# AI2Bh Modal Logic

## Michael Fourman
## Revised by Alan Bundy

## May 12, 2004

Modal logics were introduced to formalise *modalities* such as necessity and possibility.

**Syntax**  The, now standard, syntax for these modalities is

$$\Box\varphi \text{ to express, } \varphi \text{ is necessarily true}$$

$$\Diamond\varphi \text{ to express, } \varphi \text{ is possibly true}$$

These are normally considered definable in terms of each other:

$$\Diamond\varphi \Leftrightarrow \neg\Box\neg\varphi$$

We will use modal logic to formalise the assertion that an agent, $A$, *knows that $\varphi$*. Knowledge is formally similar to neccessity, so we use the (reasonably standard) syntax,

$$\boxed{\mathsf{K}_A}\,\varphi \text{ to express, } A \text{ knows that } \varphi$$

**Example**  We can use this syntax to express complex statements about knowledge, and make subtle differences explicit and precise. For example,

$$\exists x.\,\boxed{\mathsf{K}_A}\,\varphi(x) \text{ means, } \textit{for some } x, A \textit{ knows that } \varphi(x)$$

$$\boxed{\mathsf{K}_A}\,\exists x.\,\varphi(x) \text{ means, } A \textit{ knows that, for some } x, \varphi(x)$$

Suppose $\varphi(x)$ means, *x is the name of the oldest person in Edinburgh,* and you are $A$. Which of the two statements is true? (Possibly both are, but that is unlikely.)

**Semantics**  The semantics we have in mind is based on a notion of possible worlds. The idea is that in addition to the actual world, which provides a standard interpretation of predicate logic, we consider other possible worlds. Of all the possibilities that could occur, only some are compatible with $A$'s knowledge – but these are all, from $A$'s perspective, equivalent. So, given a world, $W$, for each agent, $A$, there is a set of worlds equivalent, from $A$'s perspective, to $W$. Formally, we have a set of possible worlds, and for each agent, $A$, an equivalence relation $\equiv_A$. Informally, we interpret $\boxed{\mathsf{K}_A}\,\varphi$ to mean that $\varphi$ holds in all the worlds possible for $A$. Formally,

$$W \vDash \boxed{\mathsf{K}_A}\,\varphi \textbf{ iff } \text{for all } V \text{ such that } V \equiv_A W, \text{ we have } W \vDash \varphi$$

Implicitly, this attributes to $A$ omniscience: infallible, and exhaustive powers of reasoning, since

$$\text{if } \boxed{\mathsf{K}_A}\,\varphi \text{ and } \varphi \vDash \psi \text{ then } \boxed{\mathsf{K}_A}\,\psi$$

As a model of human knowledge, this is clearly flawed; nevertheless, it gives a useful system that approximates some aspects of knowledge as we know it.

**Inference** We use the possible world semantics to justify the validity of some reasoning principles for modal logic.

Suppose that in all an agent, $A$'s possible worlds, both $\varphi \to \psi$ and $\varphi$ hold—*i.e.* suppose $\boxed{\mathsf{K}_A}(\varphi \to \psi)$, and $\boxed{\mathsf{K}_A}\varphi$ hold, actually. Then in each of these worlds $\psi$ holds—which means that $\boxed{\mathsf{K}_A}\psi$. So $\boxed{\mathsf{K}_A}(\varphi \to \psi) \wedge \boxed{\mathsf{K}_A}\varphi \to \boxed{\mathsf{K}_A}\psi$ holds. Equivalently:

$$W \vDash \boxed{\mathsf{K}_A}(\varphi \to \psi) \to \left( \boxed{\mathsf{K}_A}\varphi \to \boxed{\mathsf{K}_A}\psi \right) \qquad \text{(K)}$$

The actual world is viewed as possible by all agents, and whatever any agent knows is, in fact, *i.e. in the actual world*, true. Formally, $\equiv_A$ is reflexive: $W \equiv_A W$.

$$\textbf{So, if } W \vDash \boxed{\mathsf{K}_A}\varphi \textbf{ then we have } W \vDash \varphi.$$
$$W \vDash \boxed{\mathsf{K}_A}\varphi \to \varphi \qquad \text{(T)}$$

Furthermore, from the perspective of any agent, all its possible worlds are equivalent, so whatever it knows in the actual world, it knows in all its possible worlds — which means that, actually, it knows that it knows. Formally, since $\equiv_A$ is an equivalence relation, it is transitive: if $V \equiv_A W$ then, for each $U \equiv_A V$ we have $U \equiv_A W$.

$$\textbf{So, if } W \vDash \boxed{\mathsf{K}_A}\varphi \textbf{ and } V \equiv_A W \textbf{ then}$$
$$\textbf{for each } U \equiv_A V \textbf{ we have } U \vDash \varphi$$
$$\textbf{which means that } V \vDash \boxed{\mathsf{K}_A}\varphi$$
$$\textbf{and so, } W \vDash \boxed{\mathsf{K}_A}\boxed{\mathsf{K}_A}\varphi.$$
$$W \vDash \boxed{\mathsf{K}_A}\varphi \to \boxed{\mathsf{K}_A}\boxed{\mathsf{K}_A}\varphi \qquad \text{(4)}$$

Finally, $\equiv_A$ is symetric: if $V \equiv_A W$ then $W \equiv_A V$.

$$\textbf{If, } W \vDash \neg\boxed{\mathsf{K}_A}\varphi \textbf{ then, for some } V \equiv_A W, \textbf{ we have } V \nvDash \varphi.$$
$$\textbf{Now, for any } U \equiv_A W, \textbf{ we have } V \equiv_A U, \textbf{ so } U \nvDash \boxed{\mathsf{K}_A}\varphi; \textbf{ i.e. } U \vDash \neg\boxed{\mathsf{K}_A}\varphi.$$
$$\textbf{Which means that } W \vDash \boxed{\mathsf{K}_A}\neg\boxed{\mathsf{K}_A}\varphi.$$
$$W \vDash \neg\boxed{\mathsf{K}_A}\varphi \to \boxed{\mathsf{K}_A}\neg\boxed{\mathsf{K}_A}\varphi \qquad \text{(5)}$$

From this discussion, it is clear that if our semantics were defined in terms of an arbitrary relation, in place of the equivalence relation we have used, then different principles could be justified depending on the properties of this relation: **K** would always be valid; **T** if the relation were reflexive; **4** if transitive; **5** if transitive and symmetric. These four, **{K, T, 4, 5}**, corresponding to an equivalence relation, characterise a so-called **S5** modality. The logic for a reflexive relation, satisfying **{K, T}**, is called **KT**. For a reflexive, transitive relation, we have the logic **S4**, satisfying **{K, T, 4}**[1].

---

[1]**K, KT, S4, S5** are the traditional names for modal logics with a single neccessity operator satisfying these rules; the logics discussed here are multimodal logics.

**Common Knowledge**   Sometimes it is not enough for $A$ to know that $B$ knows, and so on; sometimes we need to know that everybody knows and that everybody knows that everybody knows. We introduce special modalities

$$\boxed{E}\,\varphi \;\textbf{everybody knows}\; \text{that}\; \varphi$$
$$\boxed{C}\,\varphi \;\text{it is}\; \textbf{common knowledge}\; \text{that}\; \varphi$$

where

$$\boxed{E}\,\varphi \Leftrightarrow \boxed{K_A}\,\varphi \wedge \boxed{K_B}\,\varphi \wedge \ldots$$
$$\boxed{C}\,\varphi \Leftrightarrow \boxed{E}\,(\varphi \wedge \boxed{C}\,\varphi)$$

Here, *everybody knows* is defined by a conjunction, over all agents. *Common knowledge* is defined recursively.

Semantically, the modality $\boxed{E}$, *everybody knows*, corresponds to the relation which is the union of all the $\equiv_A$. The union of equivalence relations is not, in general, an equivalence relation; it is not transitive. The transitive closure of this relation is an equivalence relation, it provides the model for the *common knowledge* modality, $\boxed{C}$.

We will characterise speech acts as actions that agents can plan to take, using modal logic to express their preconditions and effects.