

```

function KB-AGENT(percept) returns an action
  persistent: KB, a knowledge base
    t, a counter, initially 0, indicating time

  TELL(KB, MAKE-PERCEPT-SENTENCE(percept, t))
  action  $\leftarrow$  ASK(KB, MAKE-ACTION-QUERY(t))
  TELL(KB, MAKE-ACTION-SENTENCE(action, t))
  t  $\leftarrow$  t + 1
  return action

```

**Figure 7.1** A generic knowledge-based agent. Given a percept, the agent adds the percept to its knowledge base, asks the knowledge base for the best action, and tells the knowledge base that it has in fact taken that action.

KNOWLEDGE LEVEL

the **knowledge level**, where we need specify only what the agent knows and what its goals are, in order to fix its behavior. For example, an automated taxi might have the goal of taking a passenger from San Francisco to Marin County and might know that the Golden Gate Bridge is the only link between the two locations. Then we can expect it to cross the Golden Gate Bridge *because it knows that that will achieve its goal*. Notice that this analysis is independent of how the taxi works at the **implementation level**. It doesn't matter whether its geographical knowledge is implemented as linked lists or pixel maps, or whether it reasons by manipulating strings of symbols stored in registers or by propagating noisy signals in a network of neurons.

IMPLEMENTATION LEVEL

DECLARATIVE

A knowledge-based agent can be built simply by TELLING it what it needs to know. Starting with an empty knowledge base, the agent designer can TELL sentences one by one until the agent knows how to operate in its environment. This is called the **declarative** approach to system building. In contrast, the **procedural** approach encodes desired behaviors directly as program code. In the 1970s and 1980s, advocates of the two approaches engaged in heated debates. We now understand that a successful agent often combines both declarative and procedural elements in its design, and that declarative knowledge can often be compiled into more efficient procedural code.

We can also provide a knowledge-based agent with mechanisms that allow it to learn for itself. These mechanisms, which are discussed in Chapter 18, create general knowledge about the environment from a series of percepts. A learning agent can be fully autonomous.

## 7.2 THE WUMPUS WORLD

WUMPUS WORLD

In this section we describe an environment in which knowledge-based agents can show their worth. The **wumpus world** is a cave consisting of rooms connected by passageways. Lurking somewhere in the cave is the terrible wumpus, a beast that eats anyone who enters its room. The wumpus can be shot by an agent, but the agent has only one arrow. Some rooms contain

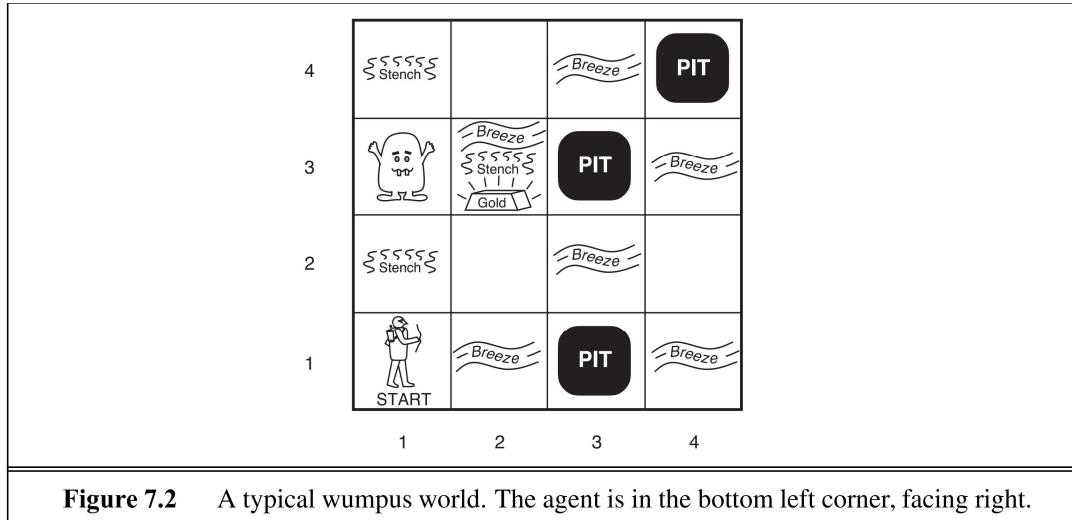
bottomless pits that will trap anyone who wanders into these rooms (except for the wumpus, which is too big to fall in). The only mitigating feature of this bleak environment is the possibility of finding a heap of gold. Although the wumpus world is rather tame by modern computer game standards, it illustrates some important points about intelligence.

A sample wumpus world is shown in Figure 7.2. The precise definition of the task environment is given, as suggested in Section 2.3, by the PEAS description:

- **Performance measure:** +1000 for climbing out of the cave with the gold, -1000 for falling into a pit or being eaten by the wumpus, -1 for each action taken and -10 for using up the arrow. The game ends either when the agent dies or when the agent climbs out of the cave.
- **Environment:** A  $4 \times 4$  grid of rooms. The agent always starts in the square labeled [1,1], facing to the right. The locations of the gold and the wumpus are chosen randomly, with a uniform distribution, from the squares other than the start square. In addition, each square other than the start can be a pit, with probability 0.2.
- **Actuators:** The agent can move *Forward*, *TurnLeft* by  $90^\circ$ , or *TurnRight* by  $90^\circ$ . The agent dies a miserable death if it enters a square containing a pit or a live wumpus. (It is safe, albeit smelly, to enter a square with a dead wumpus.) If an agent tries to move forward and bumps into a wall, then the agent does not move. The action *Grab* can be used to pick up the gold if it is in the same square as the agent. The action *Shoot* can be used to fire an arrow in a straight line in the direction the agent is facing. The arrow continues until it either hits (and hence kills) the wumpus or hits a wall. The agent has only one arrow, so only the first *Shoot* action has any effect. Finally, the action *Climb* can be used to climb out of the cave, but only from square [1,1].
- **Sensors:** The agent has five sensors, each of which gives a single bit of information:
  - In the square containing the wumpus and in the directly (not diagonally) adjacent squares, the agent will perceive a *Stench*.
  - In the squares directly adjacent to a pit, the agent will perceive a *Breeze*.
  - In the square where the gold is, the agent will perceive a *Glitter*.
  - When an agent walks into a wall, it will perceive a *Bump*.
  - When the wumpus is killed, it emits a woeful *Scream* that can be perceived anywhere in the cave.

The percepts will be given to the agent program in the form of a list of five symbols; for example, if there is a stench and a breeze, but no glitter, bump, or scream, the agent program will get [*Stench*, *Breeze*, *None*, *None*, *None*].

We can characterize the wumpus environment along the various dimensions given in Chapter 2. Clearly, it is discrete, static, and single-agent. (The wumpus doesn't move, fortunately.) It is sequential, because rewards may come only after many actions are taken. It is partially observable, because some aspects of the state are not directly perceivable: the agent's location, the wumpus's state of health, and the availability of an arrow. As for the locations of the pits and the wumpus: we could treat them as unobserved parts of the state that happen to be immutable—in which case, the transition model for the environment is completely



**Figure 7.2** A typical wumpus world. The agent is in the bottom left corner, facing right.

known; or we could say that the transition model itself is unknown because the agent doesn't know which *Forward* actions are fatal—in which case, discovering the locations of pits and wumpus completes the agent's knowledge of the transition model.

For an agent in the environment, the main challenge is its initial ignorance of the configuration of the environment; overcoming this ignorance seems to require logical reasoning. In most instances of the wumpus world, it is possible for the agent to retrieve the gold safely. Occasionally, the agent must choose between going home empty-handed and risking death to find the gold. About 21% of the environments are utterly unfair, because the gold is in a pit or surrounded by pits.

Let us watch a knowledge-based wumpus agent exploring the environment shown in Figure 7.2. We use an informal knowledge representation language consisting of writing down symbols in a grid (as in Figures 7.3 and 7.4).

The agent's initial knowledge base contains the rules of the environment, as described previously; in particular, it knows that it is in [1,1] and that [1,1] is a safe square; we denote that with an "A" and "OK," respectively, in square [1,1].

The first percept is *[None, None, None, None, None]*, from which the agent can conclude that its neighboring squares, [1,2] and [2,1], are free of dangers—they are OK. Figure 7.3(a) shows the agent's state of knowledge at this point.

A cautious agent will move only into a square that it knows to be OK. Let us suppose the agent decides to move forward to [2,1]. The agent perceives a breeze (denoted by "B") in [2,1], so there must be a pit in a neighboring square. The pit cannot be in [1,1], by the rules of the game, so there must be a pit in [2,2] or [3,1] or both. The notation "P?" in Figure 7.3(b) indicates a possible pit in those squares. At this point, there is only one known square that is OK and that has not yet been visited. So the prudent agent will turn around, go back to [1,1], and then proceed to [1,2].

The agent perceives a stench in [1,2], resulting in the state of knowledge shown in Figure 7.4(a). The stench in [1,2] means that there must be a wumpus nearby. But the

<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>1,4</td><td>2,4</td><td>3,4</td><td>4,4</td></tr> <tr><td>1,3</td><td>2,3</td><td>3,3</td><td>4,3</td></tr> <tr><td>1,2</td><td>2,2</td><td>3,2</td><td>4,2</td></tr> <tr><td style="text-align: center;">OK</td><td></td><td></td><td></td></tr> <tr><td>1,1</td><td style="text-align: center;"><span style="border: 1px solid black; padding: 2px;">A</span></td><td>2,1</td><td>3,1</td></tr> <tr><td></td><td style="text-align: center;">OK</td><td style="text-align: center;">OK</td><td></td></tr> </table>	1,4	2,4	3,4	4,4	1,3	2,3	3,3	4,3	1,2	2,2	3,2	4,2	OK				1,1	<span style="border: 1px solid black; padding: 2px;">A</span>	2,1	3,1		OK	OK		<span style="border: 1px solid black; padding: 2px;">A</span> = Agent B = Breeze G = Glitter, Gold OK = Safe square P = Pit S = Stench V = Visited W = Wumpus	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>1,4</td><td>2,4</td><td>3,4</td><td>4,4</td></tr> <tr><td>1,3</td><td>2,3</td><td>3,3</td><td>4,3</td></tr> <tr><td>1,2</td><td>2,2</td><td style="text-align: center;">P?</td><td>4,2</td></tr> <tr><td style="text-align: center;">OK</td><td></td><td></td><td></td></tr> <tr><td>1,1</td><td style="text-align: center;">V</td><td style="text-align: center;"><span style="border: 1px solid black; padding: 2px;">A</span></td><td>4,1</td></tr> <tr><td></td><td style="text-align: center;">OK</td><td style="text-align: center;">B</td><td style="text-align: center;">OK</td></tr> </table>	1,4	2,4	3,4	4,4	1,3	2,3	3,3	4,3	1,2	2,2	P?	4,2	OK				1,1	V	<span style="border: 1px solid black; padding: 2px;">A</span>	4,1		OK	B	OK
1,4	2,4	3,4	4,4																																															
1,3	2,3	3,3	4,3																																															
1,2	2,2	3,2	4,2																																															
OK																																																		
1,1	<span style="border: 1px solid black; padding: 2px;">A</span>	2,1	3,1																																															
	OK	OK																																																
1,4	2,4	3,4	4,4																																															
1,3	2,3	3,3	4,3																																															
1,2	2,2	P?	4,2																																															
OK																																																		
1,1	V	<span style="border: 1px solid black; padding: 2px;">A</span>	4,1																																															
	OK	B	OK																																															
(a)		(b)																																																

**Figure 7.3** The first step taken by the agent in the wumpus world. (a) The initial situation, after percept [None, None, None, None, None]. (b) After one move, with percept [None, Breeze, None, None, None].

<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>1,4</td><td>2,4</td><td>3,4</td><td>4,4</td></tr> <tr><td>1,3</td><td>2,3</td><td>3,3</td><td>4,3</td></tr> <tr><td>1,2</td><td style="text-align: center;"><span style="border: 1px solid black; padding: 2px;">A</span></td><td>2,2</td><td>4,2</td></tr> <tr><td style="text-align: center;">S OK</td><td style="text-align: center;">OK</td><td></td><td></td></tr> <tr><td>1,1</td><td style="text-align: center;">V</td><td style="text-align: center;">B</td><td>4,1</td></tr> <tr><td></td><td style="text-align: center;">OK</td><td style="text-align: center;">V</td><td style="text-align: center;">OK</td></tr> </table>	1,4	2,4	3,4	4,4	1,3	2,3	3,3	4,3	1,2	<span style="border: 1px solid black; padding: 2px;">A</span>	2,2	4,2	S OK	OK			1,1	V	B	4,1		OK	V	OK	<span style="border: 1px solid black; padding: 2px;">A</span> = Agent B = Breeze G = Glitter, Gold OK = Safe square P = Pit S = Stench V = Visited W = Wumpus	<table border="1" style="width: 100%; border-collapse: collapse;"> <tr><td>1,4</td><td>2,4</td><td>3,4</td><td>4,4</td></tr> <tr><td>1,3</td><td style="text-align: center;">W!</td><td>2,3</td><td>4,3</td></tr> <tr><td>1,2</td><td style="text-align: center;">S</td><td style="text-align: center;"><span style="border: 1px solid black; padding: 2px;">A</span></td><td>4,2</td></tr> <tr><td style="text-align: center;">V OK</td><td style="text-align: center;">V</td><td style="text-align: center;">S</td><td style="text-align: center;">OK</td></tr> <tr><td>1,1</td><td style="text-align: center;">V</td><td style="text-align: center;">B</td><td>4,1</td></tr> <tr><td></td><td style="text-align: center;">OK</td><td style="text-align: center;">V</td><td style="text-align: center;">OK</td></tr> </table>	1,4	2,4	3,4	4,4	1,3	W!	2,3	4,3	1,2	S	<span style="border: 1px solid black; padding: 2px;">A</span>	4,2	V OK	V	S	OK	1,1	V	B	4,1		OK	V	OK
1,4	2,4	3,4	4,4																																															
1,3	2,3	3,3	4,3																																															
1,2	<span style="border: 1px solid black; padding: 2px;">A</span>	2,2	4,2																																															
S OK	OK																																																	
1,1	V	B	4,1																																															
	OK	V	OK																																															
1,4	2,4	3,4	4,4																																															
1,3	W!	2,3	4,3																																															
1,2	S	<span style="border: 1px solid black; padding: 2px;">A</span>	4,2																																															
V OK	V	S	OK																																															
1,1	V	B	4,1																																															
	OK	V	OK																																															
(a)		(b)																																																

**Figure 7.4** Two later stages in the progress of the agent. (a) After the third move, with percept [Stench, None, None, None, None]. (b) After the fifth move, with percept [Stench, Breeze, Glitter, None, None].

wumpus cannot be in [1,1], by the rules of the game, and it cannot be in [2,2] (or the agent would have detected a stench when it was in [2,1]). Therefore, the agent can infer that the wumpus is in [1,3]. The notation W! indicates this inference. Moreover, the lack of a breeze in [1,2] implies that there is no pit in [2,2]. Yet the agent has already inferred that there must be a pit in either [2,2] or [3,1], so this means it must be in [3,1]. This is a fairly difficult inference, because it combines knowledge gained at different times in different places and relies on the lack of a percept to make one crucial step.

The agent has now proved to itself that there is neither a pit nor a wumpus in [2,2], so it is OK to move there. We do not show the agent’s state of knowledge at [2,2]; we just assume that the agent turns and moves to [2,3], giving us Figure 7.4(b). In [2,3], the agent detects a glitter, so it should grab the gold and then return home.

Note that in each case for which the agent draws a conclusion from the available information, that conclusion is *guaranteed* to be correct if the available information is correct. This is a fundamental property of logical reasoning. In the rest of this chapter, we describe how to build logical agents that can represent information and draw conclusions such as those described in the preceding paragraphs.

## 7.3 LOGIC

---

This section summarizes the fundamental concepts of logical representation and reasoning. These beautiful ideas are independent of any of logic’s particular forms. We therefore postpone the technical details of those forms until the next section, using instead the familiar example of ordinary arithmetic.

**SYNTAX** In Section 7.1, we said that knowledge bases consist of sentences. These sentences are expressed according to the **syntax** of the representation language, which specifies all the sentences that are well formed. The notion of syntax is clear enough in ordinary arithmetic: “ $x + y = 4$ ” is a well-formed sentence, whereas “ $x4y+ =$ ” is not.

**SEMANTICS** A logic must also define the **semantics** or meaning of sentences. The semantics defines the **truth** of each sentence with respect to each **possible world**. For example, the semantics for arithmetic specifies that the sentence “ $x + y = 4$ ” is true in a world where  $x$  is 2 and  $y$  is 2, but false in a world where  $x$  is 1 and  $y$  is 1. In standard logics, every sentence must be either true or false in each possible world—there is no “in between.”<sup>1</sup>

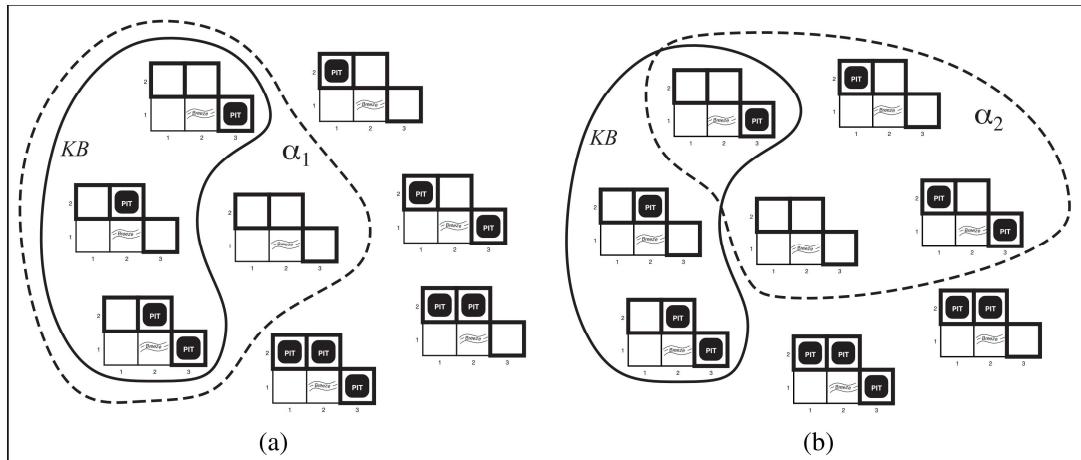
**TRUTH** **POSSIBLE WORLD** **MODEL** When we need to be precise, we use the term **model** in place of “possible world.” Whereas possible worlds might be thought of as (potentially) real environments that the agent might or might not be in, models are mathematical abstractions, each of which simply fixes the truth or falsehood of every relevant sentence. Informally, we may think of a possible world as, for example, having  $x$  men and  $y$  women sitting at a table playing bridge, and the sentence  $x + y = 4$  is true when there are four people in total. Formally, the possible models are just all possible assignments of real numbers to the variables  $x$  and  $y$ . Each such assignment fixes the truth of any sentence of arithmetic whose variables are  $x$  and  $y$ . If a sentence  $\alpha$  is true in model  $m$ , we say that  $m$  **satisfies**  $\alpha$  or sometimes  $m$  is a **model of**  $\alpha$ . We use the notation  $M(\alpha)$  to mean the set of all models of  $\alpha$ .

**SATISFACTION** **ENTAILMENT** Now that we have a notion of truth, we are ready to talk about logical reasoning. This involves the relation of logical **entailment** between sentences—the idea that a sentence *follows logically* from another sentence. In mathematical notation, we write

$$\alpha \models \beta$$

---

<sup>1</sup> **Fuzzy logic**, discussed in Chapter 14, allows for degrees of truth.



**Figure 7.5** Possible models for the presence of pits in squares [1,2], [2,2], and [3,1]. The KB corresponding to the observations of nothing in [1,1] and a breeze in [2,1] is shown by the solid line. (a) Dotted line shows models of  $\alpha_1$  (no pit in [1,2]). (b) Dotted line shows models of  $\alpha_2$  (no pit in [2,2]).

to mean that the sentence  $\alpha$  entails the sentence  $\beta$ . The formal definition of entailment is this:  $\alpha \models \beta$  if and only if, in every model in which  $\alpha$  is true,  $\beta$  is also true. Using the notation just introduced, we can write

$$\alpha \models \beta \text{ if and only if } M(\alpha) \subseteq M(\beta).$$

(Note the direction of the  $\subseteq$  here: if  $\alpha \models \beta$ , then  $\alpha$  is a *stronger* assertion than  $\beta$ : it rules out *more* possible worlds.) The relation of entailment is familiar from arithmetic; we are happy with the idea that the sentence  $x = 0$  entails the sentence  $xy = 0$ . Obviously, in any model where  $x$  is zero, it is the case that  $xy$  is zero (regardless of the value of  $y$ ).

We can apply the same kind of analysis to the wumpus-world reasoning example given in the preceding section. Consider the situation in Figure 7.3(b): the agent has detected nothing in [1,1] and a breeze in [2,1]. These percepts, combined with the agent's knowledge of the rules of the wumpus world, constitute the KB. The agent is interested (among other things) in whether the adjacent squares [1,2], [2,2], and [3,1] contain pits. Each of the three squares might or might not contain a pit, so (for the purposes of this example) there are  $2^3 = 8$  possible models. These eight models are shown in Figure 7.5.<sup>2</sup>

The KB can be thought of as a set of sentences or as a single sentence that asserts all the individual sentences. The KB is false in models that contradict what the agent knows—for example, the KB is false in any model in which [1,2] contains a pit, because there is no breeze in [1,1]. There are in fact just three models in which the KB is true, and these are

<sup>2</sup> Although the figure shows the models as partial wumpus worlds, they are really nothing more than assignments of *true* and *false* to the sentences “there is a pit in [1,2]” etc. Models, in the mathematical sense, do not need to have ‘orrible ‘airy wumpuses in them.

shown surrounded by a solid line in Figure 7.5. Now let us consider two possible conclusions:

- $\alpha_1$  = “There is no pit in [1,2].”
- $\alpha_2$  = “There is no pit in [2,2].”

We have surrounded the models of  $\alpha_1$  and  $\alpha_2$  with dotted lines in Figures 7.5(a) and 7.5(b), respectively. By inspection, we see the following:

in every model in which  $KB$  is true,  $\alpha_1$  is also true.

Hence,  $KB \models \alpha_1$ : there is no pit in [1,2]. We can also see that

in some models in which  $KB$  is true,  $\alpha_2$  is false.

Hence,  $KB \not\models \alpha_2$ : the agent *cannot* conclude that there is no pit in [2,2]. (Nor can it conclude that there *is* a pit in [2,2].)<sup>3</sup>

The preceding example not only illustrates entailment but also shows how the definition of entailment can be applied to derive conclusions—that is, to carry out **logical inference**. The inference algorithm illustrated in Figure 7.5 is called **model checking**, because it enumerates all possible models to check that  $\alpha$  is true in all models in which  $KB$  is true, that is, that  $M(KB) \subseteq M(\alpha)$ .

In understanding entailment and inference, it might help to think of the set of all consequences of  $KB$  as a haystack and of  $\alpha$  as a needle. Entailment is like the needle being in the haystack; inference is like finding it. This distinction is embodied in some formal notation: if an inference algorithm  $i$  can derive  $\alpha$  from  $KB$ , we write

$$KB \vdash_i \alpha,$$

which is pronounced “ $\alpha$  is derived from  $KB$  by  $i$ ” or “ $i$  derives  $\alpha$  from  $KB$ .”

An inference algorithm that derives only entailed sentences is called **sound** or **truth-preserving**. Soundness is a highly desirable property. An unsound inference procedure essentially makes things up as it goes along—it announces the discovery of nonexistent needles. It is easy to see that model checking, when it is applicable,<sup>4</sup> is a sound procedure.

The property of **completeness** is also desirable: an inference algorithm is complete if it can derive any sentence that is entailed. For real haystacks, which are finite in extent, it seems obvious that a systematic examination can always decide whether the needle is in the haystack. For many knowledge bases, however, the haystack of consequences is infinite, and completeness becomes an important issue.<sup>5</sup> Fortunately, there are complete inference procedures for logics that are sufficiently expressive to handle many knowledge bases.

We have described a reasoning process whose conclusions are guaranteed to be true in any world in which the premises are true; in particular, *if KB is true in the real world, then any sentence  $\alpha$  derived from KB by a sound inference procedure is also true in the real world*. So, while an inference process operates on “syntax”—internal physical configurations such as bits in registers or patterns of electrical blips in brains—the process *corresponds*

<sup>3</sup> The agent can calculate the *probability* that there is a pit in [2,2]; Chapter 13 shows how.

<sup>4</sup> Model checking works if the space of models is finite—for example, in wumpus worlds of fixed size. For arithmetic, on the other hand, the space of models is infinite: even if we restrict ourselves to the integers, there are infinitely many pairs of values for  $x$  and  $y$  in the sentence  $x + y = 4$ .

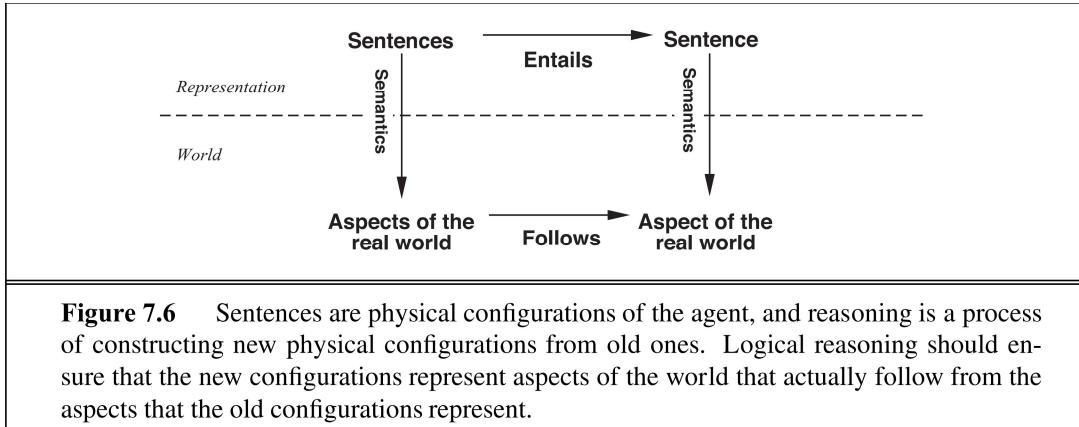
<sup>5</sup> Compare with the case of infinite search spaces in Chapter 3, where depth-first search is not complete.

LOGICAL INFERENCE  
MODEL CHECKING

SOUND  
TRUTH-PRESERVING

COMPLETENESS





to the real-world relationship whereby some aspect of the real world is the case<sup>6</sup> by virtue of other aspects of the real world being the case. This correspondence between world and representation is illustrated in Figure 7.6.



The final issue to consider is **grounding**—the connection between logical reasoning processes and the real environment in which the agent exists. In particular, *how do we know that KB is true in the real world?* (After all, *KB* is just “syntax” inside the agent’s head.) This is a philosophical question about which many, many books have been written. (See Chapter 26.) A simple answer is that the agent’s sensors create the connection. For example, our wumpus-world agent has a smell sensor. The agent program creates a suitable sentence whenever there is a smell. Then, whenever that sentence is in the knowledge base, it is true in the real world. Thus, the meaning and truth of percept sentences are defined by the processes of sensing and sentence construction that produce them. What about the rest of the agent’s knowledge, such as its belief that wumpuses cause smells in adjacent squares? This is not a direct representation of a single percept, but a general rule—derived, perhaps, from perceptual experience but not identical to a statement of that experience. General rules like this are produced by a sentence construction process called **learning**, which is the subject of Part V. Learning is fallible. It could be the case that wumpuses cause smells *except on February 29 in leap years*, which is when they take their baths. Thus, *KB* may not be true in the real world, but with good learning procedures, there is reason for optimism.

## 7.4 PROPOSITIONAL LOGIC: A VERY SIMPLE LOGIC

### PROPOSITIONAL LOGIC

We now present a simple but powerful logic called **propositional logic**. We cover the syntax of propositional logic and its semantics—the way in which the truth of sentences is determined. Then we look at **entailment**—the relation between a sentence and another sentence that follows from it—and see how this leads to a simple algorithm for logical inference. Everything takes place, of course, in the wumpus world.

<sup>6</sup> As Wittgenstein (1922) put it in his famous *Tractatus*: “The world is everything that is the case.”