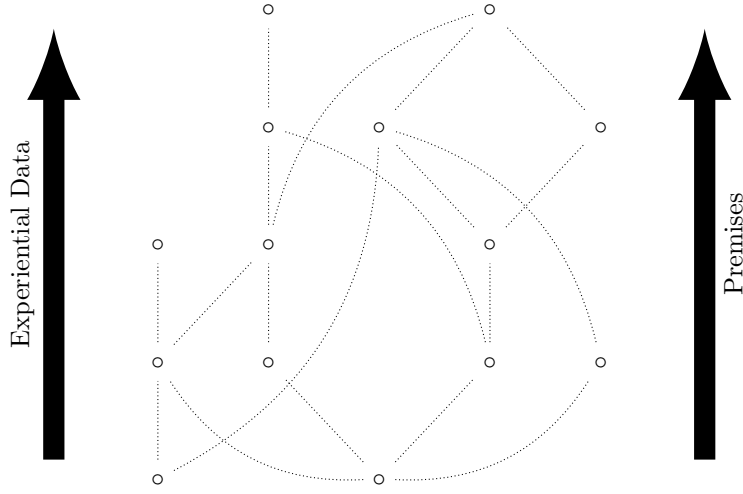


Figure 4: A network of beliefs

The above observation informs a corresponding perspective on epistemology. If an agent's world view largely rested legends about the Norse gods, we should be reluctant to say she knows various facts about nature, such as why lightning strikes. This is because all of her explanations would inevitably be based upon myths in one way or another, which would all occupy lower nodes in her belief network. This dictates that *sanity* plays a role in how much knowledge an agent can have - it is permissible to grant that an inconsistent agent has knowledge provided that the inconsistency follows only shallowly from her experiential data, and it is something she would readily give up. However, if a contradiction is intrinsic to the agent's psychology, and thus follows from a lower node in her belief poset, this suggests she does not really have knowledge. So while we should believe that irrational agents can possess knowledge, as we have argued in §1.9, we should rightfully not contend that they *always* possess knowledge. Moreover, the sort of irrationality that we are considering here need not be superficial - both mundane as well as deeply demented characters can be modeled.

Note that the above essentially presents a subjective interpretation of Quine's web of belief, which might be contentious. On the other hand, we should feel both the quote from Quine and the quote from Whitman in §1.8 suggest the following principle without too much controversy:

Quine/Whitman Principle: *Epistemic agents are compound entities, which invite compositional analysis.*

The above presents a final philosophical principle that shall be extended extend. Apart from this, from the previous discussion may extract an additional thing: Figure 4 naturally suggests that we might think of *going up* in a belief net, in a manner similar to how \boxplus allows one to *go down* as suggested in §1.7. This suggests the introduction of a new operator \boxminus . The semantics for \boxminus are

given as follows:

$$\mathfrak{M}, (a, A) \models \boxplus \phi \iff \text{for all } (b, B) \in \mathfrak{M} \text{ if } b = a \text{ and } A \subseteq B \text{ then } \mathfrak{M}, (a, A) \models \phi$$

Just as \boxminus corresponds to the agent casting assumptions into doubt, or disregarding their premises, \boxplus corresponds to the agent embracing their experience, suspending disbelief and accepting her intuitions and senses.

This concludes the presentation of novelties we propose for the practice of modelling knowledge.

1.11 Closing Remarks

The various principles extended in the previous sections are not independent - some of them are more basic than others. Their relationship is summarized in Fig. 5 - here the lower a principle is depicted, the more basic. Dotted lines indicate the philosophical justification for the higher principle supervenes on the justification of the lower principle. In addition, in further development

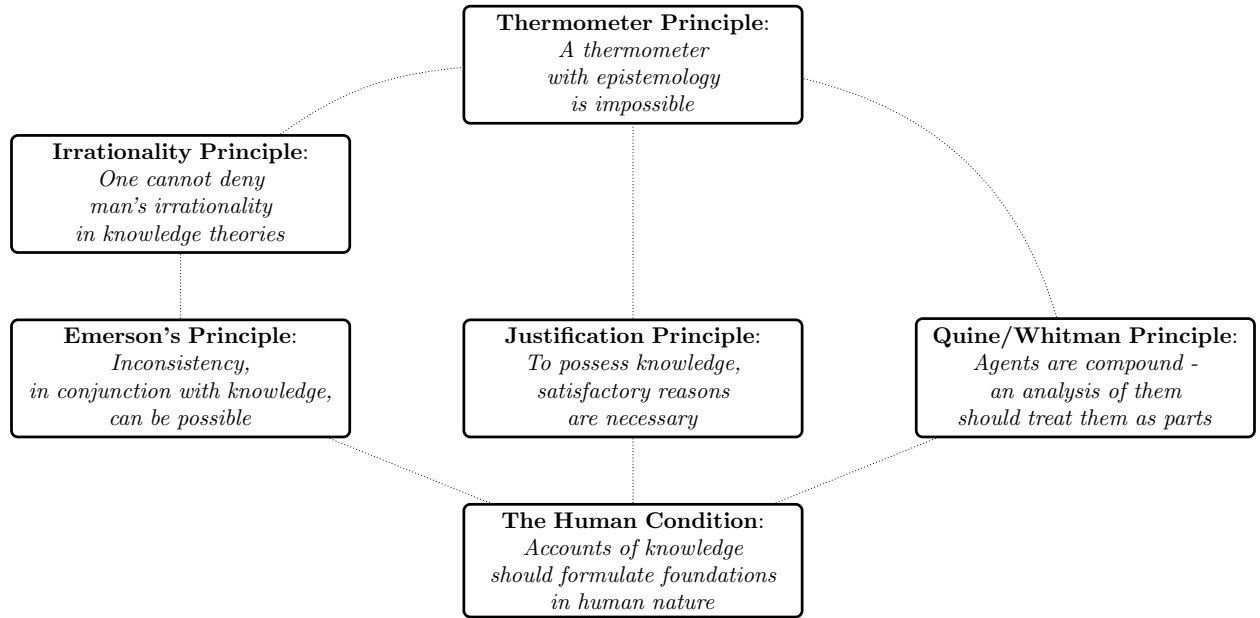


Figure 5: A visualization of the relationship of the principles presented here

of the framework sketched in §1.3, we shall want the following criteria, based on the ideas given in relevant sections:

- §1.3
- Agents shall be modelled with proofs for the things they believe.
 - To avoid paradoxes, correct foundations must be provided. Ideally, we would like our semantics to correspond to a provably terminating computation, granting certain non-deterministic operations such as a *choice operator* ε , as described in [Hil22].

For a set of beliefs A :

§1.6 It should be expressible whether everything in A is sound

§1.7 Certain subsets $B \subseteq A$ should be accessible

§1.10 Certain extensions $B \supseteq A$ could also be accessed

In line with evidentialist epistemology, as mentioned in §1.4, we have decided to call the logic presented here *Evidentialist Logic*, or **EvIL** for brevity.

2 Introduction to EvIL

From with the philosophical intuitions and scaffolding provided from §1, we shall present a precise account of the previously developed ideas. This shall be done in three movements:

§2.1 In the first section we shall provide the basic grammar and semantics for **EvIL** with a single agent; the presentation in this section will remain primarily philosophical and light.

§2.2 In the second section we develop several topics in the pure theory of **EvIL** which are considered a bit beyond the bare essentials.

§2.3 In this section, logics for **EvIL** are presented, along with completeness and decidability.

2.1 Elementary EvIL

2.1.1 Grammar & Semantics

In this section we turn to developing the formal semantics for **EvIL** with a single agent. We shall imagine the object of study in **EvIL** is an agent, which we shall call the **EvIL** agent. In §2.2.2, the semantic framework offered here is extended to incorporate multiple agents. In Appendix B, yet another framework is offered employing gamelike semantics, which avoids the grammar restriction suggested in §1.4.

The grammar restriction imposed on **EvIL** was introduced to avoid paradoxes. That being the case, we shall discard the previous definition of (\models) that was suggested in §1.4, in favor of demonstrably well-defined semantics. This shall be achieved in two steps.

Definition 2.1.1. *Let $\mathcal{L}_0(\Phi)$ be the language of classical propositional logic, defined by the following Backus-Naur form grammar:*

$$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp$$

Models for classical propositional logic can be thought of as sets $S \subseteq \Phi$; thus the truth predicate $(\models) : \wp\Phi \times \mathcal{L}_0(\Phi) \rightarrow \text{bool}$ for classical propositional logic can be given recursively as follows:

Definition 2.1.2. Define (\models) such that

$$\begin{aligned} S \models p &\iff p \in S \\ S \models \phi \rightarrow \psi &\iff S \models \phi \text{ implies } S \models \psi \\ S \models \perp &\iff \text{False} \end{aligned}$$

Further, observe that the language \mathcal{L}_0 is extended by EViL

Definition 2.1.3. Define $\mathcal{L}(\Phi)$ by the following Backus-Naur grammar:

$$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp \mid \Box \phi \mid \Box \phi \mid \Box \phi \mid \bigcirc$$

Unlike traditional modal logic, EViL employs concrete models rather than Kripke structures. EViL models are sets $\mathfrak{M} \subseteq \wp\Phi \times \wp\mathcal{L}_0(\Phi)$. Like classical propositional logic, semantics for EViL are given recursively by a predicate (\models) which:

- Takes as input:
 - An EViL model
 - A pair (a, A) where
 - ◊ $a \subseteq \Phi$ is a set of proposition letters
 - ◊ $A \subseteq \mathcal{L}_0(\Phi)$ is a set of propositional formulae.
 - A formula in the language $\mathcal{L}(\Phi)$
- Gives as output: a truth value in **bool**

More concisely, this may be written as

$$(\models) : (\wp(\wp\Phi \times \wp\mathcal{L}_0(\Phi))) \times (\wp\Phi \times \wp\mathcal{L}_0(\Phi)) \times \mathcal{L}(\Phi) \rightarrow \text{bool}.$$

Definition 2.1.4. Define (\models) recursively such that:

$$\begin{aligned} \mathfrak{M}, (a, A) \models p &\iff p \in a \\ \mathfrak{M}, (a, A) \models \phi \rightarrow \psi &\iff \mathfrak{M}, (a, A) \models \phi \text{ implies } \mathfrak{M}, (a, A) \models \psi \\ \mathfrak{M}, (a, A) \models \perp &\iff \text{False} \\ \mathfrak{M}, (a, A) \models \Box \phi &\iff \forall (b, B) \in \mathfrak{M}. (\forall \psi \in A. b \models \psi) \text{ implies } \mathfrak{M}, (b, B) \models \phi \\ \mathfrak{M}, (a, A) \models \Box \phi &\iff \forall (b, B) \in \mathfrak{M}. a = b \text{ and } B \subseteq A \text{ implies } \mathfrak{M}, (b, B) \models \phi \\ \mathfrak{M}, (a, A) \models \Box \phi &\iff \forall (b, B) \in \mathfrak{M}. a = b \text{ and } B \supseteq A \text{ implies } \mathfrak{M}, (b, B) \models \phi \\ \mathfrak{M}, (a, A) \models \bigcirc &\iff \forall \psi \in A. a \models \psi \end{aligned}$$

Remark 2.1.5. We will write $\mathfrak{M} \models \phi$ to mean $\mathfrak{M}, (a, A) \models \phi$ for all $(a, A) \in \mathfrak{M}$. Further, we will write $\models \phi$ to mean $\mathfrak{M} \models \phi$ for all \mathfrak{M} .

These semantics are well defined, since apart from relying on the semantics for propositional logic they may be observed to be compositional. Moreover, the following relationship can be observed:

Lemma 2.1.6 (Truthiness). *Let $\phi \in \mathcal{L}_0(\Phi)$. Then:*

$$a \models \phi \iff \mathfrak{M}, (a, A) \models \phi$$

for any \mathfrak{M} and A .

Proof. This may be seen immediately by induction on ϕ . QED

With this, we have the following, mirroring Prop. 1.4.4:

Definition 2.1.7. *Define the following:*

$$Th(\mathfrak{M}) := \{\phi \in \mathcal{L}(\Phi) \mid \mathfrak{M} \models \phi\}$$

Theorem 2.1.8 (Theorem Theorem). *If A is finite, then $\mathfrak{M}, (a, A) \models \Box\phi$ if and only if $Th(\mathfrak{M}) \cup A \vdash_{\text{EvIL}} \phi$.*

Proof. The proof proceeds the exactly as the proof of Proposition 1.4.4 from §1.4. QED

I shall present \vdash_{EvIL} , the logical consequence turnstile for EvIL, in §2.3.1.

I chose the name “Theorem Theorem” because it means that for every belief the EvIL agent has, it is a theorem she has derived from her premises. Theorem 2.1.8 establishes one of the central desiderata outlined in §1.11 is achieved by EvIL. With this result the foundation is set for the the central intuition driving EvIL - that beliefs are the consequences of logical deductions. It is a peculiarity of EvIL that these deductions are carried on in EvIL itself. This was achieved, primarily, by flirting heavily with paradox, as was illustrated in §1.4. As a consequence, we have tried to design EvIL to eat its own tail. It establishes that the EvIL agent is herself also a modeler just like us, using the same logic we are using to think about her herself, to think about the state space she lives in.

This sort of self referential circularity is a celebrated theme in mathematics. It is similar, in a way, to the old alchemical conception of mathematics, exemplified by the following quote due to Sir Thomas Browne:

All things began in order, so shall they end, and so shall they begin again; according to the ordainer of order and mystical Mathematicks of the City of Heaven. [Bro36], chapter 5

Another, related notion of self reference was championed by Douglas Hofstadter in his book *Gödel, Escher, Bach: An Eternal Golden Braid*, in what he calls “Strange Loops”:

The flexibility of intelligence comes from the enormous number of different rules, and levels of rules. The reason that so many rules on so many different levels must exist is that in life, a creature is faced with millions of situations of completely different types. In some situations, there are stereotyped responses which require ”just plain”

rules. Some situations are mixtures of stereotyped situations-thus they require rules for deciding which of the “just plain” rules to apply. Some situations cannot be classified-thus there must exist rules for inventing new rules ...and on and on. Without doubt, Strange Loops involving rules that change themselves, directly or indirectly, are at the core of intelligence. Sometimes the complexity of our minds seems so overwhelming that one feels that there can be no solution to the problem of understanding intelligence-that it is wrong to think that rules of any sort govern a creature’s behavior, even if one takes “rule” in the multilevel sense described above. [Hof79, pg. 24]

In the setting of EVIL, the strange loop is as follows: the rules of the logic affect, indirectly, the truths of the semantics, which in turn affect the rules of the logic. It is perhaps dangerous to engage in this sort of modeling, since it invites paradox. However, the central goal of epistemic logic is to provide a logic which models the modelers, so that we might perhaps learn something about ourselves by studying the subject. The Theorem Theorem goes some distance in EVIL towards accomplishing this goal.

It cannot be stressed enough, Theorem 2.1.8 is central to the conceptual backing behind EVIL. It provides the conceptual backbone of the perspective of epistemology this essay is intended to investigate.

2.1.2 Intuitions

In this section, we shall illustrate how we intuitively read the operators in EVIL, and provide a number of validities.

As per the traditional doxastic reading of $\Box\phi$, we read this as asserting “The EVIL agent believes ϕ ”. Because of Theorem 2.1.8, the Theorem Theorem, we shall freely conflate this with the assertion “The EVIL agent has an argument for ϕ ,” which we take to be a kind of proof.

The intuition for how to read $\Diamond\phi$ was first mentioned in §1.7 with respect to Descartes’ Meditation II – it means “If the EVIL agent were to set aside some of her beliefs, or cast some of her beliefs into doubt, then ϕ would hold.” Dually, we can read $\Box\phi$ as saying something like “For all the ways that the EVIL agent might use her imagination, ϕ holds.” One should recognize that these interpretations might seem inconsistent. These are not really an issue regard casting beliefs into doubt and embracing one’s imagination as part of the same coin. For, naturally, when one doubts more things, then for a fleeting moment their dreams take flight as the inconceivable turns around into the conceivable, if only for a little while. To give an example, if Marta sets aside for a moment her belief that

the law of gravity is an exceptionless regularity of the universe, (g)

then it seems natural that she might imagine that

a propulsion device exploiting some exception to gravitation might be constructable. (p)

In the symbology of EVIL formulae, she would code this intuition as

$$\Box(\Box\neg g \rightarrow \Diamond p). \tag{2.1.1}$$

To give another example, if Marta pretends that it is not the case that:

the canals of Amsterdam are filthy (f)

She might be able to imagine a scenario where

she may swim comfortably in the Amstel river (r)

But not really. Marta really cannot really swim at ease in the Amstel, not just because it has tons of garbage, but also because

she does not own a bathing suit, (b)

Frankly, Marta is not so bold that she could go skinny dipping in Amstel without that being awkward for her. In the language of EVIL, this thought experiment would be expressed as follows:

$$\neg \boxplus (\Box \neg f \rightarrow \Diamond r) \quad (2.1.2)$$

This is because Marta can cast into doubt the assumption of the filthiness of the canals of Amsterdam, while still retaining her belief that she does not have a bathingsuit, so swimming in Amstel would still be awkward for me. In symbols, she would write express this other sentiment as something of a refinement on (2.1.2), which is expressed as follows:

$$\Diamond (\Diamond \neg f \wedge \Box b \wedge \neg \Diamond r) \quad (2.1.3)$$

Further, the intuition for how to read $\Diamond \phi$ is “If the EVIL agent were to remember something, then ϕ would hold.” For instance, imagine a scenario where Marta wakes up and searches herself for her bike keys. To her horror, the keys are not there – and Marta immediately assumes that she might have left her keys in the lock on her bike, and figures there is a fair likelihood that

the bike has been stolen because the keys were left in the lock. (s)

But once she recalls that

she lent her bike to a friend, (l)

her fear subsides. Prior to remembering, while Marta thought it might be possible that her bike was stolen due to her own negligence, if she remembered what she had done then she no longer would have entertained this possibility. This observation is expressed as:

$$\Diamond s \wedge \boxplus (\Box l \rightarrow \Box \neg s) \quad (2.1.4)$$

We consider \boxminus and \boxplus to be inverse modalities of each other, in exactly the same way that *past* and *future* are inverse modalities in temporal logic. This is perhaps a little unusual; it is arguably more natural to think of *forgetting* as the inverse modality of remembering, and there does not appear to be an natural inverse operation corresponding to casting into doubt. Following the idea of the *web of belief* due to Quine, as presented in §1.10, we would extend a position asserting that remembering factive data is the same as embracing as much of one’s evidence as possible.

In terms of the semantics outlined, \boxminus corresponds to a subsetset relation while \boxplus corresponds to a superset relation. Because of this, we shall sometimes read $\boxminus \phi$ closer to the formal semantics, as saying something like “for all subsets of the agent’s basic beliefs or premises, ϕ holds” and dually

for $\boxplus\phi$. This is admittedly even less natural than the reading of remembering as the opposite of casting into doubt. So be it; we shall have to be comfortable with EViL agents being at best twisted cartoon versions of actual people. Is this so unnatural? Logic, in general, only affords at best cartoon versions of whatever we are trying to model with it. Consider, for instance, Peano arithmetic, which is in general not powerful enough to prove a basic number theoretic fact like Goodstein’s theorem [KP82]. Or consider the stronger system of second order arithmetic, which is in turn not strong enough to prove semantics stipulated in §2.1.1, EViL agents apparently have sets for brains, which makes an EViL agent a strange effigy for a person indeed – with the possible exception of set theorists, whose brains are typically constructed entirely of sets or urelements.

Furthermore, it is under the set theoretical reading that \odot makes the most sense. We should read it as asserting something like “the basis for the EViL agent’s beliefs is sound” or “the EViL agent’s arguments only use true premises.” It further means that the actual state of affairs is compatible with what the agent believes - reality has not been ruled out by something that the agent is taking as evidence. Moreover, sound premises intuitively exhibit the following property - any subset of them is also sound, since soundness isn’t a phenomenon that is subject to synchronicity or other failures of compositionality. A set of premises is sound if and only if all of its subsets are also sound.

2.1.3 Validities

The previous philosophical readings of EViL immediately suggest certain validities will hold in the semantics. For instance, the assertion “A set of premises is sound if and only if all of its subsets are sound.” would be expressed as

$$\models \odot \leftrightarrow \boxplus \odot \quad (2.1.5)$$

Indeed, this is a validity of EViL. Schematically, it may be tempting to think that maybe the same is true for \boxplus . However, we have that:

$$\not\models \odot \rightarrow \boxplus \odot \quad (2.1.6)$$

Nor does this make much sense. It asserts “If the agent’s basic beliefs are sound, then all extensions of her basic beliefs are sound too.” Soundness is a fragile thing – it is rather easy to think of things to add to a sound set of basic beliefs which break soundness, such as “All logicians are centaurs” and other such demonstrably false nonsense.

Related to (2.1.5), there is another related validity associated with \odot ; namely that if the EViL agent’s assumptions are sound, then anything she concludes from them is true (employing the reading which naturally arises from Theorem 2.1.8). This is expressed as

$$\models \odot \rightarrow \Box \phi \rightarrow \phi \quad (2.1.7)$$

The formula (2.1.5) expresses that the soundness of one’s premises is something *persistent* as the EViL agent carries on casting doubt on assumptions and discarding them. Another thing that is persistent this way is the EViL agent’s imagination:

$$\models \Diamond \phi \rightarrow \boxplus \Diamond \phi \quad (2.1.8)$$

One may read (2.1.8) as saying something like “If the EViL agent can conceive/imagine something, then no matter what things she casts into doubt, she can still imagine it.” One can also express

something like the dual of this, namely

$$\models \Box\phi \rightarrow \boxplus \Box\phi \quad (2.1.9)$$

We shall read the above as asserting “If the agent *can compose an argument* then she will still be able to compose that argument if she remembers more information and experiences she’s had in the world.” This should not be surprising – this is yet another expression of the Theorem 2.1.8, the Theorem Theorem, and the fact that the proof theory of EVIL is monotonic. In general, many of the assertions here exhibit interplay between \boxplus and \Box , and dually \boxminus and \Diamond – further investigation of these relationships is taken up in §2.2.1.

For better or for worse, EVIL semantics make true the following assertion: if something is achievable by repeatedly casting assumptions into doubt, then it’s achievable by casting assumptions into doubt only once:

$$\models \Diamond^+\phi \rightarrow \Diamond\phi \quad (2.1.10)$$

Here $^+$ is taken from the syntax for *regular expressions* commonly used in computer science and UNIX programming to mean “one or more” [Fri06]. Similarly, we have assumed that discarding no assumptions is, in a way, vacuously casting assumptions into doubt. In light of this EVIL also makes true the following:

$$\models \phi \rightarrow \Diamond\phi \quad (2.1.11)$$

Furthermore, it is worth mentioning some harder to understand validities of this system. The first one is that when the agent believes something, they believe it regardless of the process of doubting or embracing their beliefs:

$$\models \Box\phi \rightarrow \Box\boxminus\phi \quad (2.1.12)$$

$$\models \Box\phi \rightarrow \Box\boxplus\phi \quad (2.1.13)$$

We can observe that this generalizes to multiple agents, as specified in §2.2.2.

Another more challenging validity is the fact that if some proposition ϕ holds, then for any restriction of the EVIL agent’s beliefs (or dually, any extension), if those beliefs are sound, then ϕ must be conceivable (i.e., $\Diamond\phi$ holds). This is expressed as the following two validities:

$$\models \phi \rightarrow \boxminus(\Diamond \rightarrow \Diamond\phi) \quad (2.1.14)$$

$$\models \phi \rightarrow \boxplus(\Diamond \rightarrow \Diamond\phi) \quad (2.1.15)$$

Finally, another a peculiarity of EVIL is that not all of its validities are *schematic*. For instance, there is a kind of *Cartesian dualism* present in the semantics, where the EVIL agent’s deliberation on her evidence does not bear on brute matters of fact. For a world pair (a, A) , A and a are basically separate - an EVIL agent’s mind and the world they live are composed of different substance. This gives rise to the following four validities:

$$\models p \rightarrow \boxminus p \quad (2.1.16)$$

$$\models p \rightarrow \boxplus p \quad (2.1.17)$$

$$\models \neg p \rightarrow \boxminus \neg p \quad (2.1.18)$$

$$\models \neg p \rightarrow \boxplus \neg p \quad (2.1.19)$$

Hence, EviL is not a *normal* logic. This should admittedly be considered a wart on the semantics, since it appears that it rules out the conventional algebraic duality most modal logics exhibit (see [BRV01], chapter 5).

On the other hand, it is by the same assumption of Cartesian dualism that underlies the non-normality that (2.1.5) as is a natural consequence. By accepting non-normality, and the grammar restriction we have imposed on *basic beliefs* to avoid paradoxes, it follows as a consequence that a belief set is sound if and only if all of its subsets are sound. Hence non-normality for EviL has two aspects – it compromises the algebraic elegance of the semantics, while simultaneously giving rise to a philosophically appealing principle.

In the next section, we turn to a more systematic study of the validities of EviL . We shall see that this gives rise to an *elimination theorem*.

2.2 Basic EviL

2.2.1 Elimination

In section §2.1.3, we saw some of the structural validities of EviL from a philosophical perspective. That being the case, the manner of presentation followed intuition, which did not follow an orthodox organization. In this section, we shall look at the validities of EviL in a more systematic presentation. In doing so, we investigate an elimination theorem, which sits at the heart of EviL .

To start, the following lemma summarizes the structural validities will be studied in the subsequent discussion:

Lemma 2.2.1. *The following validities hold for all EviL models:*

$$\begin{array}{ll}
\models \Box p \leftrightarrow p & \models \Box p \leftrightarrow p \\
\models \Box \neg p \leftrightarrow \neg p & \models \Box \neg p \leftrightarrow \neg p \\
\models \Box \Diamond \phi \leftrightarrow \Diamond \phi & \models \Box \Box \phi \leftrightarrow \Box \phi \\
\models \Box \Diamond \phi \leftrightarrow \Diamond \phi & \models \Box \Diamond \phi \leftrightarrow \Diamond \phi \\
\models \Box \Box \phi \leftrightarrow \Box \phi & \models \Box \Box \phi \leftrightarrow \Box \phi \\
\models \Box \Diamond \phi \leftrightarrow \Diamond \phi & \models \Box \Diamond \phi \leftrightarrow \Diamond \phi \\
\models \Box \Diamond \phi \leftrightarrow \Diamond \phi & \models \Box \Diamond \phi \leftrightarrow \Diamond \phi \\
\models \Box \Diamond \phi \leftrightarrow \Diamond \phi & \models \Box \Diamond \phi \leftrightarrow \Diamond \phi
\end{array}$$

These validities suggest a definite interplay between the modalities of EviL ; they are highly suggestive of a general elimination theorem. To see what arises from Lemma 2.2.1, first observe that EviL makes true the usual substitution rule:

Lemma 2.2.2. *If $\models \phi \leftrightarrow \psi$ is a validity, then $\models \chi \leftrightarrow \chi[\phi/\psi]$ is a validity for any $\chi \in \mathcal{L}(\Phi)$.*

Next, consider two sublanguages of the main language of EviL :

Definition 2.2.3. *Define the following fragments:⁷*

⁷I was inspired to look at the fragment $\mathcal{L}_A(\Phi)$ by thinking about the continuous fragment of μPML [Fon08].

$\mathcal{L}_A(\Phi)$:

$$\phi ::= p \mid \neg p \mid \top \mid \perp \mid \circlearrowleft \mid \phi \wedge \psi \mid \phi \vee \psi \mid \Diamond \phi \mid \Box \phi \mid \Diamond \Box \phi$$

$\mathcal{L}_B(\Phi)$:

$$\phi ::= \neg p \mid p \mid \perp \mid \top \mid \neg \circlearrowleft \mid \phi \vee \psi \mid \phi \wedge \psi \mid \Box \phi \mid \Diamond \Box \phi \mid \Box \Diamond \phi$$

Definition 2.2.4. Define two dualizing operations $(\cdot)^A : \mathcal{L}_B(\Phi) \rightarrow \mathcal{L}_A(\Phi)$ and $(\cdot)^B : \mathcal{L}_A(\Phi) \rightarrow \mathcal{L}_B(\Phi)$, using recursion, such that:

$$\begin{array}{ll} \neg p^A := p & p^B := \neg p \\ p^A := \neg p & \neg p^B := p \\ \perp^A := \top & \top^B := \perp \\ \top^A := \perp & \perp^B := \top \\ \neg \circlearrowleft^A := \circlearrowleft & \circlearrowleft^B := \neg \circlearrowleft \\ (\phi \vee \psi)^A := \phi^A \wedge \psi^A & (\phi \wedge \psi)^B := \phi^B \vee \psi^B \\ (\phi \wedge \psi)^A := \phi^A \vee \psi^A & (\phi \vee \psi)^B := \phi^B \wedge \psi^B \\ (\Box \psi)^A := \Diamond(\psi^A) & (\Diamond \psi)^B := \Box(\psi^B) \\ (\Diamond \psi)^A := \Box(\psi^A) & (\Box \psi)^B := \Diamond(\psi^B) \\ (\Box \Diamond \psi)^A := \Diamond \Box(\psi^A) & (\Diamond \Box \psi)^B := \Box \Diamond(\psi^B) \end{array}$$

With the above definition in hand, it is straightforward to see the following duality theorem:

Theorem 2.2.5 (Duality). *Observe that for all $\phi \in \mathcal{L}_A(\Phi)$ and $\psi \in \mathcal{L}_B(\Phi)$, $(\phi^B)^A = \phi$ and $(\psi^A)^B = \psi$. Moreover, we have the following validities: $\models \neg(\phi^B) \leftrightarrow \phi$ and $\models \neg(\psi^A) \leftrightarrow \psi$.*

The above duality is convenient, since it can be leveraged to transfer results proven for the fragment $\mathcal{L}_A(\Phi)$ to $\mathcal{L}_B(\Phi)$ and vice versa.

With the above machinery in place, we can observe a natural consequence of the logical equivalences given in Lemma 2.2.1:

Definition 2.2.6. If $\phi \in \mathcal{L}_A(\Phi) \cup \mathcal{L}_B(\Phi)$ then let ϕ^* be the same formula, with all instances of \Box , \Diamond and $\Diamond \Box$ eliminated. That is, $(\cdot)^*$ has the following recursive definition:

$$\begin{array}{ll} p^* := p & (\neg p)^* := \neg p \\ \top^* := \top & \perp^* := \perp \\ \circlearrowleft^* := \circlearrowleft & (\neg \circlearrowleft)^* := \neg \circlearrowleft \\ (\phi \vee \psi)^* := (\phi^*) \vee (\psi^*) & (\phi \wedge \psi)^* := (\phi^*) \wedge (\psi^*) \\ (\Box \phi)^* := \Box(\phi^*) & (\Diamond \phi)^* := \Diamond(\phi^*) \\ (\Box \Diamond \phi)^* := \Box \Diamond(\phi^*) & (\Diamond \Box \phi)^* := \Diamond \Box(\phi^*) \end{array}$$

Theorem 2.2.7 (EVL Elimination). *For all $\phi \in \mathcal{L}_A(\Phi)$ or $\phi \in \mathcal{L}_B(\Phi)$, we have the following validity:*

$$\models \phi \leftrightarrow \phi^*$$

Proof. The proof proceeds in three steps.

Step 1: First, use induction on $\phi \in \mathcal{L}_A(\Phi)$, and show the following two facts simultaneously:

$$\models \exists \phi \leftrightarrow \phi \quad \models \Diamond \phi \leftrightarrow \phi$$

- Cases $p, \neg p, \perp, \top, \odot$: In all of these situations, the result follows directly from the validities illustrated in Lemma 2.2.1.
- Cases \wedge, \vee : For \exists the connective \wedge is simple, and dually for \Diamond for the connective \vee . This is because in each case one may simply use distribution, such as can be done here:

$$\begin{aligned} \models \exists(\phi \wedge \psi) &\leftrightarrow \exists \phi \wedge \exists \psi \\ &\leftrightarrow \phi \wedge \psi \end{aligned}$$

On the other hand, \vee is more interesting for \exists , and dually \wedge for \Diamond . Using induction, Lemma 2.2.1, and substitution, and distribution, we have the line of reasoning:

$$\begin{aligned} \models \exists(\phi \vee \psi) &\leftrightarrow \exists(\Diamond \phi \vee \Diamond \psi) \\ &\leftrightarrow \exists \Diamond(\phi \vee \psi) \\ &\leftrightarrow \Diamond(\phi \vee \psi) \\ &\leftrightarrow \Diamond \phi \vee \Diamond \psi \\ &\leftrightarrow \phi \vee \psi \end{aligned}$$

- Case \Diamond : Once again, this follows immediately from the validities of Lemma 2.2.1, namely $\models \exists \Diamond \phi \leftrightarrow \Diamond \phi$ and $\models \Diamond \Diamond \phi \leftrightarrow \Diamond \phi$
- Cases \exists, \Diamond : The final step follows from one more application of Lemma 2.2.1, namely by employing the following four validities

$$\begin{aligned} \models \exists \Diamond \phi &\leftrightarrow \Diamond \phi & \models \Diamond \Diamond \phi &\leftrightarrow \Diamond \phi \\ \models \exists \exists \phi &\leftrightarrow \exists \phi & \models \Diamond \exists \phi &\leftrightarrow \exists \phi \end{aligned}$$

Step 2: With the above, we can prove for any $\phi \in \mathcal{L}_A(\Phi)$ that $\models \phi \leftrightarrow \phi^*$. Once again, the proof proceeds by induction, the only steps worth noting involve \exists and \Diamond . In either case, these may be completed using Step 1. For instance, we know that $\models \exists \phi \leftrightarrow \phi$, hence $\models \exists \phi \leftrightarrow \phi^*$ by induction.

Step 3: With the result for $\mathcal{L}_A(\Phi)$ in hand, just observe that for $\psi \in \mathcal{L}_B(\Phi)$ we have that $(\psi^A)^* = (\psi^*)^A$. With this, substitution, and duality, we have the following chain of reasoning:

$$\begin{aligned} \models \psi &\leftrightarrow \neg(\psi^A) \\ &\leftrightarrow \neg((\psi^A)^*) \\ &\leftrightarrow \neg((\psi^*)^A) \\ &\leftrightarrow \neg(\neg(((\psi^*)^A)^B)) \\ &\leftrightarrow \neg\neg\psi^* \\ &\leftrightarrow \psi^* \end{aligned}$$

QED

Example 2.2.8. *The following validities of EVIL are consequences of Theorem 2.2.7:*

$$\begin{aligned} & \models \Box \Diamond t \vee \Box \Diamond \neg t \\ & \models ((\Box q \wedge \Box q) \vee \Box \Diamond q) \wedge ((\Box \Diamond q \vee \Box \Diamond q) \wedge \Box \Diamond q) \leftrightarrow \Box q \end{aligned}$$

IF

T

One way to read Theorem 2.2.7 is that \boxplus and \boxtimes are empty modalities on $\mathcal{L}_A(\Phi)$, and dually for $\mathcal{L}_B(\Phi)$ with \boxdot and \boxminus . Further, note that $\mathcal{L}_0(\Phi) = \mathcal{L}_A(\Phi) \cap \mathcal{L}_B(\Phi)$ (up to translation), which means that all four of \boxplus , \boxminus along with their duals \boxtimes and \boxdot vanish on the propositional language. Inspecting the semantics, this is to be expected, since neither \boxplus nor \boxminus interact with propositional truth values.

Finally, it should be mentioned that Theorem 2.2.7 reflects one of the basic themes of EVIL - the interplay between belief, reflected by \Box , and imagination, reflected by \Diamond . These two phenomena are just two sides of the same coin - furthermore, one could not have more natural opposites. Belief and imagination exemplify two warring forces dwelling within any EVIL agent's heart. Evidently soundness \odot is aligned with imagination and unsoundness $\neg \odot$ is aligned with belief.

2.2.2 Multiple Agents

In this section I turn to extending the semantics for EVIL from a single agent, as presented in §2.1.1, to accommodate multiple agents. This is primarily of interest since further results in EVIL, namely completeness, can naturally be abstracted beyond the single agent case. But I will freely admit that my EVIL intuitions are principally grounded in the single agent case – I recommend thinking about the multi-agent case as just a generalization of the single agent case.

The following provides the definition of the language of multi-agent EVIL:

Definition 2.2.9. Define $\mathcal{L}(\Phi, \mathcal{A})$ by the following Backus-Naur grammar:

$$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp \mid \Box_X \phi \mid \Box_X \phi \mid \Box_X \phi \mid \Box_X \phi$$

Here $X \in \mathcal{A}$, and \mathcal{A} is non-empty.

As in the single agent case, multi-agent EVIL models are sets $\mathfrak{M} \subseteq \wp\Phi \times (\wp\mathcal{L}_0(\Phi))^{\mathcal{A}}$ – that is, \mathfrak{M} is a set of pairs of sets of proposition letters, and indexed sets of propositional formulae.

The semantic entailment relation for multi-agent EVIL is

$$(\models) : \wp(\wp\Phi \times (\wp\mathcal{L}_0(\Phi))^{\mathcal{A}}) \times \wp\Phi \times (\wp\mathcal{L}_0(\Phi))^{\mathcal{A}} \times \mathcal{L}(\Phi, \mathcal{A}) \rightarrow \text{bool}.$$

The input/output behavior of (\models) is just as it was defined before in §2.1.1, the only difference in this setting is that instead of taking a pair as an input, where the second element is a set, it takes an indexed set.

I shall now provide a formal definition of the semantics for the multi-agent (\models) :⁸

Definition 2.2.10.

$$\begin{aligned} \mathfrak{M}, (a, A) \models p &\iff p \in a \\ \mathfrak{M}, (a, A) \models \phi \rightarrow \psi &\iff \mathfrak{M}, (a, A) \models \phi \text{ implies } \mathfrak{M}, (a, A) \models \psi \\ \mathfrak{M}, (a, A) \models \perp &\iff \text{False} \\ \mathfrak{M}, (a, A) \models \Box_X \phi &\iff \forall (b, B) \in \mathfrak{M}. (\forall \psi \in A_X. b \models \psi) \text{ implies } \mathfrak{M}, (b, B) \models \phi \\ \mathfrak{M}, (a, A) \models \Box_X \phi &\iff \forall (b, B) \in \mathfrak{M}. a = b \text{ and } B_X \subseteq A_X \text{ implies } \mathfrak{M}, (b, B) \models \phi \\ \mathfrak{M}, (a, A) \models \Box_X \phi &\iff \forall (b, B) \in \mathfrak{M}. a = b \text{ and } B_X \supseteq A_X \text{ implies } \mathfrak{M}, (b, B) \models \phi \\ \mathfrak{M}, (a, A) \models \bigcirc_X &\iff \forall \psi \in A_X. a \models \psi \end{aligned}$$

Just as in §2.1.1, Lemma 2.1.6 and Theorem 2.1.8 can be seen to obtain for the new generalized semantics. Furthermore, all of the validities mentioned in §2.1.3 and §2.2.1 hold, along with Theorem 2.2.7, where \Box , \Diamond , \Box_X , \Diamond_X , \Box_X , \Diamond_X , \Box_X , \Diamond_X , \Box_X , \Diamond_X and \bigcirc_X are all replaced with \Box_X , \Diamond_X , \Box_X , \Diamond_X , \Box_X , \Diamond_X , \Box_X , \Diamond_X , \Box_X , \Diamond_X and \bigcirc_X respectively, for any fixed $X \in \mathcal{A}$. Furthermore, compactness still fails, just as presented in §2.2.3.

Finally, there are two novel validities that these semantics give rise to:

$$\begin{aligned} \models \Box_X \phi \rightarrow \Box_X \Box_Y \phi \\ \models \Box_X \phi \rightarrow \Box_X \Diamond_Y \phi \end{aligned}$$

This is just to say, that the EVIL agent's deliberative process is opaque to other's beliefs, just as in the single agent case. This was expressed by (2.1.12) and (2.1.13) in §2.1.3. The agent cannot read anyone else's mind, nor anyone else hers.

Using the multi-agent semantics we have developed here, the proof theory for EVIL that shall be presented in §2.3 can now be given in higher generality.

⁸Where $X \in \mathcal{A}$, I use A_X to denote $A(X)$ provided that $A : \mathcal{A} \rightarrow \wp\mathcal{L}_0(\Phi)$

2.2.3 Kripke Structures

The language of EViL is evidently modal, and in previous sections the semantics have largely suggested that there are clear connections to conventional Kripke semantics. In this section, we will demonstrate that every EViL model corresponds to some highly structured Kripke model, with a minor modification on the standard definition. However, it will turn out that this correspondence is one way - the class of Kripke models for which EViL is strongly complete do not, in general, possess corresponding EViL models.

In order to understand EViL models as Kripke models, we return to the visualization technique for EViL models we introduced in §1.8. This involved thinking of the EViL models as *posets* with arrows, as we first presented in Fig. 3 in §1.8. We also saw additional examples of this visualization technique below in Figs. 6(a) and 6(b). Recall that these figures should be read as follows:

- if one point (a, A) is above another point (b, B) and connected by a densely dotted line \cdots , this means that $a = b$ and $B \subset A$.
- if one point (a, A) is connected to another point (b, B) by a line with an arrow \longrightarrow , this means that $\mathfrak{M}, (b, B) \models A$

In all of these depictions, the implicit relational structure of EViL models is given visual expression. So it seems only natural that this graphically perceived structure could also find formal expression.

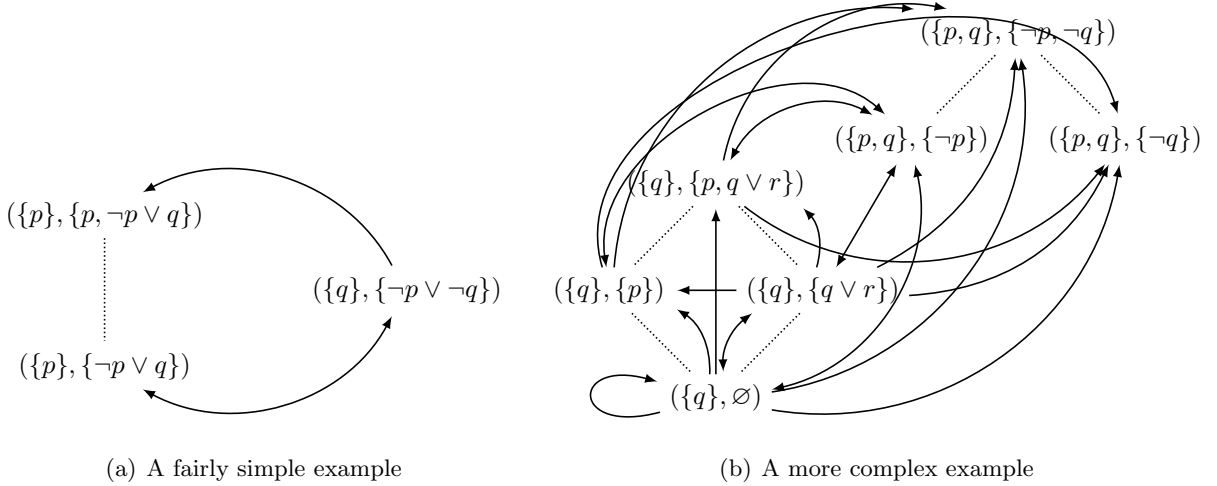


Figure 6: EViL model visualizations

Following the modified semantics provided in §2.2.2, the developments this section will assume multiple agents.

Definition 2.2.11. Let Φ be a set of letters and let \mathcal{A} be a set of agents. A **Kripke structure** is

a state transition system $\mathbb{M} = \langle W^{\mathbb{M}}, R^{\mathbb{M}}, \sqsubseteq^{\mathbb{M}}, \supseteq^{\mathbb{M}}, V^{\mathbb{M}}, P_{\odot}^{\mathbb{M}} \rangle$ where⁹:

- $W^{\mathbb{M}}$ is a set of worlds
- $R^{\mathbb{M}} : \mathcal{A} \rightarrow \wp(W \times W)$, $\sqsubseteq^{\mathbb{M}} : \mathcal{A} \rightarrow \wp(W \times W)$, and $\supseteq^{\mathbb{M}} : \mathcal{A} \rightarrow \wp(W \times W)$ are \mathcal{A} -indexed sets of relations¹⁰
- $V : \Phi \rightarrow \wp(W)$ is a predicate letter valuation
- $P_{\odot} : \mathcal{A} \rightarrow \wp(W)$ are sets of worlds indexed by agents

Let $\mathcal{K}_{\Phi, \mathcal{A}, I}$ denote the class of Kripke structures for letters Φ , agents \mathcal{A} , and where $W \subseteq I$.

Kripke semantics given by $(\Vdash) : \mathcal{K}_{\Phi, \mathcal{A}, I} \rightarrow I \rightarrow \text{bool}$ for these models are defined recursively as usual, granted the exceptional behavior of P_{\odot} .

Definition 2.2.12. Let \mathbb{M} be in the class $\mathcal{K}_{\Phi, \mathcal{A}, I}$

$$\begin{aligned}
\mathbb{M}, w \Vdash p &\iff w \in V^{\mathbb{M}}(p) \\
\mathbb{M}, w \Vdash \phi \rightarrow \psi &\iff \mathbb{M}, w \Vdash \phi \text{ implies } \mathbb{M}, w \Vdash \psi \\
\mathbb{M}, w \Vdash \perp &\iff \text{False} \\
\mathbb{M}, w \Vdash \Box_X \phi &\iff \forall v \in W^{\mathbb{M}}. w R_X^{\mathbb{M}} v \text{ implies } \mathbb{M}, v \Vdash \phi \\
\mathbb{M}, w \Vdash \Box_X \phi &\iff \forall v \in W^{\mathbb{M}}. w \sqsubseteq_X^{\mathbb{M}} v \text{ implies } \mathbb{M}, v \Vdash \phi \\
\mathbb{M}, w \Vdash \Box_X \phi &\iff \forall v \in W^{\mathbb{M}}. w \sqsubseteq_X^{\mathbb{M}} v \text{ implies } \mathbb{M}, v \Vdash \phi \\
\mathbb{M}, w \Vdash \odot_X &\iff w \in P_{\odot}^{\mathbb{M}}(X)
\end{aligned}$$

Kripke structures can be observed to typically have a lot less structure than EVIL models. On the other hand, EVIL models can be understood as Kripke structures in disguise. To illustrate this, observe the following lemma:

Definition 2.2.13 ($\mathcal{U}^{\mathfrak{M}}$ Translation). Let \mathfrak{M} be an EVIL model. Define $\mathcal{U}^{\mathfrak{M}} := \langle \mathfrak{M}, R^{\mathfrak{M}}, \sqsubseteq^{\mathfrak{M}}, \supseteq^{\mathfrak{M}}, V^{\mathfrak{M}}, P_{\odot}^{\mathfrak{M}} \rangle$, where

- $(a, A) R_X^{\mathfrak{M}}(b, B) \iff \forall \psi \in A_X. b \models \psi$
- $(a, A) \sqsubseteq_X^{\mathfrak{M}}(b, B) \iff a = b \text{ and } A_X \subseteq B_X$
- $(a, A) \supseteq_X^{\mathfrak{M}}(b, B) \iff a = b \text{ and } A_X \supseteq B_X$
- $(a, A) \in P_{\odot}^{\mathfrak{M}}(X) \iff \forall \psi \in A_X. a \models \psi$

Lemma 2.2.14. For all \mathfrak{M} and all $(a, A) \in \mathfrak{M}$, $\mathfrak{M}, (a, A) \models \phi$ if and only if $\mathcal{U}^{\mathfrak{M}}, (a, A) \Vdash \phi$

Proof. This follows from a straightforward induction on ϕ .

QED

The following summarizes the structural properties of EVIL models, when transformed into Kripke structures:

Proposition 2.2.15. For any EVIL model \mathfrak{M} , $\mathcal{U}^{\mathfrak{M}}$ has the following properties¹¹:

⁹Where the context is clear, we shall drop \mathbb{M}

¹⁰we shall abbreviate $R(X)$, $\sqsubseteq(X)$ and $\supseteq(X)$ as R_X , \sqsubseteq_X and \supseteq_X respectively.

¹¹Note that in this we have that $\{w, v\} \subseteq \wp \Phi \times \wp \mathcal{L}_0$ in the subsequent discussion

- (I) $\sqsupseteq_X^{\mathfrak{M}}$ is reflexive
- (II) $\sqsupseteq_X^{\mathfrak{M}}$ is transitive
- (III) $\sqsupseteq_X^{\mathfrak{M}}$ is anti-symmetric
- (IV) $w \sqsupseteq_X^{\mathfrak{M}} v$ if and only if $v \sqsubseteq_X^{\mathfrak{M}} w$
- (V) If $w \sqsubseteq_X^{\mathfrak{M}} v$ then $w \in V(p)$ if and only if $v \in V(p)$
- (VI) $(R_X^{\mathfrak{M}} \circ \sqsubseteq_X^{\mathfrak{M}}) \subseteq R_X^{\mathfrak{M}} \subseteq (R_X^{\mathfrak{M}} \circ \sqsupseteq_X^{\mathfrak{M}})$
- (VII) $(\sqsubseteq_Y^{\mathfrak{M}} \circ R_X^{\mathfrak{M}}) = R_X^{\mathfrak{M}} = (\sqsupseteq_Y^{\mathfrak{M}} \circ R_X^{\mathfrak{M}})$
- (VIII) $w \in P_{\odot}^{\mathfrak{M}}(X)$ if and only if $w R_X^{\mathfrak{M}} w$

... the situation in (VI) can be visualized in 7(a), while (VII) can be split into Figs. 7(b) and 7(c).

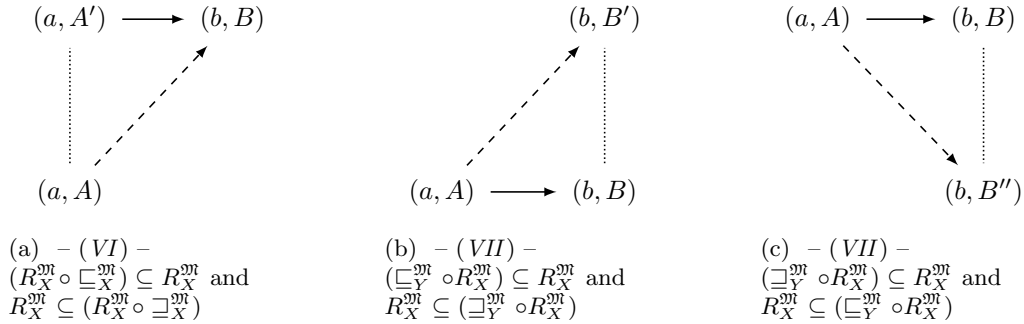


Figure 7: Visualizations of the relationships in Proposition 2.2.15

Proof. Everything except (VI) follows directly from the definitions - so we shall only demonstrate $R_X^{\mathfrak{M}} \subseteq R_X^{\mathfrak{M}} \circ \sqsupseteq_X^{\mathfrak{M}}$.

Suppose that $(a, A) R_X^{\mathfrak{M}} (b, B)$, then evidently $\forall \psi \in A_X.b \models \psi$. Now assume that $(a, A) \sqsupseteq_X^{\mathfrak{M}} (c, C)$. Then we know that $A_X \supseteq C_X$. But then $\forall \psi \in C_X.b \models \psi$, so $(c, C) R_X^{\mathfrak{M}} (b, B)$, which suffices to show the claim. QED

Definition 2.2.16. A Kripke structure is called *EvIL* if it makes true the above properties (I) through (VIII).

The Kripke semantics also serve to provide proper intuition behind EvIL models. We think of the defined relations given as follows:

- If $x R_X^{\mathfrak{M}} y$, then at world x the agent X can imagine y is true, since y is compatible with what the agent believes
- If $x \sqsubseteq_X^{\mathfrak{M}} y$, then at world x , agent X 's assumptions (or the experiences they are taking under consideration) are contained in her evidence at y

Given this perspective, the proof of (VI) can be understood in the following way - if the agent assumes fewer things, more things are imaginable, since it is easier for a world to be incompatible with an agent's evidence.

Finally, while Prop. 2.2.15 presents itself as a sort of representation lemma, the relationship between EViL semantics and Kripke semantics is not reciprocal. Not every Kripke model can be represented as an EViL model. Proposition 2.2.3 presents an elementary example of this failure of representation. It turns on the following observation:

Lemma 2.2.17. *For a given EViL model \mathfrak{M} , for any $\{(a, A), (b, B), (c, C)\} \subseteq \mathfrak{M}$, if $a = b$ then $a \models C$ if and only if $b \models C$*

Proof. This is an elementary result in the semantics of propositional logic. QED

Proposition 2.2.18 (Failure of Representation). *Consider a single agent EViL Kripke structure $\mathbb{M} := \langle W, R, \sqsubseteq, \supseteq, V, P \rangle$ where*

$$\begin{aligned} W &:= \{w, v\} & \sqsubseteq := \supseteq &:= \{(w, w), (v, v)\} \\ R &:= \{(w, v)\} & V(p) &:= \emptyset \text{ for all } p \in \Phi \\ P &:= \emptyset \end{aligned}$$

This structure is depicted in Fig. 8. No EViL model corresponds to \mathbb{M} .

Proof. Observe that the above model makes true the following:

$$\mathbb{M}, w \Vdash \Diamond \top \tag{2.2.1}$$

$$\mathbb{M}, w \Vdash \Box \neg p \text{ for all } p \in \Phi \tag{2.2.2}$$

$$\mathbb{M}, w \Vdash \neg p \text{ for all } p \in \Phi \tag{2.2.3}$$

$$\mathbb{M}, w \Vdash \neg \Diamond \Diamond \top \tag{2.2.4}$$

Armed with these observations, we can assert that it is impossible for there to be an EViL structure \mathfrak{M} with a world (a, A) such that $\mathbb{M}, w \Vdash \phi$ if and only if $\mathfrak{M}, (a, A) \models \phi$.

For suppose there were, then we could deduce the following facts, using the observations above:

- (1) From (2.2.1), there must be some pair $(b, B) \in \mathfrak{M}$ such that $b \models A$. Hence, A must be *consistent*.
- (2) From (2.2.2), we know that for the b mentioned in (1), it must be that $b = \emptyset$. This is a direct consequence of Lemma 2.1.6, the Truthiness Lemma.
- (3) From (2.2.3), evidently $a = \emptyset$
- (4) From (2.2.4), it must be that $a \not\models A$. Otherwise by the semantics of EViL as defined in §2.1.1 we would have $\mathfrak{M}, (a, A) \models \Diamond \Diamond \top$

Since $a = b = \emptyset$ and $b \models A$ then by Lemma 2.2.17 it must be that $a \models A$. But this clearly is absurd! \nmid QED

$$w \longrightarrow v$$

Figure 8: A Kripke Structure with no EViL representation

The above one way correspondence is admittedly inconvenient - it means that while EViL only enjoys some features from traditional epistemic logic, it is denied others. Despite this, EViL enjoys *most* of the benefits of basic modal logic. Indeed, we shall see in §2.3.4 that EViL is strongly complete for Kripke models with the properties described in Proposition 2.2.15.

Perhaps the most important formal feature that EViL semantics lacks in comparison to abstract Kripke semantics is that, as a consequence of the observations made in Proposition , EViL is not compact.

Theorem 2.2.19 (Failure of Compactness). *If the set of proposition letters Φ is infinite, then EViL is not compact for EViL semantics, as defined in from §2.1.1.*

Proof. We shall prove this result for the single agent case (the multiple agent case is a obvious generalization). Consider the function $\tau : \Phi \rightarrow \mathcal{L}(\Phi)$, defined as follows:

$$\tau(p) := (\Diamond \top) \wedge (\Box \neg p) \wedge (\neg p) \wedge (\neg \Diamond \Diamond \top)$$

We shall see that $\tau[\Phi]$ is finitely satisfiable, but not in its entirety.

Clearly not all of $\tau[\Phi]$ is satisfiable in EViL semantics, by the arguments presented in the proof of Proposition 2.2.3.

Now consider some finite subset of $S \subseteq_{\omega} \tau[\Phi]$. We shall construct a model that makes S true. Since τ is injective, we know there is some $\Psi \subseteq \Phi$ such that $S = \tau[\Psi]$. Since Φ is infinite, we know there is some $\rho \in \Phi \setminus \Psi$. Now consider a model $\mathfrak{M} = \{(\{\rho\}, \{\neg\rho\}), (\emptyset, \{\perp\})\}$. This is depicted in Fig. 9. It is straightforward to verify that $\mathfrak{M}, (\{\rho\}, \{\neg\rho\}) \models \tau[\Psi]$, so \mathfrak{M} is a suitable witness. QED

$$(\{\rho\}, \{\neg\rho\}) \longrightarrow (\emptyset, \{\perp\})$$

Figure 9: A model \mathfrak{M} where $\mathfrak{M}, (\{\rho\}, \{\neg\rho\}) \models \tau[\Psi]$ for $\Psi \subseteq_{\omega} \Phi$ and $\rho \notin \Psi$

A consequence of the failure of compactness, while strong completeness can be obtained for EViL using Kripke semantics, to achieve completeness for EViL semantics a finitary proof must be carried out.

We shall now turn to showing completeness for EViL.

2.3 EViL Completeness

In this section, we turn to providing a complete axiomatization of multi-agent EViL, as well as subsystems.

2.3.1 Axiom Systems

In this section, we shall present the axiom system which represents the validities of EVIL semantics as provided in .

Table 1, provides a Hilbert-style axiom system for EVIL. In addition to giving each axiom, we have also provided our own philosophical reading of what each axiom says. One unusual feature of this logic is that it is not *normal*, that is it is not closed under variable substitution.

This logic makes true a variety of relationships between the various modalities, which are given in the following lemma:

Lemma 2.3.1. *We have the following provable equivalences:*

$$\begin{array}{lll} \vdash \Box_X \phi \leftrightarrow \Box_X \Box_X \phi & \vdash \Box_X \phi \leftrightarrow \Box_X \Box_Y \phi & \vdash \Box_X \phi \leftrightarrow \Box_X \Box_Y \phi \\ \vdash \Box_X \phi \leftrightarrow \Box_X \Box_X \phi & \vdash \Box_X \phi \leftrightarrow \Box_X \Box_X \phi & \vdash \Box_X \leftrightarrow \Box_X \Box_X \end{array}$$

In addition to the main system presented above, it can be understood to contain two subsystems, corresponding to two fragments of the main grammar:

Definition 2.3.2. *Define $\mathcal{L}^\square(\Phi, \mathcal{A})$ as the fragment:*

$$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp \mid \Box_X \phi \mid \Box_X \phi \mid \Box_X$$

And define $\mathcal{L}^\boxplus(\Phi, \mathcal{A})$ as the fragment:

$$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp \mid \Box_X \phi \mid \Box_X \phi \mid \Box_X$$

Table 2 gives the axioms systems for these two fragments. For now, we shall observe that EVIL extends EVIL[□] and EVIL[⊕]. In §2.3.6 we shall make this precise.

From the definitions so far the following can be seen to hold:

Lemma 2.3.3 (Soundness). *If $\vdash \phi$ then for any model \mathfrak{M} and any $(a, A) \in \mathfrak{M}$ we have that $\mathfrak{M}, (a, A) \models \phi$*

The proof of the converse, that is *completeness*, proceeds by a three stage construction:

- The first step is to construct a Kripke model \mathfrak{K}^ϕ consisting of finite maximally consistent sets of formulae related to ϕ where $\mathfrak{K}^\phi, w \not\models \phi$ for some world $w \in W^{\mathfrak{K}^\phi}$. This model will be shown to make true nine properties.
- The second step is to construct a model $\mathfrak{K}^{\mathfrak{K}^\phi}$ which is bisimilar to \mathfrak{K}^ϕ . This model also makes true these nine properties as well as an additional tenth property.
- The final third step is to construct an EVIL model $\mathfrak{M}^{\mathfrak{K}^\phi}$. We shall show that for each $w \in W^{\mathfrak{K}^\phi}$ there is a corresponding $(a, A) \in \mathfrak{M}^{\mathfrak{K}^\phi}$ such that $\mathfrak{K}^{\mathfrak{K}^\phi}, w \models \psi$ if and only if $\mathfrak{M}^{\mathfrak{K}^\phi}, (a, A) \models \psi$ for all subformulae ψ of ϕ .

These three steps together suffice to prove completeness. we shall now proceed to demonstrate these constructions.

(1)	$\vdash \phi \rightarrow \psi \rightarrow \phi$	
(2)	$\vdash (\phi \rightarrow \psi \rightarrow \chi) \rightarrow (\phi \rightarrow \psi) \rightarrow \phi \rightarrow \chi$	<i>Axioms for basic propositional logic</i>
(3)	$\vdash (\neg\phi \rightarrow \neg\psi) \rightarrow \psi \rightarrow \phi$	
(4)	$\vdash \boxplus_X \phi \rightarrow \phi$	<i>If ϕ holds under any further evidence X considers, then ϕ holds simpliciter, since considering no additional evidence is trivially considering further evidence</i>
(5)	$\vdash \boxplus_X \phi \rightarrow \boxplus_X \boxplus_X \phi$	<i>If ϕ holds under any further evidence X considers, then ϕ holds whenever X considers even further evidence beyond that</i>
(6)	$\vdash p \rightarrow \boxplus_X p$	
(7)	$\vdash p \rightarrow \boxplus_X p$	<i>Changing one's mind does not bear on matters of fact</i>
(8)	$\vdash \Diamond_X \phi \rightarrow \boxplus_X \Diamond_X \phi$	<i>The more evidence X discards, the freer her imagination can run</i>
(9)	$\vdash \Box_X \phi \rightarrow \Box_X \boxplus_Y \phi$	
(10)	$\vdash \Box_X \phi \rightarrow \Box_X \boxplus_Y \phi$	<i>If X believes a proposition, she believes it regardless of what anyone else thinks</i>
(11)	$\vdash \bigcirc_X \rightarrow \Box_X \phi \rightarrow \phi$	<i>If X's premises are sound, then her logical conclusion are correct</i>
(12)	$\vdash \bigcirc_X \rightarrow \boxplus_X \bigcirc_X$	<i>If X's premises are sound then any subset of will be sound as well</i>
(13)	$\vdash \phi \rightarrow \boxplus_X \Diamond_X \phi$	
(14)	$\vdash \phi \rightarrow \boxplus_X \Diamond_X \phi$	<i>Embracing evidence is the inverse of discarding evidence</i>
(15)	$\vdash \Box_X (\phi \rightarrow \psi) \rightarrow \Box_X \phi \rightarrow \Box_X \psi$	
(16)	$\vdash \boxplus_X (\phi \rightarrow \psi) \rightarrow \boxplus_X \phi \rightarrow \boxplus_X \psi$	<i>Variations on axiom K</i>
(17)	$\vdash \boxplus_X (\phi \rightarrow \psi) \rightarrow \boxplus_X \phi \rightarrow \boxplus_X \psi$	
(I)	$\frac{\vdash \phi \rightarrow \psi \quad \vdash \phi}{\vdash \psi}$	<i>Modus Ponens</i>
(II)	$\frac{\vdash \phi}{\vdash \Box_X \phi}$	
(III)	$\frac{\vdash \phi}{\vdash \boxplus_X \phi}$	<i>Variations on necessitation</i>
(IV)	$\frac{\vdash \phi}{\vdash \boxplus_X \phi}$	

Table 1: A Hilbert style axiom system for EVIL

(1)	$\vdash \phi \rightarrow \psi \rightarrow \phi$	(1)	$\vdash \phi \rightarrow \psi \rightarrow \phi$
(2)	$\vdash (\phi \rightarrow \psi \rightarrow \chi) \rightarrow (\phi \rightarrow \psi) \rightarrow \phi \rightarrow \chi$	(2)	$\vdash (\phi \rightarrow \psi \rightarrow \chi) \rightarrow (\phi \rightarrow \psi) \rightarrow \phi \rightarrow \chi$
(3)	$\vdash (\neg\phi \rightarrow \neg\psi) \rightarrow \psi \rightarrow \phi$	(3)	$\vdash (\neg\phi \rightarrow \neg\psi) \rightarrow \psi \rightarrow \phi$
(4)	$\vdash \boxminus_X \phi \rightarrow \phi$	(4)	$\vdash \boxplus_X \phi \rightarrow \phi$
(5)	$\vdash \boxminus_X \phi \rightarrow \boxminus_X \boxminus_X \phi$	(5)	$\vdash \boxplus_X \phi \rightarrow \boxplus_X \boxplus_X \phi$
(6)	$\vdash p \rightarrow \boxminus_X p$	(6)	$\vdash p \rightarrow \boxplus_X p$
(7)	$\vdash \neg p \rightarrow \boxminus_X \neg p$	(7)	$\vdash \neg p \rightarrow \boxplus_X \neg p$
(8)	$\vdash \Diamond_X \phi \rightarrow \boxminus_X \Diamond_X \phi$	(8)	$\vdash \Box_X \phi \rightarrow \boxplus_X \Box_X \phi$
(9)	$\vdash \Box_X \phi \rightarrow \Box_X \boxminus_Y \phi$	(9)	$\vdash \Box_X \phi \rightarrow \Box_X \boxplus_Y \phi$
(10)	$\vdash \phi \rightarrow \boxminus_X (\bigcirc_X \rightarrow \Diamond_X \phi)$	(10)	$\vdash \phi \rightarrow \boxplus_X (\bigcirc_X \rightarrow \Diamond_X \phi)$
(11)	$\vdash \bigcirc_X \rightarrow \boxminus_X \bigcirc_X$	(11)	$\vdash \neg \bigcirc_X \rightarrow \boxplus_X \neg \bigcirc_X$
(12)	$\vdash \Box_X (\phi \rightarrow \psi) \rightarrow \Box_X \phi \rightarrow \Box_X \psi$	(12)	$\vdash \Box_X (\phi \rightarrow \psi) \rightarrow \Box_X \phi \rightarrow \Box_X \psi$
(13)	$\vdash \boxminus_X (\phi \rightarrow \psi) \rightarrow \boxminus_X \phi \rightarrow \boxminus_X \psi$	(13)	$\vdash \boxplus_X (\phi \rightarrow \psi) \rightarrow \boxplus_X \phi \rightarrow \boxplus_X \psi$
(I)	$\frac{\vdash \phi \rightarrow \psi \quad \vdash \phi}{\vdash \psi}$	(I)	$\frac{\vdash \phi \rightarrow \psi \quad \vdash \phi}{\vdash \psi}$
(II)	$\frac{\vdash \phi}{\vdash \Box_X \phi}$	(II)	$\frac{\vdash \phi}{\vdash \Box_X \phi}$
(III)	$\frac{\vdash \phi}{\vdash \boxminus_X \phi}$	(III)	$\frac{\vdash \phi}{\vdash \boxplus_X \phi}$

Table 2: Axiom system EvIL^\boxminus and EvIL^\boxplus respectively

2.3.2 Subformula Model Construction

In this section we provide definitions and lemmas related to the subformula construction \dagger^ϕ . We follow [Boo95] in our approach, as well as the “Fischer-Ladner Closure” used in the completeness theorem of PDL [BRV01].

Definition 2.3.4.

$$\sim \phi := \begin{cases} \psi & \text{if } \phi = \neg\psi \\ \neg\phi & \text{o/w} \end{cases} \quad \boxminus_X \phi := \begin{cases} \phi & \text{if } \phi = \boxminus_X \psi \\ \boxminus_X \phi & \text{o/w} \end{cases} \quad \boxplus_X \phi := \begin{cases} \phi & \text{if } \phi = \boxplus_X \psi \\ \boxplus_X \phi & \text{o/w} \end{cases}$$

Lemma 2.3.5. *By Lemma 2.3.1 we have*

$$\vdash \sim \phi \leftrightarrow \neg\phi \quad \vdash \boxminus_X \phi \leftrightarrow \boxminus_X \phi \quad \vdash \boxplus_X \phi \leftrightarrow \boxplus_X \phi$$

Moreover,

$$\boxminus_X \phi = \boxminus_X \boxminus_X \phi \quad \boxplus_X \phi = \boxplus_X \boxplus_X \phi$$

Definition 2.3.6. *Let $\delta(\phi) \subseteq \mathcal{A}$ be the set of agents that occur in ϕ ¹²*

Definition 2.3.7. *Define $\Sigma(\Delta, \phi)$ using primitive recursion as follows:*

$$\Sigma(\Delta, p) := \{p, \neg p, \perp, \neg\perp\} \cup \bigcup \{\{\boxminus_X p, \neg \boxminus_X p, \boxplus_X p, \neg \boxplus_X p\} \mid X \in \Delta\}$$

¹²In natural language, we read $\delta(\phi)$ as “the dudes mentioned by ϕ .”

$$\begin{aligned}
\Sigma(\Delta, \perp) &:= \{\perp, \neg\perp\} \\
\Sigma(\Delta, \odot_X) &:= \{\odot_X, \neg\odot_X, \boxplus_X \odot_X, \neg\boxplus_X \odot_X, \perp, \neg\perp\} \\
\Sigma(\Delta, \phi \rightarrow \psi) &:= \{\phi \rightarrow \psi, \neg(\phi \rightarrow \psi)\} \cup \Sigma(\Delta, \phi) \cup \Sigma(\Delta, \psi) \\
\Sigma(\Delta, \Box_X \phi) &:= \{\Box_X \phi, \neg\Box_X \phi, \boxplus_X \Box_X \phi, \neg\boxplus_X \Box_X \phi\} \\
&\quad \cup \bigcup \{\{\Box_X \Box_Y \phi, \neg\Box_X \Box_Y \phi, \Box_X \boxplus_Y \phi, \neg\Box_X \boxplus_Y \phi, \Box_Y \phi, \neg\Box_Y \phi, \boxplus_Y \phi, \neg\boxplus_Y \phi\} \mid Y \in \Delta\} \\
&\quad \cup \Sigma(\Delta, \phi) \\
\Sigma(\Delta, \boxplus_X \phi) &:= \{\boxplus_X \phi, \neg\boxplus_X \phi\} \cup \Sigma(\Delta, \phi) \\
\Sigma(\Delta, \boxtimes_X \phi) &:= \{\boxtimes_X \phi, \neg\boxtimes_X \phi\} \cup \Sigma(\Delta, \phi)
\end{aligned}$$

Lemma 2.3.8. $\Sigma(\delta(\phi), \phi)$ is finite. Moreover, we have the following:

- If $\psi \in \Sigma(\delta(\phi), \phi)$ then $\sim \psi \in \Sigma(\delta(\phi), \phi)$
- If $\psi \in \Sigma(\delta(\phi), \phi)$ and χ is a subformula of ψ , then $\chi \in \Sigma(\delta(\phi), \phi)$
- If $\boxplus_X \phi \in \Sigma(\delta(\phi), \phi)$ then $\Box_X \phi \in \Sigma(\delta(\phi), \phi)$
- If $\boxtimes_X \phi \in \Sigma(\delta(\phi), \phi)$ then $\boxplus_X \phi \in \Sigma(\delta(\phi), \phi)$

Definition 2.3.9. Let $At(\Psi)$ denote the maximally consistent subsets of Ψ

Lemma 2.3.10 (Lindenbaum Lemma). If $\Gamma \not\models \phi$ and $\Gamma \subseteq \Sigma(\delta(\phi), \phi)$, then there is a $\Gamma' \in At(\Sigma(\delta(\phi), \phi))$ such that $\Gamma \subseteq \Gamma'$ and $\Gamma' \not\models \phi$

Definition 2.3.11. Define $\mathbf{t}^\phi := \langle W^{\mathbf{t}^\phi}, V^{\mathbf{t}^\phi}, P_X^{\mathbf{t}^\phi}, R_{\Box_X}^{\mathbf{t}^\phi}, R_{\boxplus_X}^{\mathbf{t}^\phi}, R_{\boxtimes_X}^{\mathbf{t}^\phi} \rangle$ where:

$$\begin{aligned}
W^{\mathbf{t}^\phi} &:= At(\Sigma(\delta(\phi), \phi)) \\
V^{\mathbf{t}^\phi}(p) &:= \{w \in W^{\mathbf{t}^\phi} \mid p \in w\} \\
P_X^{\mathbf{t}^\phi} &:= \{w \in W^{\mathbf{t}^\phi} \mid \odot_X \in w\} \cup \{w \in W^{\mathbf{t}^\phi} \mid X \notin \delta(A)\} \\
R_{\Box_X}^{\mathbf{t}^\phi} &:= \{(w, v) \in W^{\mathbf{t}^\phi} \times W^{\mathbf{t}^\phi} \mid \{\psi \mid \Box_X \psi \in w\} \subseteq v\} \\
R_{\boxplus_X}^{\mathbf{t}^\phi} &:= \{(w, v) \in W^{\mathbf{t}^\phi} \times W^{\mathbf{t}^\phi} \mid \bigcup \{\{\psi, \boxplus_X \psi\} \mid \boxplus_X \psi \in w\} \subseteq v \wedge \bigcup \{\{\psi, \boxtimes_X \psi\} \mid \boxtimes_X \psi \in v\} \subseteq w\} \\
R_{\boxtimes_X}^{\mathbf{t}^\phi} &:= \{(v, w) \in W^{\mathbf{t}^\phi} \times W^{\mathbf{t}^\phi} \mid \bigcup \{\{\psi, \boxtimes_X \psi\} \mid \boxtimes_X \psi \in w\} \subseteq v \wedge \bigcup \{\{\psi, \boxplus_X \psi\} \mid \boxplus_X \psi \in v\} \subseteq w\}
\end{aligned}$$

Lemma 2.3.12 (Truth Lemma). For any subformula $\psi \in \Sigma(\delta(\phi), \phi)$ and any $w \in W^{\mathbf{t}^\phi}$, we have that $\mathbf{t}^\phi, w \models \psi$ if and only if $\psi \in w$

Proof. The proof proceeds by induction on ψ . Most of the steps are routine, with the exception of the right to left directions for the boxes.

We shall demonstrate the right to left direction for \boxplus_X . Assume that $\boxplus_X \psi \notin w$, then $w \not\models \boxplus_X \psi$. By Lemma 2.3.5 this is true if and only if $w \not\models \Box_X \psi$. Now abbreviate:

$$\begin{aligned}
A &:= \bigcup \{\{\chi, \boxplus_X \chi\} \mid \boxplus_X \chi \in w\} \\
B &:= \{\sim \boxtimes_X \chi \mid \boxtimes_X \chi \in \Sigma(\delta(\phi), \phi) \wedge \sim \chi \in w\}
\end{aligned}$$

Now suppose towards a contradiction that $\{\sim \psi\} \cup A \cup B \vdash \perp$. Then $A \cup B \vdash \psi$, and furthermore by Lemma 2.3.5 and rule (III) from the axioms we have that $\Box_X A \cup \Box_X B \vdash \Box_X \psi$.¹³ But then let

$$\begin{aligned} A' &:= \{\Box_X \chi \mid \Box_X \chi \in w\} \\ B' &:= \{\sim \chi \mid \sim \chi \in w\} \end{aligned}$$

Since $\Box_X \Box_X \chi = \Box_X \chi$ by Lemma 2.3.5, we have $A' = \Box_X A$. Moreover, by Lemma 2.3.5, axiom 13, and classical logic we can see that

$$\vdash \sim \chi \rightarrow \Box_X \sim \Box_X \psi$$

Thus for every $\beta \in \Box_X B$ we have that $B' \vdash \beta$. Hence by n applications of the Cut rule we can arrive at

$$A' \cup B' \vdash \Box_X \chi$$

However, evidently $A' \cup B' \subseteq w$, hence $w \vdash \Box_X \psi$, which contradicts what has been stipulated. \nmid

Hence it must be that $\{\sim \psi\} \cup A \cup B \not\vdash \perp$. In addition, from the fact that $w \subseteq \Sigma(\delta(\phi), \phi)$ with Lemma 2.3.8 and the hypothesis we have that $\{\sim \psi\} \cup A \cup B \subseteq \Sigma(\delta(\phi), \phi)$. Hence by the Lindenbaum Lemma we have that there is some $v \in At(\Sigma(\delta(\phi), \phi))$ such that $\{\sim \psi\} \cup A \cup B \subseteq v$. By the inductive hypothesis we have that $\mathfrak{t}^\phi, v \not\models \psi$.

To complete the argument, we must show that $w R_{\Box_X}^{\mathfrak{t}^\phi} v$. Since $A \subseteq v$ we just need to check that $\bigcup\{\{\psi, \Box_X \psi\} \mid \Box_X \psi \in v\} \subseteq w$. Suppose that $\Box_X \psi \in v$ but $\psi \notin w$. Since w is maximally consistent we have then that $\neg \psi \in w$. Thus $\sim \Box_X \psi \in w$, which contradicts that v is consistent. \nmid Now suppose that $\Box_X \psi \in v$ but $\Box_X \psi \notin w$, hence $\sim \Box_X \psi \in w$ and thus $\sim \Box_X \Box_X \psi \in w$. However we know from Lemma 2.3.5 that $\Box_X \Box_X \psi = \Box_X \psi$, which once again implies that v is inconsistent. \nmid QED

Lemma 2.3.13 (\mathfrak{t}^ϕ is Partly EViL). \mathfrak{t}^ϕ makes true the following properties:

- (1) $R_{\Box_X}^{\mathfrak{t}^\phi} \subseteq W^{\mathfrak{t}^\phi} \times W^{\mathfrak{t}^\phi}$
- (2) $W^{\mathfrak{t}^\phi}$ is finite
- (3) For all $w \in W^{\mathfrak{t}^\phi}$ we have $w R_{\Box_X}^{\mathfrak{t}^\phi} w$
- (4) If $w R_{\Box_X}^{\mathfrak{t}^\phi} v$ and $v R_{\Box_X}^{\mathfrak{t}^\phi} z$ then $w R_{\Box_X}^{\mathfrak{t}^\phi} z$
- (5) $R_{\Box_X}^{\mathfrak{t}^\phi} = (R_{\Box_X}^{\mathfrak{t}^\phi})^{-1}$
- (6) If $w R_{\Box_X}^{\mathfrak{t}^\phi} v$ then $w \in V^{\mathfrak{t}^\phi}(p)$ if and only if $v \in V^{\mathfrak{t}^\phi}(p)$
- (7) If $w R_{\Box_X}^{\mathfrak{t}^\phi} v$ and $v R_{\Box_X}^{\mathfrak{t}^\phi} u$ then $w R_{\Box_X}^{\mathfrak{t}^\phi} u$
- (8) If $w R_{\Box_X}^{\mathfrak{t}^\phi} v$ then $u R_{\Box_X}^{\mathfrak{t}^\phi} w$ if and only if $u R_{\Box_X}^{\mathfrak{t}^\phi} v$
- (9) If $w \in P_X^{\mathfrak{t}^\phi}$ then $w R_{\Box_X}^{\mathfrak{t}^\phi} w$

...for all $\{X, Y\} \subseteq \mathcal{A}$. Any model with the same modal similarity type as \mathfrak{t}^ϕ that makes the above true is said to be **partly EViL**

¹³Here $\Box_X S$ is shorthand for $\{\Box_X \chi \mid \chi \in S\}$.

Unfortunately, while \dagger^ϕ is nearly what is necessary to derive completeness for our semantics, it is not perfect. Another stage of the construction is necessary.

2.3.3 Bisimulation

We first introduce a Backus-Naur form grammar for the **Either** type constructor, which may be viewed as a coproduct in category theory (in the category of Sets)¹⁴:

$$\text{Either } a \text{ } b ::= a_l \mid b_r$$

Definition 2.3.14. Let \mathbb{M} be a Kripke model, then define $\dagger^\mathbb{M}$ as a model

$$\langle W^{\dagger^\mathbb{M}}, V^{\dagger^\mathbb{M}}, P_X^{\dagger^\mathbb{M}}, R_{\square_X}^{\dagger^\mathbb{M}}, R_{\boxplus_X}^{\dagger^\mathbb{M}}, R_{\boxminus_X}^{\dagger^\mathbb{M}} \rangle$$

Where:

$$\begin{aligned} W^{\dagger^\mathbb{M}} &:= \bigcup \{ \{w_l, w_r\} \mid w \in W^\mathbb{M} \} \\ V^{\dagger^\mathbb{M}}(p) &:= \bigcup \{ \{w_l, w_r\} \mid w \in V^\mathbb{M}(p) \} \\ P_X^{\dagger^\mathbb{M}} &:= \bigcup \{ \{w_l, w_r\} \mid w \in P_X^\mathbb{M} \} \\ R_{\square_X}^{\dagger^\mathbb{M}} &:= \bigcup \{ \{ (w_l, v_r), (w_r, v_l) \} \mid w R_{\square_X}^\mathbb{M} v \wedge w \notin P_X^\mathbb{M} \} \cup \bigcup \{ \{w_l, w_r\} \times \{v_l, v_r\} \mid w R_{\square_X}^\mathbb{M} v \wedge w \in P_X^\mathbb{M} \} \\ R_{\boxplus_X}^{\dagger^\mathbb{M}} &:= \bigcup \{ \{ (w_l, v_l), (w_r, v_r) \} \mid w R_{\boxplus_X}^\mathbb{M} v \} \\ R_{\boxminus_X}^{\dagger^\mathbb{M}} &:= \bigcup \{ \{ (w_l, v_l), (w_r, v_r) \} \mid w R_{\boxminus_X}^\mathbb{M} v \} \end{aligned}$$

Lemma 2.3.15. For any Kripke model $\mathbb{M} = \langle W, V, P_X, R_{\square_X}, R_{\boxplus_X}, R_{\boxminus_X} \rangle$, we have the following bisimulation Z between \mathbb{M} and $\dagger^\mathbb{M}$:

$$w Z w_l \quad \& \quad w Z w_r$$

Lemma 2.3.16. If \mathbb{M} is partly **EvIL** then $\dagger^\mathbb{M}$ is partly **EvIL** as well. It also makes true another, novel property:

$$(10) \text{ If } w R_{\square_X}^{\dagger^\mathbb{M}} w \text{ then } w \in P_X^{\dagger^\mathbb{M}}$$

Any partly **EvIL** Kripke model that makes true this tenth property is said to be **completely EvIL**.

2.3.4 Abstract Completeness

2.3.5 Translation

In the subsequent discussion, it will be useful to exploit certain properties of *partly EvIL* models. To this end we introduce the concept of a *column*.

Definition 2.3.17. Let \mathbb{M} be a partly **EvIL** Kripke structure. Let

$$\lceil w \rceil^\mathbb{M} := \{v \mid w(R_{\boxplus_X}^\mathbb{M} \cup R_{\boxminus_X}^\mathbb{M})^* v\}$$

Here R^* is the reflexive transitive closure of R .

¹⁴Either is taken from the functional programming language **Haskell**

Lemma 2.3.18 (Column Lemma). *The following hold if \mathbb{M} is partly EVIL:*

- (1) *For all w we have $w \in \ulcorner w \urcorner^{\mathbb{M}}$*
- (2) *If $w \in \ulcorner v \urcorner^{\mathbb{M}}$ then $\ulcorner w \urcorner^{\mathbb{M}} = \ulcorner v \urcorner^{\mathbb{M}}$*
- (3) *If $wR_{\Box_X}^{\mathbb{M}} v$ then for all $u \in \ulcorner v \urcorner^{\mathbb{M}}$ we have $wR_{\Box_X}^{\mathbb{M}} u$*
- (4) *If $w \in \ulcorner v \urcorner^{\mathbb{M}}$ then $w \in V^{\mathbb{M}}(p)$ if and only if $v \in V^{\mathbb{M}}(p)$ for all $p \in \Phi$*

Definition 2.3.19. *Let $L(\phi) := \{p \in \Phi \mid p \text{ is a subformula of } \phi\}$*

Let $\Lambda^{\mathbb{M}} := \bigcup \{\{\{w\}, \ulcorner w \urcorner^{\mathbb{M}}\} \mid w \in W^{\mathbb{M}}\}$

Let $\rho_{\phi}^{\mathbb{M}} : \Lambda^{\mathbb{M}} \rightarrow \Phi \setminus L(\phi)$ be an injection

Let $\theta_{\phi}^{\mathbb{M}} : W^{\mathbb{M}} \rightarrow \wp\Phi \times \wp(\mathcal{L}|_{Prop(\Phi)})$ be defined such that:

$$\begin{aligned} \theta_{\phi}^{\mathbb{M}}(w) := & (\{p \in L(\phi) \mid \mathbb{M}, w \Vdash p\} \cup \{\rho_{\phi}^{\mathbb{M}}(\ulcorner w \urcorner^{\mathbb{M}})\}, \lambda X. \{\neg \rho_{\phi}^{\mathbb{M}}(\ulcorner v \urcorner^{\mathbb{M}}) \mid \neg wR_{\Box_X}^{\mathbb{M}} v\} \\ & \cup \{\perp \rightarrow \rho_{\phi}^{\mathbb{M}}(\{v\}) \mid wR_{\Box_X}^{\mathbb{M}} v\}) \end{aligned}$$

Let $\boxtimes_{\phi}^{\mathbb{M}} := \theta_{\phi}^{\mathbb{M}}[W^{\mathbb{M}}]$

Lemma 2.3.20. *Let \mathbb{M} be a completely EVIL Kripke structure. Then for any subformula ψ of ϕ and any $w \in W^{\mathbb{M}}$, we have $\mathbb{M}, w \Vdash \psi$ if and only if $\boxtimes_{\phi}^{\mathbb{M}}, \theta_{\phi}^{\mathbb{M}}(w) \models \psi$*

Proof. Apply induction. The only challenging cases involve the boxes, so we shall illustrate $\Box_X \psi$.

Assume that $\mathbb{M}, w \not\Vdash \Box_X \psi$, then there's some $v \in W^{\mathbb{M}}$ such that $wR_{\Box_X}^{\mathbb{M}} v$ and $\mathbb{M}, v \not\Vdash \psi$. Let $(a, A) := \theta^{\mathbb{M}}(w)$ and $(b, B) := \theta^{\mathbb{M}}(v)$. By the inductive hypothesis it suffices to show that $\boxtimes_{\phi}^{\mathbb{M}}, (b, B) \models A_X$. But the only things in A_X are tautologies or formulae of the form $\neg \rho_{\phi}^{\mathbb{M}}(\ulcorner u \urcorner^{\mathbb{M}})$ where $\neg wR_{\Box_X}^{\mathbb{M}} u$. But then Lemma 2.3.18 it can't be that $\neg \rho_{\phi}^{\mathbb{M}}(\ulcorner v \urcorner^{\mathbb{M}}) \in A_X$, and this suffices.

Now assume that $\boxtimes_{\phi}^{\mathbb{M}}, (a, A) \not\models \Box_X \psi$ where $(a, A) = \theta^{\mathbb{M}}(w)$, so there must be some $v \in W^{\mathbb{M}}$ such that $\boxtimes_{\phi}^{\mathbb{M}}, (b, B) \not\models \psi$ where $(b, B) = \theta^{\mathbb{M}}(v)$ and $\boxtimes_{\phi}^{\mathbb{M}}, (b, B) \models A_X$. By the inductive hypothesis it suffices to show that $wR_{\Box_X}^{\mathbb{M}} v$, but this must be the case for otherwise $\neg \rho_{\phi}^{\mathbb{M}}(\ulcorner v \urcorner^{\mathbb{M}}) \in A_X$ and then it couldn't be that $\boxtimes_{\phi}^{\mathbb{M}}, (b, B) \models A_X$ since $\rho_{\phi}^{\mathbb{M}}(\ulcorner v \urcorner^{\mathbb{M}}) \in B$.

The inductive steps for the other boxes follow by similar reasoning. QED

2.3.6 Completeness

Theorem 2.3.21. *If $\not\models \phi$ then there is some model \mathfrak{M} and some $(a, A) \in \mathfrak{M}$ such that $\mathfrak{M}, (a, A) \not\models \phi$*

Proof. Assume $\not\models \phi$, then by Lemmas 2.3.12 and 2.3.13 we have some partly EVIL Kripke structure and world such that $\mathbb{A}, a \not\Vdash \phi$. By Lemma 2.3.15 we have that there is a completely EVIL Kripke structure \mathbb{B} such that $\mathbb{A} \leftrightarrow \mathbb{B}$, thus there is some world b such that $\mathbb{B}, b \not\Vdash \phi$. Finally by Lemma 2.3.20 we have that there's a model \mathfrak{C} in EVIL semantics and a pair $(c, C) \in \mathfrak{C}$ such that $\mathfrak{C}, (c, C) \not\models \phi$. QED

2.3.7 Conservativity, Decidability & Complexity

In this section, we discuss basic computability results for EvIL . All of the fragments of EvIL are decidable; furthermore, we shall establish a lower bound on the computational complexity.

We shall first prove the following lemma:

Lemma 2.3.22. *EvIL , EvIL^\Box and EvIL^\Box with a single agent are all conservative extensions of the basic modal logic with just axiom K . That is, if $\not\vdash_K \phi$ then $\not\vdash_{\text{EvIL}} \phi$ and similarly for the fragments EvIL^\Box and EvIL^\Box .*

EvIL with $m > n$ agents is a conservative extension of EvIL with n agents, and likewise for the fragments EvIL^\Box and EvIL^\Box

Proof. Assume that $\not\vdash_K \phi$, then we know from modal logic that there's a finite Kripke Structure $\mathbb{M} := \langle W, V, R \rangle$ such and a world $w \in W$ such that $\mathbb{M}, w \not\vdash \phi$. Now extend \mathbb{M} to $\mathbb{M}' := \langle W, V, P, R_\Box, R_\Box, R_\Box \rangle$ where

- $P := \{(v, v) \mid vRv\}$
- $R_\Box := R_\Box := \{(w, w) \mid w \in W\}$

This model is trivially completely EvIL . Moreover we know that \mathbb{M} is an elementary submodel of \mathbb{M}' , so $\mathbb{M}', w \not\vdash \phi$. Hence by the Lemma 2.3.20 we have a model \mathfrak{M} and $(a, A) \in \mathfrak{M}$ such that $\mathfrak{M}, (a, A) \not\vdash \phi$; so by soundness for EvIL we have the desired result.

Similarly, if we $\not\vdash_{\text{EvIL}_A} \phi$ then by completeness can find a witnessing \mathfrak{M} and $(a, A) \in \mathfrak{M}$ such that $\mathfrak{M}, (a, A) \not\vdash \phi$. But then we can embed \mathfrak{M} into \mathfrak{M}' for agents $\mathcal{B} \supseteq \mathcal{A}$ where $\mathfrak{M}' := \{(a, A') \mid (a, A) \in \mathfrak{M}\}$ and

$$A'_X := \begin{cases} A_X & X \in \mathcal{A} \\ \emptyset & X \notin \mathcal{A} \end{cases}$$

QED

By similar arguments, EvIL is a conservative extension of EvIL^\Box and EvIL^\Box , and that all three of these are conservative extensions of K . This is summerized in the Fig. 10.

Lemma 2.3.23. *EvIL is PSPACE hard*

Proof. This follows trivially from the fact that EvIL is a conservative extension of basic modal logic, and the decision problem for basic modal logic is PSPACE complete. QED

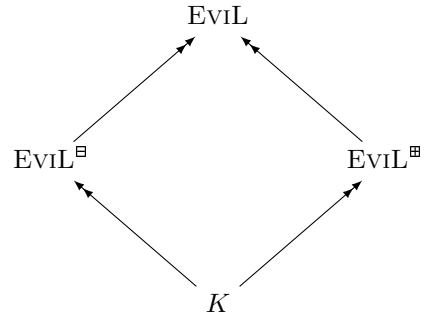


Figure 10: EViL conservative extensions of K

3 Applications

3.1 Collapse

3.2 Epistemic Plurality

3.2.1 Different Kinds of Knowledge

3.2.2 Moore's Paradox

3.2.3 Fitch's Paradox

3.3 Intuitionistic Logic

3.3.1 The Gödel Tarski McKinsey Embedding

3.3.2 Knowledge

3.3.3 Imagination

3.3.4 van Benthem $S4$

3.3.5 ImK_{\Box}

4 Epilogue

4.1 Comparison to Other Approaches

4.2 Failures

There are several points of failure of EVIL that I feel must be addressed:

- (I) EVIL is not really a logic, because it is non-normal and non-compact, so it therefore any kind of reasonable algebraic duality is impossible [for details on this, see BRV01, chapter 5]
- (II) EVIL is not dynamic and therefore fails to conform to the prevailing paradigm for epistemic logics
- (III) EVIL is not completely computer verified - only the completeness theorem for the central axiom system for EVIL has been produced; none of the many auxiliary results have been verified
- (IV) EVIL only partly accommodates irrationality

- (V) EViL is inhuman - the assumptions it makes for the nature of knowledge and EViL agent's cognitive abilities are unrealistic

A Grammars

$\mathcal{L}_{\text{therm}}$	$\phi ::= x \text{ Pascals} \mid y \text{ moles} \mid z \text{ Kelvin} \mid \phi \rightarrow \psi \mid \perp \mid \Box \phi$	pg. 6
\mathcal{L}_0	$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp$	pg. 10
$\mathcal{L}_K(\Phi)$	$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp \mid \Box \phi$	pg. 8
$\mathcal{L}_A(\Phi)$	$\phi ::= p \mid \neg p \mid \top \mid \perp \mid \circlearrowleft \mid \phi \wedge \psi \mid \phi \vee \psi \mid \Diamond \phi \mid \Box \phi \mid \Diamond \phi$	pg. 26
$\mathcal{L}_B(\Phi)$	$\phi ::= \neg p \mid p \mid \perp \mid \top \mid \neg \circlearrowleft \mid \phi \vee \psi \mid \phi \wedge \psi \mid \Box \phi \mid \Diamond \phi \mid \Box \phi$	pg. 27
	$\text{Either } a \mid b ::= a_l \mid b_r$	pg. 41

B Alternate Semantics

In this section, we shall present an alternative work to the framework proposed in §1.3. These semantics are inspired by game semantics for modal logic, such as those in [vB10], chapter 2.

First, recall the basic modal grammar $\mathcal{L}_K(\Phi)$:

$$\phi ::= p \in \Phi \mid \phi \rightarrow \psi \mid \perp \mid \Box \phi$$

Next, consider structures of the form $\langle W, V, \beta, \iota \rangle$ consisting of:

- A set of worlds W
- A propositional valuation function $V : \Phi \rightarrow \wp W$
- An belief function $\beta : W \rightarrow \wp \mathcal{L}_K(\Phi)$
- An imagination function $\iota : W \rightarrow \wp W$

We shall call these *belief-imagination models*. One can think of a model \mathfrak{M} sort of like a of tuples like in §2; however in this case evidently it would have to be $\mathfrak{M} \subseteq \wp \Phi \times \wp \mathcal{L}_K(\Phi) \times \wp \mathfrak{M}$, so apparently it would have to be a non-wellfounded set. This is somewhat natural, given a modal logic setting - see for instance [BM96] for an elaboration on these connections.

Definition B.0.1. *Define by recursion the following two truth relations:*

First relation:

$$\mathfrak{M}, w \Vdash p \iff p \in V(w)$$

$$\mathfrak{M}, w \Vdash \phi \wedge \psi \iff \text{both } \mathfrak{M}, w \Vdash \phi \text{ and } \mathfrak{M}, w \Vdash \psi$$

$\mathfrak{M}, w \Vdash \phi \vee \psi \iff \text{either } \mathfrak{M}, w \Vdash \phi \text{ or } \mathfrak{M}, w \Vdash \psi$

$\mathfrak{M}, w \Vdash \neg \phi \iff \mathfrak{M}, w \Vdash \phi$

$\mathfrak{M}, w \Vdash \Box \phi \iff \beta(w) \vdash^* \phi$

Where \vdash^* is a sequent that is closed under reflection and resolution:

$$\frac{\phi \in \Gamma}{\Gamma \vdash^* \phi} \quad \frac{\Gamma \vdash^* \neg \phi \vee \psi \quad \Delta \vdash^* \phi}{\Gamma \cup \Delta \vdash^* \psi}$$

Second relation:

$\mathfrak{M}, w \Vdash p \iff p \notin V(w)$

$\mathfrak{M}, w \Vdash \phi \wedge \psi \iff \text{either } \mathfrak{M}, w \Vdash \phi \text{ or } \mathfrak{M}, w \Vdash \psi$

$\mathfrak{M}, w \Vdash \phi \vee \psi \iff \text{both } \mathfrak{M}, w \Vdash \phi \text{ and } \mathfrak{M}, w \Vdash \psi$

$\mathfrak{M}, w \Vdash \neg \phi \iff \mathfrak{M}, w \Vdash \phi$

$\mathfrak{M}, w \Vdash \Box \phi \iff \text{there is some } v \in \iota(w) \text{ such that } \mathfrak{M}, v \Vdash \phi$

It is necessary to motivate the intuition behind these semantics. Informally, we think of these two truth relations correspond to two players, whom we shall call the *logician* and the *philosopher*. The logician wields a set beliefs given by β and tries to compose compelling arguments, and the philosopher employs a corpus of thought experiments given by ι to thwart the logician's arguments. Of course, the logician and the philosopher are really just two aspects of a single epistemic agent we are trying to model; we shall imagine epistemic agents modeled by this system to be embroiled in internal conflict. This sort of dissension between reason and imagination rages on within us all – it is fundamental to human nature.

These semantics are not naturally bivalent; that is it does not hold that either $\mathfrak{M}, w \Vdash \phi$ or $\mathfrak{M}, w \Vdash \neg \phi$, exclusively. To see this consider a model where $\beta(w) = \iota(w) = \emptyset$; then evidently $\mathfrak{M}, w \not\Vdash \Box p$ and $\mathfrak{M}, w \not\Vdash \neg \Box p$.

However, bivalence has a convenient semantic characterization:

Proposition B.0.2. *Let $\mathbb{M}^{\mathfrak{M}} = \langle W^{\mathfrak{M}}, V^{\mathfrak{M}}, R^{\mathfrak{M}} \rangle$ be a model for basic modal logic model based on a belief/imagination model \mathfrak{M} , where $wR^{\mathfrak{M}}v := v \in \iota(w)$, and let \Vdash_{\Box} be the modal truth predicate. We have that \Vdash and \Vdash_{\Box} are bivalent if and only if $\mathfrak{M}, w \Vdash \phi \iff \mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \phi$.*

Proof. (\implies) Assume that \Vdash and \Vdash_{\Box} are bivalent and consider any $\phi \in \mathcal{L}_K(\Phi)$. The proof that $\mathfrak{M}, w \Vdash \phi$ is equivalent to $\mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \phi$ proceeds by induction. The case for proposition letters, conjunction and disjunction are straightforward, so we shall only consider negation and modality.

Negation: We have the following chain of equivalences:

$$\begin{aligned}
\mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \neg\phi &\iff \mathbb{M}^{\mathfrak{M}}, w \nVdash_{\Box} \phi \\
&\iff \mathfrak{M}, w \nVdash \phi && \text{(inductive step)} \\
&\iff \mathfrak{M}, w \Vdash \phi && \text{(bivalence)} \\
&\iff \mathfrak{M}, w \Vdash \neg\phi
\end{aligned}$$

Modality: We have another chain of equivalences:

$$\begin{aligned}
\mathfrak{M}, w \Vdash \Box\phi &\iff \mathfrak{M}, w \nVdash \Box\phi && \text{(bivalence)} \\
&\iff \forall v \in \iota(w). \mathfrak{M}, w \nVdash \phi && \text{(definition)} \\
&\iff \forall v \in \iota(w). \mathfrak{M}, w \Vdash \phi && \text{(bivalence)} \\
&\iff \forall v \in \iota(w). \mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \phi && \text{(inductive step)} \\
&\iff \forall v. wR^{\mathfrak{M}}v \implies \mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \phi && \text{(definition)} \\
&\iff \mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \Box\phi
\end{aligned}$$

This completes the induction.

(\Leftarrow) Assume that $\mathfrak{M}, w \Vdash \phi$ and $\mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \phi$ are always equivalent. We have:

$$\begin{aligned}
\mathfrak{M}, w \Vdash \phi &\iff \mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \phi \\
&\iff \mathfrak{M}, w \nVdash_{\Box} \neg\phi \\
&\iff \mathfrak{M}, w \nVdash \neg\phi && \text{(hypothesis)} \\
&\iff \mathfrak{M}, w \Vdash \phi
\end{aligned}$$

QED

Corollary B.0.3. *If \Vdash and \Vdash are bivalent, then $\beta(w) \vdash^* \phi$ for all $\phi \in \text{Th}_{\Vdash}(\mathfrak{M})$ for all $w \in W^{\mathfrak{M}}$, where $\text{Th}_{\Vdash}(\mathfrak{M}) = \{\phi \in \mathcal{L}_K(\Phi) \mid \mathfrak{M}, w \Vdash \phi \text{ for all } w \in W^{\mathfrak{M}}\}$.*

Evidently bivalence of \Vdash and \Vdash gives rise to semantics where the agent has a proof for every proposition they believe. Furthermore, we can take any modal logic model $\mathbb{M} := \langle W^{\mathbb{M}}, V^{\mathbb{M}}, R^{\mathbb{M}} \rangle$ and define an equivalent belief/imagination model $\mathfrak{M}^{\mathbb{M}} := \langle W^{\mathbb{M}}, V^{\mathbb{M}}, \beta^{\mathbb{M}}, \iota^{\mathbb{M}} \rangle$ where:

$$\begin{aligned}
\beta^{\mathbb{M}}(w) &:= \{\phi \in \mathcal{L}_K(\Phi) \mid \mathbb{M}, w \Vdash_{\Box} \Box\phi\} \\
\iota^{\mathbb{M}}(w) &:= \{v \in W^{\mathbb{M}} \mid wR^{\mathbb{M}}v\}
\end{aligned}$$

We can immediately leverage this to give the a characterization of these semantics:

Proposition B.0.4. *The basic modal logic K is sound and strongly complete for bivalent belief/imagination models.*

Proof. Soundness is trivial given the previous lemma, strong completeness follows by considering the canonical model \mathbb{K} and looking at $\mathfrak{M}^{\mathbb{K}}$. QED

However, recalling on the remarks presented in §1.2, it is wrong for agents to be able to have everything they believe in their minds; this is about as bad as the thermometer theory of knowledge. However, this is evidently not entirely necessary. Call a belief/imagination model *reasonable* if the following two constraints are satisfied:

- $\beta(w) \vdash^* \phi$ for all $\phi \in \text{Th}_{\Vdash}(\mathfrak{M})$ for all $w \in W^{\mathfrak{M}}$, where $\text{Th}_{\Vdash}(\mathfrak{M}) = \{\phi \in \mathcal{L}_K(\Phi) \mid \mathfrak{M}, w \Vdash \phi \text{ for all } w \in W^{\mathfrak{M}}\}$
- $\text{Mod}_{\nVdash}^{\mathfrak{M}}(\beta(w)) \subseteq \iota(w)$, where $\text{Mod}_{\nVdash}^{\mathfrak{M}}(\beta(w)) = \{v \in W^{\mathfrak{M}} \mid \mathfrak{M}, v \nVdash \phi \text{ for all } \phi \in \beta(w)\}$
- $\beta(w) \setminus \text{Th}_{\Vdash}(\mathfrak{M})$ is finite

Evidently, forcing these requirements suffices to force bivalence:

Proposition B.0.5. *Let $\mathbb{M}^{\mathfrak{M}}$ be defined as in Prop. B.0.2. For any reasonable model \mathfrak{M} and any $w \in W^{\mathfrak{M}}$, we have:*

- (i) *If $\mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \phi$ then $\mathfrak{M}, w \Vdash \phi$*
- (ii) *If $\mathbb{M}^{\mathfrak{M}}, w \nVdash_{\Box} \phi$ then $\mathfrak{M}, w \Vdash \phi$*

Hence we have \Vdash and \Vdash are bivalent.

Proof. The propositional, disjunctive and conjunctive cases are all straightforward; we shall focus on negation and modality.

Negation: In the case of (i), we know that

$$\begin{aligned} \mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \neg \phi &\iff \mathbb{M}^{\mathfrak{M}}, w \nVdash_{\Box} \phi \\ &\implies \mathfrak{M}, w \Vdash \phi \quad (\text{by the inductive step}) \\ &\iff \mathfrak{M}, w \Vdash \neg \phi \end{aligned}$$

The proof for (ii) is similar.

Modality: In the case of (i), assume that $\mathbb{M}^{\mathfrak{M}}, w \Vdash_{\Box} \Box \phi$. Using the definition of reasonableness and the inductive step we know for all $v \in W^{\mathfrak{M}}$ that if $\mathfrak{M}, v \nVdash \psi$ for all $\psi \in \beta(w) \setminus \text{Th}(\mathfrak{M})$ then $\mathfrak{M}, v \Vdash \phi$.

From this and the fact that \mathfrak{M} is reasonable we can infer that $\bigvee_{\psi \in \beta(w) \setminus \text{Th}(\mathfrak{M})} \neg \psi \vee \phi \in \text{Th}_{\Vdash}(\mathfrak{M})$. We know further from reasonableness that we have $\text{Th}_{\Vdash}(\mathfrak{M}) \subseteq \beta(w)$. So we can prove by induction that repeatedly applying resolution gets $\beta(w) \vdash^* \phi$, which just means that $\mathfrak{M}, w \Vdash \Box \phi$, as desired.

The case of (ii) follows trivially by induction. QED

We may continue to obtain weak completeness for these semantics:

Proposition B.0.6. $\vdash_K \phi$ if and only if $\mathfrak{M}, w \Vdash \phi$ for all reasonable models \mathfrak{M} and $w \in W^{\mathfrak{M}}$

Proof. Left to right follows straightforwardly, so we just need to prove right to left.

Assume $\not\vdash_K \phi$. As before, let $\mathbb{M} = \langle W^{\mathbb{M}}, V^{\mathbb{M}}, R^{\mathbb{M}} \rangle$ be a finite model and with a world $w \in W^{\mathbb{M}}$ such that $\mathbb{M}, w \not\vdash_{\square} \phi$. Now consider a slightly modified model $\mathbb{M}' := \langle W^{\mathbb{M}}, V', R^{\mathbb{M}} \rangle$ where

$$V'(p) := \begin{cases} \{v\} & p = \rho(v) \\ V(p) & o/w \end{cases}$$

A proof by induction on subformulae ψ of ϕ verifies that $\mathbb{M}, w \vdash_{\square} \psi$ if and only if $\mathbb{M}', w \vdash_{\square} \psi$.

So now consider $\mathfrak{M} := \langle W^{\mathbb{M}'}, V^{\mathbb{M}'}, \tau, \lambda x. R^{\mathbb{M}'}[x] \rangle$ such that

$$\tau(w) := \text{Th}(\mathbb{M}') \cup \left\{ \bigvee_{v \in R^{\mathbb{M}'}[w]} \rho(v) \right\},$$

where $\text{Th}(\mathbb{M}') := \{\psi \in \mathcal{L}_K(\Phi) \mid \mathbb{M}', v \vdash \psi \text{ for all } v \in W^{\mathbb{M}'}\}$. A proof by induction on ψ shows that $\mathbb{M}', w \vdash_{\square} \psi$, $\mathfrak{M}, w \vdash \psi$ and $\mathfrak{M}, v \Vdash \psi$ are equivalent for all $\psi \in \mathcal{L}_K(\Phi)$. Thus we have that for all $v \in W^{\mathfrak{M}}$ that $\mathfrak{M}, v \Vdash \psi$ for all $\psi \in \text{Th}(\mathbb{M}')$. Moreover, evidently $w R^{\mathbb{M}'} v$ if and only if $\mathbb{M}', v \vdash_{\square} \bigvee_{u \in R^{\mathbb{M}'}[w]} \rho(u)$, whence we have that $w R^{\mathbb{M}'} v$ if and only if $\mathfrak{M}, v \not\vdash \chi$ for all $\chi \in \tau(w)$. With this we can employ induction and establish that $\mathbb{M}', w \vdash_{\square} \psi$ if and only if $\mathfrak{M}, w \models \psi$ for all $\psi \in \mathcal{L}_K(\Phi)$. Since $\mathbb{M}', w \not\vdash_{\square} \phi$, we have that $\mathfrak{M}, w \not\models \phi$. Finally, note that in this model we have that $\text{Mod}_{\Vdash}^{\mathfrak{M}}(\beta(w)) = R^{\mathbb{M}'}[w]$. With this and the definition of \mathfrak{M} , we can see that \mathfrak{M} is evidently reasonable, and thus we may complete the proof.

QED

Now, while reasonable models attain the goal of modeling agents that have proofs for the things they believe, they should not be considered adequate. These models are only reasonable in the sense that they indeed model agents providing nontrivial proofs for their beliefs. However, they are not reasonable in the sense that they are simple to reckon with. So while the semantics provided in §2 requires a grammar restriction, it should be preferred over the formulation given above, precisely because it is more manageable.

C An Application of Pure Model Theory to EviL Semantics

Recall that (VI), presented in Prop. 2.2.15 in §2.2.3 states:

Proposition C.0.7. *For any EviL model \mathfrak{M} , $\cup^{\mathfrak{M}}$ has the following property:*

$$(R_X^{\mathfrak{M}} \circ \sqsubseteq_X^{\mathfrak{M}}) \subseteq R_X^{\mathfrak{M}} \subseteq (R_X^{\mathfrak{M}} \circ \supseteq_X^{\mathfrak{M}}) \quad (VI)$$

In other words, if $(a, A) R_X^{\mathfrak{M}}(c, C)$ and $(a, A) \supseteq_X^{\mathfrak{M}}(b, B)$, then $(b, B) \supseteq_X^{\mathfrak{M}}(c, C)$

Recall that along with this principle, the following philosophical reading was offered:

“If the agent assumes fewer things, more things are imaginable, since it’s easier for a world to be incompatible with an agent’s evidence.”

In fact, in light of Theorem 2.1.8, the Theorem Theorem, the interplay expressed in (VI) follows from a general model theoretic relationship. For a given Kripke structure \mathbb{M} , define two operators $Mod^{\mathbb{M}} : \wp\mathcal{L}(\Phi, \mathcal{A}) \rightarrow \wp(W^{\mathbb{M}})$ and $Th^{\mathbb{M}} : \wp(W^{\mathbb{M}}) \rightarrow \wp\mathcal{L}(\Phi, \mathcal{A})$

$$\begin{aligned} Mod^{\mathbb{M}}(\Delta) &= \{x \in W \mid \forall \psi \in \Delta. \mathbb{M}, x \Vdash \psi\} \\ Th^{\mathbb{M}}(\nabla) &= \{\psi \in \mathcal{L}(\Phi, \mathcal{A}) \mid \forall x \in \nabla. \mathbb{M}, x \Vdash \psi\} \end{aligned}$$

We then have, for any $\Delta \in \wp\mathcal{L}(\Phi, \mathcal{A})$ and $\nabla \in \wp(W^{\mathbb{M}})$:

$$\nabla \subseteq Mod^{\mathbb{M}}(\Delta) \text{ if and only if } \Delta \subseteq Th^{\mathbb{M}}(\nabla)$$

From this, we may observe that these two operations form what is referred as an *antitone Galois connection*, between the lattice $\wp(W^{\mathbb{M}})$ and the lattice $\wp\mathcal{L}(\Phi, \mathcal{A})$. It follows from the theory of Galois connections [Rom08, chapter 3] that we have the following two properties:

$$\text{If } \nabla \supseteq \nabla' \text{ then } Th^{\mathbb{M}}(\nabla) \subseteq Th^{\mathbb{M}}(\nabla') \quad (\text{C.0.1})$$

$$\text{If } \Delta \supseteq \Delta' \text{ then } Mod^{\mathbb{M}}(\Delta) \subseteq Mod^{\mathbb{M}}(\Delta') \quad (\text{C.0.2})$$

We can see that (VI) follows from (C.0.2). To see this, assume that $(a, A) \sqsupseteq_X^{\mathfrak{M}} (b, B)$. Then observe:

$$\begin{aligned} (a, A) \sqsupseteq_X^{\mathfrak{M}} (b, B) &\implies a = b \text{ and } A_X \supseteq B_X && \text{by the definition of } \sqsupseteq_X^{\mathfrak{M}} \\ &\implies A_X \supseteq B_X && \text{weakening} \\ &\implies Mod^{\mathfrak{M}}(A_X) \subseteq Mod^{\mathfrak{M}}(B_X) && \text{from (C.0.2)} \\ &\implies \text{if } \mathfrak{M}, (c, C) \models A_X \text{ then } \mathfrak{M}, (c, C) \models B_X && \text{by the definition of } Mod^{\mathfrak{M}} \\ &\implies \text{if } (a, A) R_X^{\mathfrak{M}}(c, C) \text{ then } (b, B) R_X^{\mathfrak{M}}(c, C) && \text{by the definition of } R_X^{\mathfrak{M}} \end{aligned}$$

The above line of reasoning illustrates that structural features of EVIL models are consequences of the decision to set $\mathfrak{M}, (a, A) \models \Box\phi \iff Th(\mathfrak{M}) \cup A \vdash \phi$.

References

- [AGM85] Carlos E. Alchourron, Peter Gardenfors, and David Makinson. On the logic of theory change: Partial meet contraction and revision functions. *The Journal of Symbolic Logic*, 50(2):510–530, June 1985. ArticleType: primary_article / Full publication date: Jun., 1985 / Copyright © 1985 Association for Symbolic Logic.
- [AHV02] N. Agray, W. Van Der Hoek, and E. De Vink. On BAN logics for industrial security protocols. *Lecture notes in computer science*, page 29–36, 2002.
- [AN05] S. Artemov and E. Nogina. Introducing justification into epistemic logic. *Journal of Logic and Computation*, 15(6):1059, 2005.
- [Art07] S. N Artemov. Justification logic. *CUNY Graduate Center, New York*, 2007.
- [BM96] Jon Barwise and Lawrence Stuart Moss. *Vicious circles*. CSLI Publications, 1996.
- [Boo95] George Boolos. *The logic of provability*. Cambridge University Press, 1995.
- [Bro36] Sir Thomas Browne. *The garden of Cyrus*. printed in the year, 1736.
- [BRV01] P. Blackburn, M. De Rijke, and Y. Venema. *Modal logic*. Cambridge Univ Pr, 2001.
- [BS08] Alexandru Baltag and Sonja Smets. Probabilistic dynamic belief revision. *Synthese*, 165(2):179–202, November 2008.
- [CF04] Earl Conee and Richard Feldman. *Evidentialism*. Oxford University Press, USA, June 2004.
- [Cla34] É Clapeyron. Mémoire sur la puissance motrice de la chaleur. *J. l’école polytechnique*, 14:153–190, 1834.
- [Den98] Daniel Dennett. *The Intentional Stance*. MIT Press, Cambridge Mass, 7th edition, 1998.
- [DeP01] Michael Raymond DePaul. *Resurrecting old-fashioned foundationalism*. Rowman & Littlefield, 2001.
- [Eme08] Ralph Waldo Emerson. *Essays — First Series*. December 2008. LoC Class PS: Language and Literatures: American and Canadian literature.
- [Fit04] M. Fitting. A logic of explicit knowledge. *Logica Yearbook*, page 11–22, 2004.
- [Fit05] M. Fitting. The logic of proofs, semantically. *Annals of Pure and Applied Logic*, 132(1):1–25, 2005.
- [Fon08] Gaëlle Fontaine. Continuous fragment of the mu-Calculus. In *Computer Science Logic*, pages 139–153. 2008.
- [Fri06] J. Friedl. *Mastering regular expressions*. O’Reilly Media, Inc., 2006.

- [GG02] Dov M. Gabbay and F. Guenther. *Handbook of Philosophical Logic: Volume 6*. Springer, Dordrecht [u.a.], 2nd edition, May 2002.
- [Hal99] Joseph Y. Halpern. Set-theoretic completeness for epistemic and conditional logic. *Annals of Mathematics and Artificial Intelligence*, 26(1-4):1–27, 1999.
- [Hil22] David Hilbert. Neubegründung der mathematik. erste mitteilung. *Abhandlungen aus dem Mathematischen Seminar der Universität Hamburg*, 1(1):157–177, December 1922.
- [Hin69] Jaakko K. Hintikka. *Knowledge and Belief*. Cornell Univ. Pr., 1969.
- [HMdV05] A. Hommersom, J. J Meyer, and E. de Vink. Toward reasoning about security protocols: A semantic approach. *Electronic Notes in Theoretical Computer Science*, 126:53–75, 2005.
- [HMOV04] A. Hommersom, J. J Meyer, and E. De Vink. Update semantics of security protocols. *Synthese*, 142(2):229–267, 2004.
- [Hof79] Douglas Hofstadter. *Go?del, Escher, Bach : an eternal golden braid*. Basic Books, New York, 1979.
- [HS06] V. F Hendricks and J. Symons. Where’s the bridge? epistemology and epistemic logic. *Philosophical Studies*, 128(1):137–167, 2006.
- [KL86] Sarit Kraus and Daniel Lehmann. Knowledge, belief and time. In *Automata, Languages and Programming*, pages 186–195. 1986.
- [Koo03] Barteld P. Kooi. Probabilistic dynamic epistemic logic. *Journal of Logic, Language and Information*, 12(4):381–408, 2003.
- [KP82] L. Kirby and J. Paris. Accessible independence results for peano arithmetic. *Bulletin of the London Mathematical Society*, 14(4):285, 1982.
- [Len78] Wolfgang Lenzen. *Recent work in epistemic logic*. North-Holland, 1978.
- [Lev84] H. J Levesque. A logic of implicit and explicit belief. In *Proceedings of the National Conference on Artificial Intelligence*, pages 198–202, 1984.
- [LH59] E. J. Lemmon and G. P. Henderson. Symposium: Is there only one correct system of modal logic? *Proceedings of the Aristotelian Society, Supplementary Volumes*, 33:23–56, 1959. ArticleType: primary_article / Full publication date: 1959 / Copyright © 1959 The Aristotelian Society.
- [LL51] Clarence Irving Lewis and Cooper Harold Langford. *Symbolic Logic*. Dover Publications, 1951.
- [MvdH95] John-Jules Ch Meyer and Wiebe van der Hoek. *Epistemic Logic for AI and Computer Science*. Cambridge University Press, 1995.
- [Pla98] Plato. *The Republic*. October 1998. LoC Class PA: Language and Literatures: Classical Languages and Literature.

- [Pri06] Graham Priest. *Doubt truth to be a liar*. Clarendon Press; Oxford University Press, Oxford ; New York, 2006.
- [Qui51] W. V. Quine. Two dogmas of empiricism. *The Philosophical Review*, 60(1):20–43, January 1951. ArticleType: primary_article / Full publication date: Jan., 1951 / Copyright © 1951 Cornell University.
- [Ran82] V. Rantala. Impossible worlds semantics and logical omniscience. *Intensional Logic: Theory and Applications*, 1982.
- [Rom08] Steven Roman. *Lattices and Ordered Sets*. Springer, 1 edition, September 2008.
- [Rub98] Ariel Rubinstein. *Modeling Bounded Rationality*. MIT Press, 1998.
- [Rum02] Donald H. Rumsfeld. Defense.gov news transcript: DoD news briefing - secretary rumsfeld and gen. myers. <http://www.defense.gov/transcripts/transcript.aspx?transcriptid=2636>, February 2002.
- [vB91] J. van Benthem. Reflections on epistemic logic. *Logique Anal., Nouv. Sér.*, 34(133-134):5–14, 1991.
- [vB03] Johan van Benthem. Conditional probability meets update logic. *Journal of Logic, Language and Information*, 12(4):409–421, 2003.
- [vB10] Johan van Benthem. *Modal Logic for Open Minds*. Center for the Study of Language and Information, February 2010.
- [vBGK09] Johan van Benthem, Jelle Gerbrandy, and Barteld Kooi. Dynamic update with probabilities. *Studia Logica*, 93(1):67–96, October 2009.
- [vBV09] J. van Benthem and F. R Velázquez-Quesada. Inference, promotion, and the dynamics of awareness. *ILLC Amsterdam. To appear in Knowledge, Rationality and Action*, 2009.
- [Vel09] F. R Velázquez-Quesada. Inference and update. *Synthese*, 169(2):283–300, 2009.
- [VMTD05] John Vietch, David B. Manley, Charles S. Taylor, and René Descartes. Descartes’ meditations. <http://www.wright.edu/cola/descartes/>, July 2005.
- [Whi08] Walt Whitman. *Leaves of Grass*. August 2008. LoC Class PS: Language and Literatures: American and Canadian literature.