

## An Introduction to Modal Logic

Ordinary logic studies the partition of sentences<sup>1</sup> into two categories, true and false. Modal logic investigates a finer classification. A sentence can be either necessary (true, and it couldn't have been otherwise), contingently true (true, but it might have been false), contingently false (false, but it might have been true), or impossible (couldn't have been true). The informal study of the logical properties of modality goes back at least to Aristotle, but the advent of formal systems of symbolic modal logic dates from the publication of C. I. Lewis's *Survey of Symbolic Logic*<sup>2</sup> in 1918. Lewis started with the sentential calculus, and added symbols to represent "it is necessary that," "it is possible that," and "implies."<sup>3</sup> He then developed deductive calculi by adopting various appealing axioms and deriving their consequences.

What redeemed formal modal logic from empty symbol pushing was the development of possible-world semantics. The idea comes from Leibniz, who thought of God as surveying all possible worlds, selecting the one that was best, and making it actual. The idea that this is the best of all possible worlds is a nutty idea, even by philosophers' standards – see Voltaire's *Candide* – but it proved fruitful. Possible-world semantics for modal sentential calculus were developed by J.C.C. MacKensie and Alfred Tarski in the 1940s,<sup>4</sup> and they were extended to

---

<sup>1</sup> I am restricting attention to sentences used to make assertion, setting aside sentences used to ask questions or make requests or issue promises. A general theory that encompasses these other ways of employing language is developed by John Searle in *Speech Acts* (Cambridge: Cambridge University Press, 1969).

<sup>2</sup> Berkeley, Calif.: University of California Press.

<sup>3</sup> Recall my polemics in Logic I against reading "→" as "implies." The confusion of "implies," which is a transitive verb, with a sentential connective originates with Alfred North Whitehead and Bertrand Russell's *Principia Mathematica* (Cambridge: Cambridge University Press, 1910). See W. V. Quine, "Reply to Professor Marcus," *Synthese* 20 (1961), reprinted in Quine, *The Ways of Paradox* (Cambridge, Mass.: Harvard University Press, 1966).

<sup>4</sup> "On Closed Elements in Closure Algebra," *Annals of Mathematics* 47 (1946): 122-162.

modal predicate calculus by Saul Kripke in the 1960s.<sup>5</sup> We'll only talk about modal sentential calculus here, and that only superficially.<sup>6</sup>

We start with a version of the sentential calculus with infinitely many atomic sentences and add an operator " $\Box$ ," read "it is necessary that"; if  $\phi$  is a sentence, so is  $\Box\phi$ . We treat Lewis's symbol for possibility, " $\Diamond$ ," as defined:  $\Diamond\phi =_{\text{Def}} \sim \Box \sim \phi$ . We'll have no need for Lewis's symbol for implication.

Our interest here is in the interpretation of the modal sentential calculus in which " $\Box$ " means "It is provable that," so that " $\Diamond$ " means "It is consistent that."

A *Kripke model* is an ordered quadruple  $\langle W, R, I, a \rangle$ , where  $W$ , the set of *worlds*, is a nonempty set;  $R$ , the *accessibility relation* is a binary relation on  $W$ ;  $I$ , the *interpretation function*, is a function that assigns to each pair  $\langle \phi, w \rangle$  with  $\phi$  a sentence and  $w$  a world either the value 0 or the value 1; and  $a \in W$  is the *actual world*. The triple  $\langle W, R, I \rangle$  is a *frame*. For  $\phi$  and atomic sentence and  $w$  a world,  $\phi$  is true in  $w$  if and only if  $I(\phi, w) = 1$ . A conjunction is true in  $w$  iff both conjuncts are true in  $w$ , a disjunction is true in  $w$  if and only if one or both disjuncts are true in  $w$ , and so on.  $\Box\phi$  is true in  $w$  iff  $\phi$  is true in every world  $v$  with  $Rwv$ . A sentence is *true* in the model iff it's true in  $a$ .

A sentence is *valid* for a frame or set of frames iff it's true at every world in every member of the set.

---

<sup>5</sup> "Semantical Considerations on Modal Logic," *Acta Philosophica Fennica* 16 (1963): 83-94.

<sup>6</sup> For a fuller treatment, see Brian Chellas, *Modal Logic* (Cambridge: Cambridge University Press, 1980), G. E. Hughes and Max Cresswell, *A New Introduction to Modal Logic* (Routledge, 1968), or J. C. Beall and Bas van Fraassen, *Possibilities and Paradox* (Oxford: Oxford University Press, 2003).

A *normal modal system* is a set of sentences  $\Gamma$  with the following properties:

*Tautological consequence:* Every tautological consequence of  $\Gamma$  is in  $\Gamma$ .

*Necessitation:* If  $\phi$  is in  $\Gamma$ , so is  $\Box\phi$ .

*Schema (K):* Each instance of the axiom schema

$$(K) \quad (\Box(\phi \rightarrow \psi) \rightarrow (\Box\phi \rightarrow \Box\psi))$$

is in  $\Gamma$ .

**Theorem.** For  $\Gamma$  a set of sentences, the following are equivalent:

- (i)  $\Gamma$  is a normal modal system.
- (ii) There is a class of frames such that  $\Gamma$  is the set of sentences valid for every member of the class
- (iii) Either  $\Gamma$  is the set of all sentences or there is a frame  $\langle W, R, I \rangle$  such that  $\Gamma$  is the set of sentences valid for  $\langle W, R, I \rangle$ .

**Proof:** That (ii) implies (i) is easy to check. That (iii) implies (ii) is immediate; if  $\Gamma$  is the set of all sentences, our class of frames will be the empty class. So we only need to worry about showing that (i) implies (iii). Given  $\Gamma$  a normal modal system, let's say a set of sentences is  $\Gamma$ -consistent iff it contains all the members of  $\Gamma$  and it is consistent by the sentential calculus. A *maximal  $\Gamma$ -consistent* set is a  $\Gamma$ -consistent set such that, for every sentence, either the sentence or its negation is in the set. Let  $\mathcal{W}$  be the class of all maximal  $\Gamma$ -consistent sets of sentences; unless  $\Gamma$  is the set of all sentences,  $\mathcal{W}$  will be nonempty. For  $u$  and  $v$  elements of  $\mathcal{W}$ , define  $Ruv$  iff  $\phi$  is in  $v$  whenever  $\Box\phi$  is in  $u$ . Define  $I(\phi, w) = 1$  iff  $\phi \in w$ , for  $\phi$  atomic. We want to show that, for any sentence  $\phi$ ,  $\phi$  is true in  $w$  in the frame  $\langle \mathcal{W}, R, I \rangle$  iff  $\phi \in w$ . This will tell us that the following are equivalent:

$\phi \in \Gamma$

$\phi$  is an element of every maximal  $\Gamma$ -consistent set

$\phi$  is an element of every world  $w$  in  $W$

$\phi$  is true in every world  $w$  in  $W$

$\phi$  is valid for the frame  $\langle W, R, I \rangle$

The proof that, for any sentence  $\phi$ ,  $\phi$  is true in  $w$  if and only if  $\phi \in w$  proceeds by induction on the complexity of  $\phi$ . The only part of this that isn't entirely routine is to show that  $\Box\psi$  is true in  $w$  iff it's an element of  $w$ . Here's the proof of the right-to-left direction: If  $\Box\psi$  is an element of  $w$ , then, by definition of  $R$ ,  $\psi$  is an element of every world accessible from  $w$ . It follows by inductive hypothesis that  $\psi$  is true in every world accessible from  $w$ , that is, that  $\Box\psi$  is true in  $w$ .

For the other direction, suppose that  $\Box\psi$  isn't an element of  $w$ . We want to see that there is a world accessible from  $w$  in which  $\psi$  isn't true. This means, according to the inductive hypothesis, that we want a world accessible from  $w$  that contains  $\psi$ . That is, given the definition of  $R$ , we want a maximal  $\Gamma$ -consistent set of sentences that includes all the sentences  $\theta$  with  $\Box\theta$  in  $w$  but that doesn't include  $\psi$ . To get this, it will suffice to show that  $\Gamma \cup \{\text{sentences } \theta: \Box\theta \in w\} \cup \{\sim\psi\}$  is tautologically consistent. If it is, we can expand  $\Gamma \cup \{\theta: \Box\theta \in w\} \cup \{\sim\psi\}$  to a maximal  $\Gamma$ -consistent set by the familiar technique of marching through the sentences one by one, for each sentence when we come to it adding either it or its negation to the set, preserving  $\Gamma$ -consistency at every stage. Because of Necessitation, if  $\gamma$  is in  $\Gamma$ ,  $\Box\gamma$  is in  $\Gamma$ , and so  $\Box\gamma$  is in  $w$  and  $\gamma$  is in  $\{\theta: \Box\theta \in w\}$ . So it will be enough to show that  $\{\theta: \Box\theta \in w\} \cup \{\sim\psi\}$  is tautologically consistent. If not, then there exist sentences  $\theta_1, \theta_2, \dots, \theta_n$  such that each  $\Box\theta_i$  is in

w and such that  $(\theta_1 \rightarrow (\theta_2 \rightarrow \dots \rightarrow (\theta_n \rightarrow \psi)\dots))$  is a tautology. It follows by Tautological Consequence and Necessitation that  $\Box(\theta_1 \rightarrow (\theta_2 \rightarrow \dots \rightarrow (\theta_n \rightarrow \psi)\dots))$  is a member of  $\Gamma$ , and hence, by multiple applications of Tautological Consequence and Schema (K), that  $(\Box\theta_1 \rightarrow (\Box\theta_2 \rightarrow \dots \rightarrow (\Box\theta_n \rightarrow \Box\psi)\dots))$  is in  $\Gamma$ , and hence in w. Because w is closed under *modus ponens*, it follows that  $\Box\psi$  is in w, contrary to our assumption.

The frame  $\langle W, R, I \rangle$  that we just constructed is called the *canonical frame* for  $\Gamma$ . The principal moral of the theorem is that, if a sentence is outside  $\Gamma$ , then there is a world in the canonical frame in which it is false.  $\square$

Let me write down some axioms schemata; the schemata were named by different people at different times, so the nomenclature is annoyingly haphazard:

- (T)             $(\Box\phi \rightarrow \phi)$
- (4)             $(\Box\phi \rightarrow \Box\Box\phi)$
- (B)             $(\phi \rightarrow \Box\Diamond\phi)$
- (5)             $(\Diamond\phi \rightarrow \Box\Diamond\phi)$

Let me also write down some notable properties of binary relations:

R is a *reflexive* relation on W iff, for each w in W, we have Rww.

R is *transitive* iff, for each u, v, and w, if Ruv and Rvw, then Ruw.

R is *symmetric* iff, for each u and v, if Ruv, then Rvu.

R is *Euclidean* iff, for each u, v, and w, if Ruv and Ruw, then Rvw.

K is defined to be the smallest normal modal system (that is, every other normal modal system includes K), so that a sentence is an element of K iff it is derivable from instances of schema (K) by the rules TC and Necessitation. A sentence is in K iff it is true in every world in

every frame. Why? The set of sentences valid for every frame is a normal modal system, so it includes K. If  $\phi$  isn't in K, then there is a frame in which there is a world in which  $\phi$  is false, namely, the canonical model for K.

KT is defined to be the smallest normal modal system that includes (T), so that a sentence is an element of KT iff it is derivable from (K) and (T) by the rules TC and Necessitation. A sentence is a element of KT iff it is true in every world in every reflexive frame. Why? Given a model  $\langle W, R, I, a \rangle$ , with R reflexive, if  $\Box\phi$  is true in a, then  $\phi$  is true in every world accessible from a; in particular,  $\phi$  is true in a itself; so all instances of schema (T) are true in the model. Consequently, the set of sentences valid for every reflexive frame is a normal modal system that includes (T). Moreover, the canonical frame for KT is reflexive; for any world w in the canonical frame, if  $\Box\phi$  is in w,  $\phi$  is in w, so we have  $Rww$ . Thus, if  $\phi$  isn't in KT, then there is a reflexive frame in which there is a world in which  $\phi$  is false, namely, the canonical frame for KT.

K4 is defined to be the smallest normal modal system that includes (4). A sentence is an element of K4 iff it's true in every world in every transitive frame. Why? Given a model  $\langle W, R, I, a \rangle$  with R transitive, if  $\Box\phi$  is true in a and w is a world accessible from a, then every world accessible from w is accessible from a.. Since  $\phi$  is true in every world accessible from a,  $\phi$  must be true in every world accessible from w, so that  $\Box\phi$  is true in w. We have shown that  $\Box\phi$  is true in every world accessible from a, so that  $\Box\Box\phi$  is true in a. Thus we see that all instances of schema (4) are true in the model, so that the set of sentences true in every transitive model will be a normal modal system that includes (4). Moreover, the canonical frame for K4 is transitive. If u, v, and w are worlds in the canonical frame for K4 with  $Ruv$  and  $Rvw$ , then if  $\Box\phi$  is in u,  $\Box\Box\phi$  is in u, so that  $\Box\phi$  is in v and  $\phi$  is in w. Consequently,  $Ruw$ . Thus, if  $\psi$  is not in

KT, then the canonical frame for KT will be a transitive frame in which there is a world in which  $\phi$  is false.

KB is the smallest normal system that contains (B). A sentence is in KB iff it's valid for the class of symmetric frames.

K5 is the smallest normal modal system that includes (5). A sentence is in K5 iff it's valid for the class of Euclidean frames.

KT4, which Lewis called "S4," is the smallest normal modal system that includes both (T) and (4). A sentence is in KT4 iff it's valid for the class of reflexive, transitive frames.

KTB is the smallest normal modal system that includes both (T) and (B). A sentence is in KTB iff it's valid for the class of reflexive, symmetric frames.

KT5, which Lewis called (S5), is the smallest normal modal system that includes both (T) and (5). A sentence is in KT5 iff it's valid for the class of reflexive, Euclidean frames. Since a binary relation that is reflexive and Euclidean will also be transitive and symmetric, KT5 is the same as KT4B5.

I could keep doing this for a long time, but you get the point.

The mystery component of the story is the accessibility relation. I have never heard a remotely satisfying explanation of why one possible world should or should not be accessible

from another.<sup>7</sup> As far as I can tell, the accessibility relation is something we tack on *ad hoc* so as to get the pretty relations between frames and axiom systems.

---

<sup>7</sup> We sometimes think of the march of history as following a forked path through a continually unfolding array of branching possibilities, so that what happens now can constrain what will be possible tomorrow. It could happen that, as of now, it's possible that I should fly to Jamaica tomorrow, but that some untoward event could happen tonight that would render it impossible for me to fly to Jamaica tomorrow; I might, for example, be eaten by a tiger escaped from a circus. So, even though it's not possible for me to fly to Jamaica tomorrow, it might become impossible for it to be possible for me to fly to Jamaica tomorrow, so that some instances of schema (5) can fail. To represent this conception formally, we take a "possible world" to be, not (as we would have expected) a possible complete course of history, but rather an ordered pair consisting of a complete course of history and a time. A statement will be possible at  $\langle h, t \rangle$  if its truth is compatible with the course of history according to  $h$  as it has unfolded up to time  $t$ . A statement not involving modality will be true in  $\langle h, t \rangle$  iff it's true in  $h$ . A pair  $\langle h', t' \rangle$  will be accessible from  $\langle h, t \rangle$  if  $t'$  is later than or equal to  $t$  and if  $h$  and  $h'$  agree in their depiction of the history of the world up to time  $t$ . The appropriate modal logic will be KT4.