

# Wumpus World Problem

Reference: Logical Agents (Chapter 7)  
Norvig and Russell

This chapter introduces knowledge-based agents. The concepts that we discuss—the *representation* of knowledge and the *reasoning* processes that bring knowledge to life—are central to the entire field of artificial intelligence.

Humans, it seems, know things and do reasoning. Knowledge and reasoning are also important for artificial agents because they enable successful behaviors that would be very hard to achieve otherwise. We have seen that knowledge of action outcomes enables problem-solving agents to perform well in complex environments. A reflex agents could only find its way from Arad to Bucharest by dumb luck. The knowledge of problem-solving agents is, however, very specific and inflexible. A chess program can calculate the legal moves of its king, but does not know in any useful sense that no piece can be on two different squares at the same time. Knowledge-based agents can benefit from knowledge expressed in very general forms, combining and recombining information to suit myriad purposes. Often, this process can be quite far removed from the needs of the moment—as when a mathematician proves a theorem or an astronomer calculates the earth’s life expectancy.

Knowledge and reasoning also play a crucial role in dealing with *partially observable* environments. A knowledge-based agent can combine general knowledge with current percepts to infer hidden aspects of the current state prior to selecting actions. For example, a physician diagnoses a patient—that is, infers a disease state that is not directly observable—prior to choosing a treatment. Some of the knowledge that the physician uses in the form of rules learned from textbooks and teachers, and some is in the form of patterns of association that the physician may not be able to consciously describe. If its inside the physician’s head, it counts as knowledge.

Understanding natural language also requires inferring hidden state, namely, the intention of the speaker. When we hear, “John saw the diamond through the window and coveted it,” we know “it” refers to the diamond and not the window—we reason, perhaps unconsciously, with our knowledge of relative value. Similarly, when we hear, “John threw the brick through the window and broke it,” we know “it” refers to the window. Reasoning allows

us to cope with the virtually infinite variety of utterances using a finite store of commonsense knowledge. Problem-solving agents have difficulty with this kind of ambiguity because their representation of contingency problems is inherently exponential.

Our final reason for studying knowledge-based agents is their flexibility. They are able to accept new tasks in the form of explicitly described goals, they can achieve competence quickly by being told or learning new knowledge about the environment, and they can adapt to changes in the environment by updating the relevant knowledge.

We begin in Section 7.1 with the overall agent design. Section 7.2 introduces a simple new environment, the wumpus world, and illustrates the operation of a knowledge-based agent without going into any technical detail. Then, in Section 7.3, we explain the general principles of **logic**. Logic will be the primary vehicle for representing knowledge throughout Part III of the book. The knowledge of logical agents is always *definite*—each proposition is either true or false in the world, although the agent may be agnostic about some propositions.

Logic has the pedagogical advantage of being simple example of a representation for knowledge-based agents, but logic has some severe limitations. Clearly, a large portion of the reasoning carried out by humans and other agents in partially observable environments depends on handling knowledge that is *uncertain*. Logic cannot represent this uncertainty well, so in Part V we cover probability, which can. In Part VI and Part VII we cover many representations, including some based on continuous mathematics such as mixtures of Gaussians, neural networks, and other representations.

Section 7.4 of this chapter defines a simple logic called **propositional logic**. While much less expressive than **first-order logic** (Chapter 8), propositional logic serves to illustrate all the basic concepts of logic. There is also a well-developed technology for reasoning in propositional logic, which we describe in sections 7.5 and 7.6. Finally, Section 7.7 combines the concept of logical agents with the technology of propositional logic to build some simple agents for the wumpus world. Certain shortcomings in propositional logic are identified, motivating the development of more powerful logics in subsequent chapters.

## 7.1 KNOWLEDGE-BASED AGENTS

---

|                                   |   |
|-----------------------------------|---|
| KNOWLEDGE BASE                    | The central component of a knowledge-based agent is its <b>knowledge base</b> , or KB. Informally, a knowledge base is a set of <b>sentences</b> . (Here “sentence” is used as a technical term. It is related but is not identical to the sentences of English and other natural languages.) Each sentence is expressed in a language called a <b>knowledge representation language</b> and represents some assertion about the world.   |
| SENTENCE                          |   |
| KNOWLEDGE REPRESENTATION LANGUAGE |   |
| INFERENCE                         | There must be a way to add new sentences to the knowledge base and a way to query what is known. The standard names for these tasks are TELL and ASK, respectively. Both tasks may involve <b>inference</b> —that is, deriving new sentences from old. In <b>logical agents</b> , which are the main subject of study in this chapter, inference must obey the fundamental requirement that when one ASKs a question of the knowledge base, the answer should follow from what has been told (or rather, TELLED) to the knowledge base previously. Later in the |
| LOGICAL AGENTS                    |   |

```

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

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

```

**Figure 7.1** A generic knowledge-based agent.

chapter, we will be more precise about the crucial word “follow.” For now, take it to mean that the inference process should not just make things up as it goes along.

Figure 7.1 shows the outline of a knowledge-based agent program. Like all our agents, it takes a percept as input and returns an action. The agent maintains a knowledge base, *KB*, which may initially contain some **background knowledge**. Each time the agent program is called, it does two things. First, it TELLS the knowledge base what it perceives. Second, it ASKS the knowledge base what action it should perform. In the process of answering this query, extensive reasoning may be done about the current state of the world, about the outcomes of possible action sequences, and so on. Once the action is chosen, the agent records its choice with TELL and executes the action. The second TELL is necessary to let the knowledge base know that the hypothetical *action* has actually been executed.

The details of the representation language are hidden inside two functions that implement the interface between the sensors and actuators and the core representation and reasoning system. MAKE-PERCEPT-SENTENCE takes a percept and a time and returns a sentence asserting that the agent perceived the percept at the given time. MAKE-ACTION-QUERY takes a time as input and returns a sentence that asks what action should be performed at that time. The details of the inference mechanisms are hidden inside TELL and ASK. Later sections will reveal these details.

The agent in Figure 7.1 appears quite similar to the agents with internal state described in Chapter 2. Because of the definitions of TELL and ASK, however, the knowledge-based agent is not an arbitrary program for calculating actions. It is amenable to a description at 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 delivering a passenger to Marin County and might know that it is in San Francisco and 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.

BACKGROUND  
KNOWLEDGE

KNOWLEDGE LEVEL

IMPLEMENTATION  
LEVEL



DECLARATIVE

As we mentioned in the introduction to the chapter, *one can build a knowledge-based agent simply by TELLing it what it needs to know*. The agent's initial program, before it starts to receive percepts, is built by adding one by one the sentences that represent the designer's knowledge of the environment. Designing the representation language to make it easy to express this knowledge in the form of sentences simplifies the construction problem enormously. This is called the **declarative** approach to system building. In contrast, the **procedural** approach encodes desired behaviors directly as program code; minimizing the role of explicit representation and reasoning can result in a much more efficient system. We will see agents of both kinds in Section 7.7. In the 1970s and 1980s, advocates of the two approaches engaged in heated debates. We now understand that a successful agent must combine both declarative and procedural elements in its design.

In addition to TELLing it what it needs to know, we can 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 out of a series of percepts. This knowledge can be incorporated into the agent's knowledge base and used for decision making. In this way, the agent can be fully autonomous.

All these capabilities—representation, reasoning, and learning—rest on the centuries-long development of the theory and technology of logic. Before explaining that theory and technology, however, we will create a simple world with which to illustrate them.

## 7.2 THE WUMPUS WORLD

---

WUMPUS WORLD

The **wumpus world** is a cave consisting of rooms connected by passageways. Lurking somewhere in the cave is the 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 living in this environment is the possibility of finding a heap of gold. Although the wumpus world is rather tame by modern computer game standards, it makes an excellent testbed environment for intelligent agents. Michael Genesereth was the first to suggest this.

A sample wumpus world is shown in Figure 7.2. The precise definition of the task environment is given, as suggested in Chapter 2, by the PEAS description:

- ◇ **Performance measure:** +1000 for picking up 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.
- ◇ **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, turn left by  $90^\circ$ , or turn right 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.) Moving forward has no

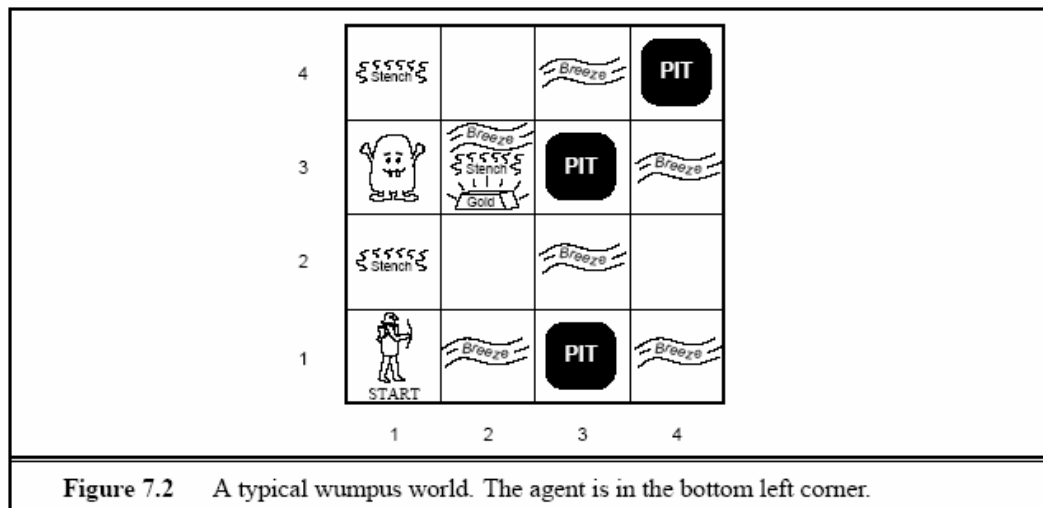
effect if there is a wall in front of the agent. The action *Grab* can be used to pick up an object that is in the same square as the agent. The action *Shoot* can be used to fire an arrow in a straight line in the direction the agent is facing. The arrow continues until it either hits (and hence kills) the wumpus or hits a wall. The agent only has one arrow, so only the first *Shoot* action has any effect.

- ◇ **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 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 will receive the percept [*Stench*, *Breeze*, *None*, *None*, *None*].

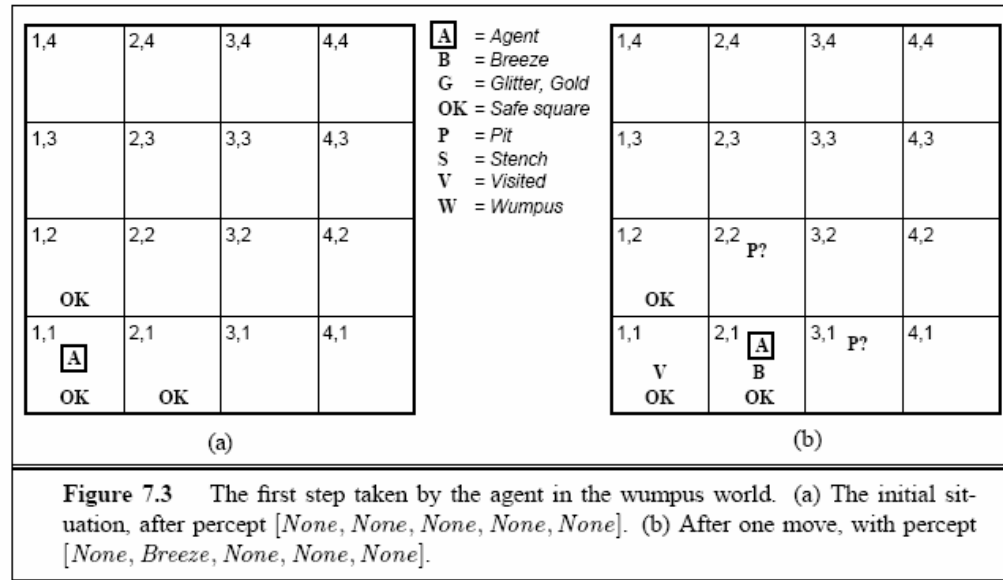
Exercise 7.1 asks you to define the wumpus environment along the various dimensions given in Chapter 2. The principal difficulty for the agent is its initial ignorance of the configuration of the environment; overcoming this ignorance seems to require logical reasoning. In most instances of the wumpus 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. The agent's initial knowledge base contains the rules of the environment, as listed



previously; in particular, it knows that it is in [1,1] and that [1,1] is a safe square. We will see how its knowledge evolves as new percepts arrive and actions are taken.

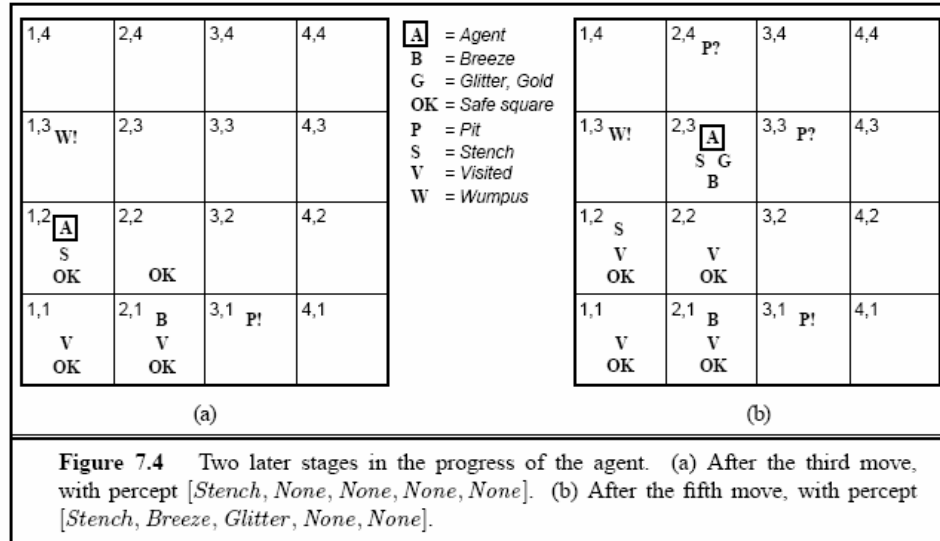
The first percept is  $[None, None, None, None, None]$ , from which the agent can conclude that its neighboring squares are safe. Figure 7.3(a) shows the agent's state of knowledge at this point. We list (some of) the sentences in the knowledge base using letters such as  $B$  (breezy) and  $OK$  (safe, neither pit nor wumpus) marked in the appropriate squares. Figure 7.2, on the other hand, depicts the world itself.



From the fact that there was no stench or breeze in [1,1], the agent can infer that [1,2] and [2,1] are free of dangers. They are marked with an  $OK$  to indicate this. A cautious agent will move only into a square that it knows is  $OK$ . Let us suppose the agent decides to move forward to [2,1], giving the scene in Figure 7.3(b).

The agent detects a breeze 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 has not been visited yet. So the prudent agent will turn around, go back to [1,1], and then proceed to [1,2].

The new percept in [1,2] is  $[Stench, None, None, None, None]$ , 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 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. Moreover, the lack of a  $Breeze$  in [1,2] implies that there is no pit in [2,2]. Yet we 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 inference is beyond the abilities of most animals, but it is typical of the kind of reasoning that a logical agent does.

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 will 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 thereby end the game.



*In each case where 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 the necessary information and draw the conclusions that were described in the preceding paragraphs.

### 7.3 LOGIC

This section provides an overview of all the fundamental concepts of logical representation and reasoning. We postpone the technical details of any particular form of logic until the next section. We will instead use informal examples from the wumpus world and from the familiar realm of arithmetic. We adopt this rather unusual approach because the ideas of logic are far more general and beautiful than is commonly supposed.

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 " $x2y+ =$ " is not. The syntax of logical

SYNTAX



languages (and of arithmetic, for that matter) is usually designed for writing papers and books. There are literally dozens of different syntaxes, some with lots of Greek letters and exotic mathematical symbols, and some with rather visually appealing diagrams with arrows and bubbles. In all cases, however, sentences in an agent's knowledge base are real physical configurations of (parts of) the agent. Reasoning will involve generating and manipulating those configurations.

SEMANTICS

A logic must also define the **semantics** of the language. Loosely speaking, semantics has to do with the “meaning” of sentences. In logic, the definition is more precise. The semantics of the language defines the **truth** of each sentence with respect to each **possible world**. For example, the usual semantics adopted 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.<sup>1</sup> In standard logics, every sentence must be either true or false in each possible world—there is no “in between.”<sup>2</sup>

TRUTH

POSSIBLE WORLD

MODEL

When we need to be precise, we will use the term **model** in place of “possible world.” (We will also use the phrase “ $m$  is a model of  $\alpha$ ” to mean that sentence  $\alpha$  is true in model  $m$ .) 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  $x$  and  $y$  as the number of men and women sitting at a table playing bridge, for example, and the sentence  $x + y = 4$  is true when there are four in total; formally, the possible models are just all possible assignments of numbers to the variables  $x$  and  $y$ . Each such assignment fixes the truth of any sentence of arithmetic whose variables are  $x$  and  $y$ .

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 as

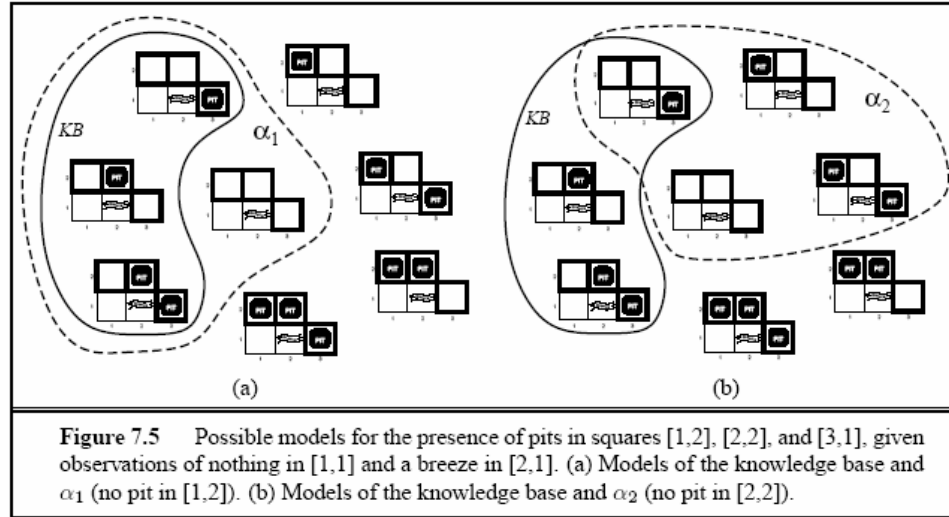
$$\alpha \models \beta$$

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. Another way to say this is that if  $\alpha$  is true, then  $\beta$  *must* also be true. Informally, the truth of  $\beta$  is “contained” in the truth of  $\alpha$ . The relation of entailment is familiar from arithmetic; we are happy with the idea that the sentence  $x + y = 4$  entails the sentence  $4 = x + y$ . Obviously, in any model where  $x + y = 4$ —such as the model in which  $x$  is 2 and  $y$  is 2—it is the case that  $4 = x + y$ . We will see shortly that a knowledge base can be considered a statement, and we often talk of a knowledge base entailing a sentence.

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 (the PEAS description on page 197), constitute the KB. The

<sup>1</sup> The reader will no doubt have noticed the similarity between the notion of truth of sentences and the notion of satisfaction of constraints in Chapter 5. This is no accident—constraint languages are indeed logics and constraint solving is a form of logical reasoning.

<sup>2</sup> Fuzzy logic, discussed in Chapter 14, allows for degrees of truth.



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 are shown in Figure 7.5.<sup>3</sup>

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 shown as a subset of the models 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 marked the models of  $\alpha_1$  and  $\alpha_2$  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>4</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.

LOGICAL INFERENCE  
MODEL CHECKING

<sup>3</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.

<sup>4</sup> The agent can calculate the *probability* that there is a pit in [2,2]; Chapter 13 shows how.

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  $KB$ .”

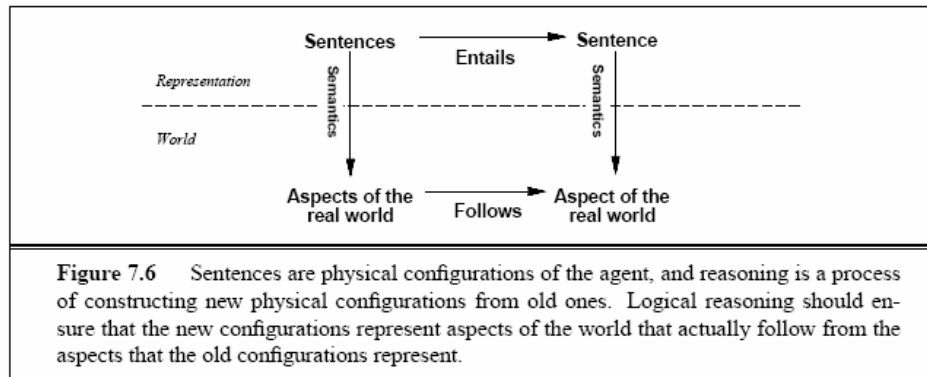
SOUND  
TRUTH-PRESERVING  
  
COMPLETENESS

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>5</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>6</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* to the real-world relationship whereby some aspect of the real world is the case<sup>7</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.



<sup>5</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>6</sup> Compare with the case of infinite search spaces in Chapter 3, where depth-first search is not complete.

<sup>7</sup> As Wittgenstein (1922) put it in his famous *Tractatus*: “The world is everything that is the case.”



The final issue that must be addressed by an account of logical agents is that of **grounding**—the connection, if any, 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 VI. 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 very simple logic called **propositional logic**.<sup>8</sup> 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.

### Syntax

#### ATOMIC SENTENCES PROPOSITION SYMBOL

The **syntax** of propositional logic defines the allowable sentences. The **atomic sentences**—the indivisible syntactic elements—consist of a single **proposition symbol**. Each such symbol stands for a proposition that can be true or false. We will use uppercase names for symbols:  $P$ ,  $Q$ ,  $R$ , and so on. The names are arbitrary but are often chosen to have some mnemonic value to the reader. For example, we might use  $W_{1,3}$  to stand for the proposition that the wumpus is in  $[1,3]$ . (Remember that symbols such as  $W_{1,3}$  are *atomic*, i.e.,  $W$ ,  $1$ , and  $3$  are not meaningful parts of the symbol.) There are two proposition symbols with fixed meanings: *True* is the always-true proposition and *False* is the always-false proposition.

#### COMPLEX SENTENCES LOGICAL CONNECTIVES

**Complex sentences** are constructed from simpler sentences using **logical connectives**. There are five connectives in common use:

#### NEGATION LITERAL

$\neg$  (not). A sentence such as  $\neg W_{1,3}$  is called the **negation** of  $W_{1,3}$ . A **literal** is either an atomic sentence (a **positive literal**) or a negated atomic sentence (a **negative literal**).

<sup>8</sup> Propositional logic is also called **Boolean logic**, after the logician George Boole (1815–1864).

|               |   |
|---------------|---|
| CONJUNCTION   | $\wedge$ (and). A sentence whose main connective is $\wedge$ , such as $W_{1,3} \wedge P_{3,1}$ , is called a <b>conjunction</b> ; its parts are the <b>conjuncts</b> . (The $\wedge$ looks like an “A” for “And.”)   |
| DISJUNCTION   | $\vee$ (or). A sentence using $\vee$ , such as $(W_{1,3} \wedge P_{3,1}) \vee W_{2,2}$ , is a <b>disjunction</b> of the <b>disjuncts</b> $(W_{1,3} \wedge P_{3,1})$ and $W_{2,2}$ . (Historically, the $\vee$ comes from the Latin “vel,” which means “or.” For most people, it is easier to remember as an upside-down $\wedge$ .)   |
| IMPLICATION   | $\Rightarrow$ (implies). A sentence such as $(W_{1,3} \wedge P_{3,1}) \Rightarrow \neg W_{2,2}$ is called an <b>implication</b> (or conditional). Its <b>premise</b> or <b>antecedent</b> is $(W_{1,3} \wedge P_{3,1})$ , and its <b>conclusion</b> or <b>consequent</b> is $\neg W_{2,2}$ . Implications are also known as <b>rules</b> or <b>if–then</b> statements. The implication symbol is sometimes written in other books as $\supset$ or $\rightarrow$ . |
| BICONDITIONAL | $\Leftrightarrow$ (if and only if). The sentence $W_{1,3} \Leftrightarrow \neg W_{2,2}$ is a <b>biconditional</b> .   |

Figure 7.7 gives a formal grammar of propositional logic; see page 984 if you are not familiar with the BNF notation.

|                        |               |   |
|------------------------|---------------|---|
| <i>Sentence</i>        | $\rightarrow$ | <i>AtomicSentence</i>   <i>ComplexSentence</i>        |
| <i>AtomicSentence</i>  | $\rightarrow$ | <b>True</b>   <b>False</b>   <i>Symbol</i>            |
| <i>Symbol</i>          | $\rightarrow$ | <b>P</b>   <b>Q</b>   <b>R</b>   ...                  |
| <i>ComplexSentence</i> | $\rightarrow$ | $\neg$ <i>Sentence</i>                                |
|                        |               | ( <i>Sentence</i> $\wedge$ <i>Sentence</i> )          |
|                        |               | ( <i>Sentence</i> $\vee$ <i>Sentence</i> )            |
|                        |               | ( <i>Sentence</i> $\Rightarrow$ <i>Sentence</i> )     |
|                        |               | ( <i>Sentence</i> $\Leftrightarrow$ <i>Sentence</i> ) |

**Figure 7.7** A BNF (Backus–Naur Form) grammar of sentences in propositional logic.

Notice that the grammar is very strict about parentheses: every sentence constructed with binary connectives must be enclosed in parentheses. This ensures that the syntax is completely unambiguous. It also means that we have to write  $((A \wedge B) \Rightarrow C)$  instead of  $A \wedge B \Rightarrow C$ , for example. To improve readability, we will often omit parentheses, relying instead on an order of precedence for the connectives. This is similar to the precedence used in arithmetic—for example,  $ab + c$  is read as  $((ab) + c)$  rather than  $a(b + c)$  because multiplication has higher precedence than addition. The order of precedence in propositional logic is (from highest to lowest):  $\neg$ ,  $\wedge$ ,  $\vee$ ,  $\Rightarrow$ , and  $\Leftrightarrow$ . Hence, the sentence

$$\neg P \vee Q \wedge R \Rightarrow S$$

is equivalent to the sentence

$$((\neg P) \vee (Q \wedge R)) \Rightarrow S.$$

Precedence does not resolve ambiguity in sentences such as  $A \wedge B \wedge C$ , which could be read as  $((A \wedge B) \wedge C)$  or as  $(A \wedge (B \wedge C))$ . Because these two readings mean the same thing according to the semantics given in the next section, sentences such as  $A \wedge B \wedge C$  are allowed. We also allow  $A \vee B \vee C$  and  $A \Leftrightarrow B \Leftrightarrow C$ . Sentences such as  $A \Rightarrow B \Rightarrow C$  are not

allowed because the two readings have different meanings; we insist on parentheses in this case. Finally, we will sometimes use square brackets instead of parentheses when it makes the sentence clearer.

## Semantics

Having specified the syntax of propositional logic, we now specify its semantics. The semantics defines the rules for determining the truth of a sentence with respect to a particular model. In propositional logic, a model simply fixes the truth value—*true* or *false*—for every proposition symbol. For example, if the sentences in the knowledge base make use of the proposition symbols  $P_{1,2}$ ,  $P_{2,2}$ , and  $P_{3,1}$ , then one possible model is

$$m_1 = \{P_{1,2} = \text{false}, P_{2,2} = \text{false}, P_{3,1} = \text{true}\}.$$

With three proposition symbols, there are  $2^3 = 8$  possible models—exactly those depicted in Figure 7.5. Notice, however, that because we have pinned down the syntax, the models become purely mathematical objects with no necessary connection to wumpus worlds.  $P_{1,2}$  is just a symbol; it might mean “there is a pit in [1,2]” or “I’m in Paris today and tomorrow.”

The semantics for propositional logic must specify how to compute the truth value of *any* sentence, given a model. This is done recursively. All sentences are constructed from atomic sentences and the five connectives; therefore, we need to specify how to compute the truth of atomic sentences and how to compute the truth of sentences formed with each of the five connectives. Atomic sentences are easy:

- *True* is true in every model and *False* is false in every model.
- The truth value of every other proposition symbol must be specified directly in the model. For example, in the model  $m_1$  given earlier,  $P_{1,2}$  is false.

For complex sentences, we have rules such as

- For any sentence  $s$  and any model  $m$ , the sentence  $\neg s$  is true in  $m$  if and only if  $s$  is false in  $m$ .

Such rules reduce the truth of a complex sentence to the truth of simpler sentences. The rules for each connective can be summarized in a **truth table** that specifies the truth value of a complex sentence for each possible assignment of truth values to its components. Truth tables for the five logical connectives are given in Figure 7.8. Using these tables, the truth value of any sentence  $s$  can be computed with respect to any model  $m$  by a simple process of recursive evaluation. For example, the sentence  $\neg P_{1,2} \wedge (P_{2,2} \vee P_{3,1})$ , evaluated in  $m_1$ , gives  $\text{true} \wedge (\text{false} \vee \text{true}) = \text{true} \wedge \text{true} = \text{true}$ . Exercise 7.3 asks you to write the algorithm  $\text{PL-TRUE?}(s, m)$ , which computes the truth value of a propositional logic sentence  $s$  in a model  $m$ .

Previously we said that a knowledge base consists of a set of sentences. We can now see that a logical knowledge base is a conjunction of those sentences. That is, if we start with an empty  $KB$  and do  $\text{TELL}(KB, S_1) \dots \text{TELL}(KB, S_n)$  then we have  $KB = S_1 \wedge \dots \wedge S_n$ . This means that we can treat knowledge bases and sentences interchangeably.

The truth tables for “and,” “or,” and “not” are in close accord with our intuitions about the English words. The main point of possible confusion is that  $P \vee Q$  is true when  $P$  is true

TRUTH TABLE

| $P$   | $Q$   | $\neg P$ | $P \wedge Q$ | $P \vee Q$ | $P \Rightarrow Q$ | $P \Leftrightarrow Q$ |
|-------|-------|----------|--------------|------------|-------------------|-----------------------|
| false | false | true     | false        | false      | true              | true                  |
| false | true  | true     | false        | true       | true              | false                 |
| true  | false | false    | false        | true       | false             | false                 |
| true  | true  | false    | true         | true       | true              | true                  |

**Figure 7.8** Truth tables for the five logical connectives. To use the table to compute, for example, the value of  $P \vee Q$  when  $P$  is true and  $Q$  is false, first look on the left for the row where  $P$  is true and  $Q$  is false (the third row). Then look in that row under the  $P \vee Q$  column to see the result: *true*. Another way to look at this is to think of each row as a model, and the entries under each column for that row as saying whether the corresponding sentence is true in that model.

or  $Q$  is true *or both*. There is a different connective called “exclusive or” (“xor” for short) that yields false when both disjuncts are true.<sup>9</sup> There is no consensus on the symbol for exclusive or; two choices are  $\dot{\vee}$  and  $\oplus$ .

The truth table for  $\Rightarrow$  may seem puzzling at first, because it might not quite fit one’s intuitive understanding of “ $P$  implies  $Q$ ” or “if  $P$  then  $Q$ .” For one thing, propositional logic does not require any relation of causation or relevance between  $P$  and  $Q$ . The sentence “5 is odd implies Tokyo is the capital of Japan” is a true sentence of propositional logic (under the normal interpretation), even though it is a decidedly odd sentence of English. Another point of confusion is that any implication is true whenever its antecedent is false. For example, “5 is even implies Sam is smart” is true, regardless of whether Sam is smart. This seems bizarre, but it makes sense if you think of “ $P \Rightarrow Q$ ” as saying, “If  $P$  is true, then I am claiming that  $Q$  is true. Otherwise I am making no claim.” The only way for this sentence to be *false* is if  $P$  is true but  $Q$  is false.

The truth table for a biconditional,  $P \Leftrightarrow Q$ , shows that it is true whenever both  $P \Rightarrow Q$  and  $Q \Rightarrow P$  are true. In English, this is often written as “ $P$  if and only if  $Q$ ” or “ $P$  iff  $Q$ .” The rules of the wumpus world are best written using  $\Leftrightarrow$ . For example, a square is breezy *if* a neighboring square has a pit, and a square is breezy *only if* a neighboring square has a pit. So we need biconditionals such as

$$B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1}) ,$$

where  $B_{1,1}$  means that there is a breeze in  $[1,1]$ . Notice that the one-way implication

$$B_{1,1} \Rightarrow (P_{1,2} \vee P_{2,1})$$

is true in the wumpus world, but incomplete. It does not rule out models in which  $B_{1,1}$  is false and  $P_{1,2}$  is true, which would violate the rules of the wumpus world. Another way of putting it is that the implication requires the presence of pits if there is a breeze, whereas the biconditional also requires the absence of pits if there is no breeze.

<sup>9</sup> Latin has a separate word, *aut*, for exclusive or.

## A simple knowledge base

Now that we have defined the semantics for propositional logic, we can construct a knowledge base for the wumpus world. For simplicity, we will deal only with pits; the wumpus itself is left as an exercise. We will provide enough knowledge to carry out the inference that was done informally in Section 7.3.

First, we need to choose our vocabulary of proposition symbols. For each  $i, j$ :

- Let  $P_{i,j}$  be true if there is a pit in  $[i, j]$ .
- Let  $B_{i,j}$  be true if there is a breeze in  $[i, j]$ .

The knowledge base includes the following sentences, each one labeled for convenience:

- There is no pit in  $[1,1]$ :

$$R_1 : \neg P_{1,1} .$$

- A square is breezy if and only if there is a pit in a neighboring square. This has to be stated for each square; for now, we include just the relevant squares:

$$R_2 : B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1}) .$$

$$R_3 : B_{2,1} \Leftrightarrow (P_{1,1} \vee P_{2,2} \vee P_{3,1}) .$$

- The preceding sentences are true in all wumpus worlds. Now we include the breeze percepts for the first two squares visited in the specific world the agent is in, leading up to the situation in Figure 7.3(b).

$$R_4 : \neg B_{1,1} .$$

$$R_5 : B_{2,1} .$$

The knowledge base, then, consists of sentences  $R_1$  through  $R_5$ . It can also be considered as a single sentence—the conjunction  $R_1 \wedge R_2 \wedge R_3 \wedge R_4 \wedge R_5$ —because it asserts that all the individual sentences are true.

## Inference

Recall that the aim of logical inference is to decide whether  $KB \models \alpha$  for some sentence  $\alpha$ . For example, is  $P_{2,2}$  entailed? Our first algorithm for inference will be a direct implementation of the definition of entailment: enumerate the models, and check that  $\alpha$  is true in every model in which  $KB$  is true. For propositional logic, models are assignments of *true* or *false* to every proposition symbol. Returning to our wumpus-world example, the relevant proposition symbols are  $B_{1,1}$ ,  $B_{2,1}$ ,  $P_{1,1}$ ,  $P_{1,2}$ ,  $P_{2,1}$ ,  $P_{2,2}$ , and  $P_{3,1}$ . With seven symbols, there are  $2^7 = 128$  possible models; in three of these,  $KB$  is true (Figure 7.9). In those three models,  $\neg P_{1,2}$  is true, hence there is no pit in  $[1,2]$ . On the other hand,  $P_{2,2}$  is true in two of the three models and false in one, so we cannot yet tell whether there is a pit in  $[2,2]$ .

Figure 7.9 reproduces in a more precise form the reasoning illustrated in Figure 7.5. A general algorithm for deciding entailment in propositional logic is shown in Figure 7.10. Like the BACKTRACKING-SEARCH algorithm on page 76, TT-ENTAILS? performs a recursive enumeration of a finite space of assignments to variables. The algorithm is **sound**, because it



| $B_{1,1}$ | $B_{2,1}$ | $P_{1,1}$ | $P_{1,2}$ | $P_{2,1}$ | $P_{2,2}$ | $P_{3,1}$ | $R_1$    | $R_2$    | $R_3$    | $R_4$    | $R_5$    | $KB$        |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|----------|----------|----------|----------|----------|-------------|
| false     | false     | false     | false     | false     | false     | false     | true     | true     | true     | true     | false    | false       |
| false     | false     | false     | false     | false     | false     | true      | true     | true     | false    | true     | false    | false       |
| $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$    |
| false     | true      | false     | false     | false     | false     | false     | true     | true     | false    | true     | true     | false       |
| false     | true      | false     | false     | false     | false     | true      | true     | true     | true     | true     | true     | <u>true</u> |
| false     | true      | false     | false     | false     | true      | true      | true     | true     | true     | true     | true     | <u>true</u> |
| false     | true      | false     | false     | true      | false     | false     | true     | false    | false    | true     | true     | false       |
| $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$  | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$    |
| true      | true      | true      | true      | true      | true      | true      | false    | true     | true     | false    | true     | false       |

**Figure 7.9** A truth table constructed for the knowledge base given in the text.  $KB$  is true if  $R_1$  through  $R_5$  are true, which occurs in just 3 of the 128 rows. In all 3 rows,  $P_{1,2}$  is false, so there is no pit in  $[1,2]$ . On the other hand, there might (or might not) be a pit in  $[2,2]$ .

```

function TT-ENTAILS?( $KB, \alpha$ ) returns true or false
  inputs:  $KB$ , the knowledge base, a sentence in propositional logic
            $\alpha$ , the query, a sentence in propositional logic

   $symbols \leftarrow$  a list of the proposition symbols in  $KB$  and  $\alpha$ 
  return TT-CHECK-ALL( $KB, \alpha, symbols, []$ )



---


function TT-CHECK-ALL( $KB, \alpha, symbols, model$ ) returns true or false
  if EMPTY?( $symbols$ ) then
    if PL-TRUE?( $KB, model$ ) then return PL-TRUE?( $\alpha, model$ )
    else return true
  else do
     $P \leftarrow$  FIRST( $symbols$ );  $rest \leftarrow$  REST( $symbols$ )
    return TT-CHECK-ALL( $KB, \alpha, rest, EXTEND(P, true, model)$ ) and
           TT-CHECK-ALL( $KB, \alpha, rest, EXTEND(P, false, model)$ )

```

**Figure 7.10** A truth-table enumeration algorithm for deciding propositional entailment. TT stands for truth table. PL-TRUE? returns true if a sentence holds within a model. The variable  $model$  represents a partial model—an assignment to only some of the variables. The function call  $EXTEND(P, true, model)$  returns a new partial model in which  $P$  has the value *true*.

implements directly the definition of entailment, and **complete**, because it works for any  $KB$  and  $\alpha$  and always terminates—there are only finitely many models to examine.

Of course, “finitely many” is not always the same as “few.” If  $KB$  and  $\alpha$  contain  $n$  symbols in all, then there are  $2^n$  models. Thus, the time complexity of the algorithm is  $O(2^n)$ . (The space complexity is only  $O(n)$  because the enumeration is depth-first.) Later in this

|  |  |
|--|--|
| $(\alpha \wedge \beta) \equiv (\beta \wedge \alpha)$   | commutativity of $\wedge$              |
| $(\alpha \vee \beta) \equiv (\beta \vee \alpha)$   | commutativity of $\vee$                |
| $((\alpha \wedge \beta) \wedge \gamma) \equiv (\alpha \wedge (\beta \wedge \gamma))$                   | associativity of $\wedge$              |
| $((\alpha \vee \beta) \vee \gamma) \equiv (\alpha \vee (\beta \vee \gamma))$                           | associativity of $\vee$                |
| $\neg(\neg\alpha) \equiv \alpha$   | double-negation elimination            |
| $(\alpha \Rightarrow \beta) \equiv (\neg\beta \Rightarrow \neg\alpha)$                                 | contraposition                         |
| $(\alpha \Rightarrow \beta) \equiv (\neg\alpha \vee \beta)$  | implication elimination                |
| $(\alpha \Leftrightarrow \beta) \equiv ((\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha))$ | biconditional elimination              |
| $\neg(\alpha \wedge \beta) \equiv (\neg\alpha \vee \neg\beta)$   | de Morgan                              |
| $\neg(\alpha \vee \beta) \equiv (\neg\alpha \wedge \neg\beta)$   | de Morgan                              |
| $(\alpha \wedge (\beta \vee \gamma)) \equiv ((\alpha \wedge \beta) \vee (\alpha \wedge \gamma))$       | distributivity of $\wedge$ over $\vee$ |
| $(\alpha \vee (\beta \wedge \gamma)) \equiv ((\alpha \vee \beta) \wedge (\alpha \vee \gamma))$         | distributivity of $\vee$ over $\wedge$ |

**Figure 7.11** Standard logical equivalences. The symbols  $\alpha$ ,  $\beta$ , and  $\gamma$  stand for arbitrary sentences of propositional logic.



chapter, we will see algorithms that are much more efficient in practice. Unfortunately, *every known inference algorithm for propositional logic has a worst-case complexity that is exponential in the size of the input*. We do not expect to do better than this because propositional entailment is co-NP-complete. (See Appendix A.)

### Equivalence, validity, and satisfiability

Before we plunge into the details of logical inference algorithms, we will need some additional concepts related to entailment. Like entailment, these concepts apply to all forms of logic, but they are best illustrated for a particular logic, such as propositional logic.

LOGICAL  
EQUIVALENCE

The first concept is **logical equivalence**: two sentences  $\alpha$  and  $\beta$  are logically equivalent if they are true in the same set of models. We write this as  $\alpha \Leftrightarrow \beta$ . For example, we can easily show (using truth tables) that  $P \wedge Q$  and  $Q \wedge P$  are logically equivalent; other equivalences are shown in Figure 7.11. They play much the same role in logic as arithmetic identities do in ordinary mathematics. An alternative definition of equivalence is as follows: for any two sentences  $\alpha$  and  $\beta$ ,

$$\alpha \equiv \beta \quad \text{if and only if} \quad \alpha \models \beta \text{ and } \beta \models \alpha.$$

(Recall that  $\models$  means entailment.)

VALIDITY

The second concept we will need is **validity**. A sentence is valid if it is true in *all* models. For example, the sentence  $P \vee \neg P$  is valid. Valid sentences are also known as **tautologies**—they are *necessarily* true and hence vacuous. Because the sentence *True* is true in all models, every valid sentence is logically equivalent to *True*.

TAUTOLOGY

What good are valid sentences? From our definition of entailment, we can derive the **deduction theorem**, which was known to the ancient Greeks:

DEDUCTION  
THEOREM



*For any sentences  $\alpha$  and  $\beta$ ,  $\alpha \models \beta$  if and only if the sentence  $(\alpha \Rightarrow \beta)$  is valid.*

(Exercise 7.4 asks for a proof.) We can think of the inference algorithm in Figure 7.10 as

checking the validity of  $(KB \Rightarrow \alpha)$ . Conversely, every valid implication sentence describes a legitimate inference.

SATISFIABILITY

The final concept we will need is **satisfiability**. A sentence is satisfiable if it is true in *some* model. For example, the knowledge base given earlier,  $(R_1 \wedge R_2 \wedge R_3 \wedge R_4 \wedge R_5)$ , is satisfiable because there are three models in which it is true, as shown in Figure 7.9. If a sentence  $\alpha$  is true in a model  $m$ , then we say that  $m$  **satisfies**  $\alpha$ , or that  $m$  is a **model of**  $\alpha$ . Satisfiability can be checked by enumerating the possible models until one is found that satisfies the sentence. Determining the satisfiability of sentences in propositional logic was the first problem proved to be NP-complete.

SATISFIES

Many problems in computer science are really satisfiability problems. For example, all the constraint satisfaction problems in Chapter 5 are essentially asking whether the constraints are satisfiable by some assignment. With appropriate transformations, search problems can also be solved by checking satisfiability. Validity and satisfiability are of course connected:  $\alpha$  is valid iff  $\neg\alpha$  is unsatisfiable; contrapositively,  $\alpha$  is satisfiable iff  $\neg\alpha$  is not valid. We also have the following useful result:



REDUCTIO AD  
ABSURDUM  
REFUTATION

$\alpha \models \beta$  if and only if the sentence  $(\alpha \wedge \neg\beta)$  is unsatisfiable.

Proving  $\beta$  from  $\alpha$  by checking the unsatisfiability of  $(\alpha \wedge \neg\beta)$  corresponds exactly to the standard mathematical proof technique of *reductio ad absurdum* (literally, “reduction to an absurd thing”). It is also called proof by **refutation** or proof by **contradiction**. One assumes a sentence  $\beta$  to be false and shows that this leads to a contradiction with known axioms  $\alpha$ . This contradiction is exactly what is meant by saying that the sentence  $(\alpha \wedge \neg\beta)$  is unsatisfiable.

## 7.5 REASONING PATTERNS IN PROPOSITIONAL LOGIC

This section covers standard patterns of inference that can be applied to derive chains of conclusions that lead to the desired goal. These patterns of inference are called **inference rules**. The best-known rule is called **Modus Ponens** and is written as follows:

INFERENCE RULES

MODUS PONENS

$$\frac{\alpha \Rightarrow \beta, \quad \alpha}{\beta}$$

The notation means that, whenever any sentences of the form  $\alpha \Rightarrow \beta$  and  $\alpha$  are given, then the sentence  $\beta$  can be inferred. For example, if  $(WumpusAhead \wedge WumpusAlive) \Rightarrow Shoot$  and  $(WumpusAhead \wedge WumpusAlive)$  are given, then *Shoot* can be inferred.

AND-ELIMINATION

Another useful inference rule is **And-Elimination**, which says that, from a conjunction, any of the conjuncts can be inferred:

$$\frac{\alpha \wedge \beta}{\alpha}$$

For example, from  $(WumpusAhead \wedge WumpusAlive)$ , *WumpusAlive* can be inferred.

By considering the possible truth values of  $\alpha$  and  $\beta$ , one can show easily that Modus Ponens and And-Elimination are sound once and for all. These rules can then be used in any particular instances where they apply, generating sound inferences without the need for enumerating models.

All of the logical equivalences in Figure 7.11 can be used as inference rules. For example, the equivalence for biconditional elimination yields the two inference rules

$$\frac{\alpha \Leftrightarrow \beta}{(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)} \quad \text{and} \quad \frac{(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)}{\alpha \Leftrightarrow \beta}.$$

Not all inference rules work in both directions like this. For example, we cannot run Modus Ponens in the opposite direction to obtain  $\alpha \Rightarrow \beta$  and  $\alpha$  from  $\beta$ .

Let us see how these inference rules and equivalences can be used in the wumpus world. We start with the knowledge base containing  $R_1$  through  $R_5$ , and show how to prove  $\neg P_{1,2}$ , that is, there is no pit in [1,2]. First, we apply biconditional elimination to  $R_2$  to obtain

$$R_6 : (B_{1,1} \Rightarrow (P_{1,2} \vee P_{2,1})) \wedge ((P_{1,2} \vee P_{2,1}) \Rightarrow B_{1,1}).$$

Then we apply And-Elimination to  $R_6$  to obtain

$$R_7 : ((P_{1,2} \vee P_{2,1}) \Rightarrow B_{1,1}).$$

Logical equivalence for contrapositives gives

$$R_8 : (\neg B_{1,1} \Rightarrow \neg(P_{1,2} \vee P_{2,1})).$$

Now we can apply Modus Ponens with  $R_8$  and the percept  $R_4$  (i.e.,  $\neg B_{1,1}$ ), to obtain

$$R_9 : \neg(P_{1,2} \vee P_{2,1}).$$

Finally, we apply de Morgan's rule, giving the conclusion

$$R_{10} : \neg P_{1,2} \wedge \neg P_{2,1}.$$

That is, neither [1,2] nor [2,1] contains a pit.

PROOF

The preceding derivation—a sequence of applications of inference rules—is called a **proof**. Finding proofs is exactly like finding solutions to search problems. In fact, if the successor function is defined to generate all possible applications of inference rules, then all of the search algorithms in Chapters 3 and 4 can be applied to find proofs. Thus, searching for proofs is an alternative to enumerating models. The search can go forward from the initial knowledge base, applying inference rules to derive the goal sentence, or it can go backward from the goal sentence, trying to find a chain of inference rules leading from the initial knowledge base. Later in this section, we will see two families of algorithms that use these techniques.



The fact that inference in propositional logic is NP-complete suggests that, in the worst case, searching for proofs is going to be no more efficient than enumerating models. In many practical cases, however, *finding a proof can be highly efficient simply because it can ignore irrelevant propositions, no matter how many of them there are*. For example, the proof given earlier leading to  $\neg P_{1,2} \wedge \neg P_{2,1}$  does not mention the propositions  $B_{2,1}$ ,  $P_{1,1}$ ,  $P_{2,2}$ , or  $P_{3,1}$ . They can be ignored because the goal proposition,  $P_{1,2}$ , appears only in sentence  $R_4$ ; the other propositions in  $R_4$  appear only in  $R_4$  and  $R_2$ ; so  $R_1$ ,  $R_3$ , and  $R_5$  have no bearing on the proof. The same would hold even if we added a million more sentences to the knowledge base; the simple truth-table algorithm, on the other hand, would be overwhelmed by the exponential explosion of models.

MONOTONICITY

This property of logical systems actually follows from a much more fundamental property called **monotonicity**. Monotonicity says that the set of entailed sentences can only *in-*

crease as information is added to the knowledge base.<sup>10</sup> For any sentences  $\alpha$  and  $\beta$ ,

$$\text{if } KB \models \alpha \text{ then } KB \wedge \beta \models \alpha .$$

For example, suppose the knowledge base contains the additional assertion  $\beta$  stating that there are exactly eight pits in the world. This knowledge might help the agent draw *additional* conclusions, but it cannot invalidate any conclusion  $\alpha$  already inferred—such as the conclusion that there is no pit in [1,2]. Monotonicity means that inference rules can be applied whenever suitable premises are found in the knowledge base—the conclusion of the rule must follow *regardless of what else is in the knowledge base*.

## Resolution

We have argued that the inference rules covered so far are *sound*, but we have not discussed the question of *completeness* for the inference algorithms that use them. Search algorithms such as iterative deepening search (page 78) are complete in the sense that they will find any reachable goal, but if the available inference rules are inadequate, then the goal is not reachable—no proof exists that uses only those inference rules. For example, if we removed the biconditional elimination rule, the proof in the preceding section would not go through. The current section introduces a single inference rule, **resolution**, that yields a complete inference algorithm when coupled with any complete search algorithm.

We begin by using a simple version of the resolution rule in the wumpus world. Let us consider the steps leading up to Figure 7.4(a): the agent returns from [2,1] to [1,1] and then goes to [1,2], where it perceives a stench, but no breeze. We add the following facts to the knowledge base:

$$\begin{aligned} R_{11} : & \neg B_{1,2} . \\ R_{12} : & B_{1,2} \Leftrightarrow (P_{1,1} \vee P_{2,2} \vee P_{1,3}) . \end{aligned}$$

By the same process that led to  $R_{10}$  earlier, we can now derive the absence of pits in [2,2] and [1,3] (remember that [1,1] is already known to be pitless):

$$\begin{aligned} R_{13} : & \neg P_{2,2} . \\ R_{14} : & \neg P_{1,3} . \end{aligned}$$

We can also apply biconditional elimination to  $R_3$ , followed by modus ponens with  $R_5$ , to obtain the fact that there is a pit in [1,1], [2,2], or [3,1]:

$$R_{15} : P_{1,1} \vee P_{2,2} \vee P_{3,1} .$$

Now comes the first application of the resolution rule: the literal  $\neg P_{2,2}$  in  $R_{13}$  *resolves with* the literal  $P_{2,2}$  in  $R_{15}$  to give

$$R_{16} : P_{1,1} \vee P_{3,1} .$$

In English; if there's a pit in one of [1,1], [2,2], and [3,1], and it's not in [2,2], then it's in [1,1] or [3,1]. Similarly, the literal  $\neg P_{1,1}$  in  $R_1$  resolves with the literal  $P_{1,1}$  in  $R_{16}$  to give

$$R_{17} : P_{3,1} .$$

---

<sup>10</sup> **Nonmonotonic** logics, which violate the monotonicity property, capture a common property of human reasoning: changing one's mind. They are discussed in Section 10.7.

UNIT RESOLUTION In English: if there's a pit in [1,1] or [3,1], and it's not in [1,1], then it's in [3,1]. These last two inference steps are examples of the **unit resolution** inference rule,

$$\frac{\ell_1 \vee \dots \vee \ell_k, \quad m}{\ell_1 \vee \dots \vee \ell_{i-1} \vee \ell_{i+1} \vee \dots \vee \ell_k},$$

COMPLEMENTARY  
LITERALS  
CLAUSE

where each  $\ell$  is a literal and  $\ell_i$  and  $m$  are **complementary literals** (i.e., one is the negation of the other). Thus, the unit resolution rule takes a **clause**—a disjunction of literals—and a literal and produces a new clause. Note that a single literal can be viewed as a disjunction of one literal, also known as a **unit clause**.

UNIT CLAUSE  
RESOLUTION

The unit resolution rule can be generalized to the full **resolution** rule,

$$\frac{\ell_1 \vee \dots \vee \ell_k, \quad m_1 \vee \dots \vee m_n}{\ell_1 \vee \dots \vee \ell_{i-1} \vee \ell_{i+1} \vee \dots \vee \ell_k \vee m_1 \vee \dots \vee m_{j-1} \vee m_{j+1} \vee \dots \vee m_n},$$

where  $\ell_i$  and  $m_j$  are complementary literals. If we were dealing only with clauses of length two we could write this as

$$\frac{\ell_1 \vee \ell_2, \quad \neg \ell_2 \vee \ell_3}{\ell_1 \vee \ell_3}.$$

That is, resolution takes two clauses and produces a new clause containing all the literals of the two original clauses *except* the two complementary literals. For example, we have

$$\frac{P_{1,1} \vee P_{3,1}, \quad \neg P_{1,1} \vee \neg P_{2,2}}{P_{3,1} \vee \neg P_{2,2}}.$$

FACTORING

There is one more technical aspect of the resolution rule: the resulting clause should contain only one copy of each literal.<sup>11</sup> The removal of multiple copies of literals is called **factoring**. For example, if we resolve  $(A \vee B)$  with  $(A \vee \neg B)$ , we obtain  $(A \vee A)$ , which is reduced to just  $A$ .

The *soundness* of the resolution rule can be seen easily by considering the literal  $\ell_i$ . If  $\ell_i$  is true, then  $m_j$  is false, and hence  $m_1 \vee \dots \vee m_{j-1} \vee m_{j+1} \vee \dots \vee m_n$  must be true, because  $m_1 \vee \dots \vee m_n$  is given. If  $\ell_i$  is false, then  $\ell_1 \vee \dots \vee \ell_{i-1} \vee \ell_{i+1} \vee \dots \vee \ell_k$  must be true because  $\ell_1 \vee \dots \vee \ell_k$  is given. Now  $\ell_i$  is either true or false, so one or other of these conclusions holds—exactly as the resolution rule states.



REFUTATION  
COMPLETENESS

What is more surprising about the resolution rule is that it forms the basis for a family of *complete* inference procedures. *Any complete search algorithm, applying only the resolution rule, can derive any conclusion entailed by any knowledge base in propositional logic.* There is a caveat: resolution is complete in a specialized sense. Given that  $A$  is true, we cannot use resolution to automatically generate the consequence  $A \vee B$ . However, we can use resolution to answer the question of whether  $A \vee B$  is true. This is called **refutation completeness**, meaning that resolution can always be used to either confirm or refute a sentence, but it cannot be used to enumerate true sentences. The next two subsections explain how resolution accomplishes this.

<sup>11</sup> If a clause is viewed as a *set* of literals, then this restriction is automatically respected. Using set notation for clauses makes the resolution rule much cleaner, at the cost of introducing additional notation.



CONJUNCTIVE  
NORMAL FORM  
K-CNF

### Conjunctive normal form

The resolution rule applies only to disjunctions of literals, so it would seem to be relevant only to knowledge bases and queries consisting of such disjunctions. How, then, can it lead to a complete inference procedure for all of propositional logic? The answer is that *every sentence of propositional logic is logically equivalent to a conjunction of disjunctions of literals*. A sentence expressed as a conjunction of disjunctions of literals is said to be in **conjunctive normal form** or **CNF**. We will also find it useful later to consider the restricted family of **k-CNF** sentences. A sentence in **k-CNF** has exactly  $k$  literals per clause:

$$(\ell_{1,1} \vee \dots \vee \ell_{1,k}) \wedge \dots \wedge (\ell_{n,1} \vee \dots \vee \ell_{n,k}) .$$

It turns out that every sentence can be transformed into a 3-CNF sentence that has an equivalent set of models.

Rather than prove these assertions (see Exercise 7.10), we describe a simple conversion procedure. We illustrate the procedure by converting  $R_2$ , the sentence  $B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1})$ , into CNF. The steps are as follows:

1. Eliminate  $\Leftrightarrow$ , replacing  $\alpha \Leftrightarrow \beta$  with  $(\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha)$ .

$$(B_{1,1} \Rightarrow (P_{1,2} \vee P_{2,1})) \wedge ((P_{1,2} \vee P_{2,1}) \Rightarrow B_{1,1}) .$$

2. Eliminate  $\Rightarrow$ , replacing  $\alpha \Rightarrow \beta$  with  $\neg\alpha \vee \beta$ :

$$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge (\neg(P_{1,2} \vee P_{2,1}) \vee B_{1,1}) .$$

3. CNF requires  $\neg$  to appear only in literals, so we “move  $\neg$  inwards” by repeated application of the following equivalences from Figure 7.11:

$$\neg(\neg\alpha) \equiv \alpha \quad (\text{double-negation elimination})$$

$$\neg(\alpha \wedge \beta) \equiv (\neg\alpha \vee \neg\beta) \quad (\text{de Morgan})$$

$$\neg(\alpha \vee \beta) \equiv (\neg\alpha \wedge \neg\beta) \quad (\text{de Morgan})$$

In the example, we require just one application of the last rule:

$$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge ((\neg P_{1,2} \wedge \neg P_{2,1}) \vee B_{1,1}) .$$

4. Now we have a sentence containing nested  $\wedge$  and  $\vee$  operators applied to literals. We apply the distributivity law from Figure 7.11, distributing  $\vee$  over  $\wedge$  wherever possible.

$$(\neg B_{1,1} \vee P_{1,2} \vee P_{2,1}) \wedge (\neg P_{1,2} \vee B_{1,1}) \wedge (\neg P_{2,1} \vee B_{1,1}) .$$

The original sentence is now in CNF, as a conjunction of three clauses. It is much harder to read, but it can be used as input to a resolution procedure.

### A resolution algorithm

Inference procedures based on resolution work by using the principle of proof by contradiction discussed at the end of Section 7.4. That is, to show that  $KB \models \alpha$ , we show that  $(KB \wedge \neg\alpha)$  is unsatisfiable. We do this by proving a contradiction.

A resolution algorithm is shown in Figure 7.12. First,  $(KB \wedge \neg\alpha)$  is converted into CNF. Then, the resolution rule is applied to the resulting clauses. Each pair that contains complementary literals is resolved to produce a new clause, which is added to the set if it is not already present. The process continues until one of two things happens:

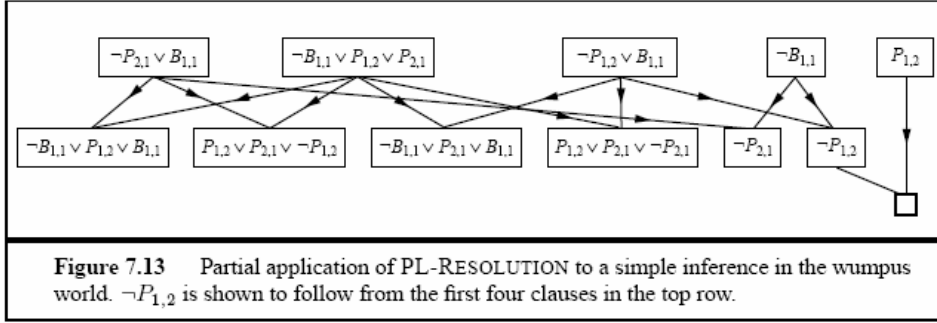
```

function PL-RESOLUTION( $KB, \alpha$ ) returns true or false
inputs:  $KB$ , the knowledge base, a sentence in propositional logic
          $\alpha$ , the query, a sentence in propositional logic

 $clauses \leftarrow$  the set of clauses in the CNF representation of  $KB \wedge \neg\alpha$ 
 $new \leftarrow \{ \}$ 
loop do
  for each  $C_i, C_j$  in  $clauses$  do
     $resolvents \leftarrow$  PL-RESOLVE( $C_i, C_j$ )
    if  $resolvents$  contains the empty clause then return true
     $new \leftarrow new \cup resolvents$ 
  if  $new \subseteq clauses$  then return false
   $clauses \leftarrow clauses \cup new$ 

```

**Figure 7.12** A simple resolution algorithm for propositional logic. The function PL-RESOLVE returns the set of all possible clauses obtained by resolving its two inputs.



**Figure 7.13** Partial application of PL-RESOLUTION to a simple inference in the wumpus world.  $\neg P_{1,2}$  is shown to follow from the first four clauses in the top row.

- there are no new clauses that can be added, in which case  $\alpha$  does not entail  $\beta$  or
- an application of the resolution rule derives the *empty* clause, in which case  $\alpha$  entails  $\beta$ .

The empty clause—a disjunction of no disjuncts—is equivalent to *False* because a disjunction is true only if at least one of its disjuncts is true. Another way to see that an empty clause represents a contradiction is to observe that it arises only from resolving two complementary unit clauses such as  $P$  and  $\neg P$ .

We can apply the resolution procedure to a very simple inference in the wumpus world. When the agent is in  $[1,1]$ , there is no breeze, so there can be no pits in neighboring squares. The relevant knowledge base is

$$KB = R_2 \wedge R_4 = (B_{1,1} \Leftrightarrow (P_{1,2} \vee P_{2,1})) \wedge \neg B_{1,1}$$

and we wish to prove  $\alpha$  which is, say,  $\neg P_{1,2}$ . When we convert  $(KB \wedge \neg\alpha)$  into CNF, we obtain the clauses shown at the top of Figure 7.13. The second row of the figure shows all the clauses obtained by resolving pairs in the first row. Then, when  $P_{1,2}$  is resolved with  $\neg P_{1,2}$ , we obtain the empty clause, shown as a small square. Inspection of Figure 7.13 reveals that



many resolution steps are pointless. For example, the clause  $B_{1,1} \vee \neg B_{1,1} \vee P_{1,2}$  is equivalent to  $\text{True} \vee P_{1,2}$  which is equivalent to  $\text{True}$ . Deducing that  $\text{True}$  is true is not very helpful. Therefore, any clause in which two complementary literals appear can be discarded.

### Completeness of resolution

RESOLUTION  
CLOSURE

To conclude our discussion of resolution, we now show why PL-RESOLUTION is complete. To do this, it will be useful to introduce the **resolution closure**  $RC(S)$  of a set of clauses  $S$ , which is the set of all clauses derivable by repeated application of the resolution rule to clauses in  $S$  or their derivatives. The resolution closure is what PL-RESOLUTION computes as the final value of the variable *clauses*. It is easy to see that  $RC(S)$  must be finite, because there are only finitely many distinct clauses that can be constructed out of the symbols  $P_1, \dots, P_k$  that appear in  $S$ . (Notice that this would not be true without the factoring step that removes multiple copies of literals.) Hence, PL-RESOLUTION always terminates.

GROUND  
RESOLUTION  
THEOREM

The completeness theorem for resolution in propositional logic is called the **ground resolution theorem**:

If a set of clauses is unsatisfiable, then the resolution closure of those clauses contains the empty clause.

We prove this theorem by demonstrating its contrapositive: if the closure  $RC(S)$  does *not* contain the empty clause, then  $S$  is satisfiable. In fact, we can construct a model for  $S$  with suitable truth values for  $P_1, \dots, P_k$ . The construction procedure is as follows:

For  $i$  from 1 to  $k$ ,

- If there is a clause in  $RC(S)$  containing the literal  $\neg P_i$  such that all its other literals are false under the assignment chosen for  $P_1, \dots, P_{i-1}$ , then assign *false* to  $P_i$ .
- Otherwise, assign *true* to  $P_i$ .

It remains to show that this assignment to  $P_1, \dots, P_k$  is a model of  $S$ , provided that  $RC(S)$  is closed under resolution and does not contain the empty clause. The proof of this is left as an exercise.

In this section, we bring together what we have learned so far in order to construct agents that operate using propositional logic. We will look at two kinds of agents: those which use inference algorithms and a knowledge base, like the generic knowledge-based agent in Figure 7.1, and those which evaluate logical expressions directly in the form of circuits. We will demonstrate both kinds of agents in the wumpus world, and will find that both suffer from serious drawbacks.

### Finding pits and wumpuses using logical inference

Let us begin with an agent that reasons logically about the location of pits, wumpuses, and safe squares. It begins with a knowledge base that states the “physics” of the wumpus world. It knows that  $[1,1]$  does not contain a pit or a wumpus; that is,  $\neg P_{1,1}$  and  $\neg W_{1,1}$ . For every square  $[x, y]$ , it knows a sentence stating how a breeze arises:

$$B_{x,y} \Leftrightarrow (P_{x,y+1} \vee P_{x,y-1} \vee P_{x+1,y} \vee P_{x-1,y}) . \quad (7.1)$$

For every square  $[x, y]$ , it knows a sentence stating how a stench arises:

$$S_{x,y} \Leftrightarrow (W_{x,y+1} \vee W_{x,y-1} \vee W_{x+1,y} \vee W_{x-1,y}) . \quad (7.2)$$

Finally, it knows that there is exactly one wumpus. This is expressed in two parts. First, we have to say that there is *at least one* wumpus:

$$W_{1,1} \vee W_{1,2} \vee \dots \vee W_{4,3} \vee W_{4,4} .$$

Then, we have to say that there is *at most one* wumpus. One way to do this is to say that for any two squares, one of them must be wumpus-free. With  $n$  squares, we get  $n(n-1)/2$

```

function PL-WUMPUS-AGENT(percept) returns an action
  inputs: percept, a list, [stench, breeze, glitter]
  static: KB, a knowledge base, initially containing the “physics” of the wumpus world
           x, y, orientation, the agent’s position (initially 1,1) and orientation (initially right)
           visited, an array indicating which squares have been visited, initially false
           action, the agent’s most recent action, initially null
           plan, an action sequence, initially empty

  update x, y, orientation, visited based on action
  if stench then TELL(KB, Sx,y) else TELL(KB,  $\neg S_{x,y}$ )
  if breeze then TELL(KB, Bx,y) else TELL(KB,  $\neg B_{x,y}$ )
  if glitter then action  $\leftarrow$  grab
  else if plan is nonempty then action  $\leftarrow$  POP(plan)
  else if for some fringe square [i,j], ASK(KB, ( $\neg P_{i,j} \wedge \neg W_{i,j}$ )) is true or
           for some fringe square [i,j], ASK(KB, (Pi,j  $\vee$  Wi,j)) is false then do
           plan  $\leftarrow$  A*-GRAPH-SEARCH(ROUTE-PROBLEM([x,y], orientation, [i,j], visited))
           action  $\leftarrow$  POP(plan)
  else action  $\leftarrow$  a randomly chosen move
  return action

```

**Figure 7.19** A wumpus-world agent that uses propositional logic to identify pits, wumpuses, and safe squares. The subroutine ROUTE-PROBLEM constructs a search problem whose solution is a sequence of actions leading from [*x,y*] to [*i,j*] and passing through only previously visited squares.

sentences such as  $\neg W_{1,1} \vee \neg W_{1,2}$ . For a  $4 \times 4$  world, then, we begin with a total of 155 sentences containing 64 distinct symbols.

The agent program, shown in Figure 7.19, TELLS its knowledge base about each new breeze and stench percept. (It also updates some ordinary program variables to keep track of where it is and where it has been—more on this later.) Then, the program chooses where to look next among the fringe squares—that is, the squares adjacent to those already visited. A fringe square [*i,j*] is *provably safe* if the sentence  $(\neg P_{i,j} \wedge \neg W_{i,j})$  is entailed by the knowledge base. The next best thing is a *possibly safe* square, for which the agent cannot prove that there is a pit or a wumpus—that is, for which  $(P_{i,j} \vee W_{i,j})$  is *not* entailed.

The entailment computation in ASK can be implemented using any of the methods described earlier in the chapter. TT-ENTAILS? (Figure 7.10) is obviously impractical, since it would have to enumerate  $2^{64}$  rows. DPLL (Figure 7.16) performs the required inferences in a few milliseconds, thanks mainly to the unit propagation heuristic. WALKSAT can also be used, with the usual caveats about incompleteness. In wumpus worlds, failures to find a model, given 10,000 flips, invariably correspond to unsatisfiability, so no errors are likely due to incompleteness.

PL-WUMPUS-AGENT works quite well in a small wumpus world. There is, however, something deeply unsatisfying about the agent’s knowledge base. *KB* contains “physics” sentences of the form given in Equations (7.1) and (7.2) for *every single square*. The larger

the environment, the larger the initial knowledge base needs to be. We would much prefer to have just two sentences that say how breezes and stench arise in *all* squares. These are beyond the powers of propositional logic to express. In the next chapter, we will see a more expressive logical language in which such sentences are easy to express.

### Keeping track of location and orientation

The agent program in Figure 7.19 “cheats” because it keeps track of location *outside* the knowledge base, instead of using logical reasoning.<sup>13</sup> To do it “properly,” we will need propositions for location. One’s first inclination might be to use a symbol such as  $L_{1,1}$  to mean that the agent is in  $[1,1]$ . Then the initial knowledge base might include sentences like

$$L_{1,1} \wedge \text{FacingRight} \wedge \text{Forward} \Rightarrow L_{2,1}.$$

Instantly, we see that this won’t work. If the agent starts in  $[1,1]$  facing right and moves forward, the knowledge base will entail both  $L_{1,1}$  (the original location) and  $L_{2,1}$  (the new location). Yet these propositions cannot both be true! The problem is that the location propositions should refer to two different times. We need  $L_{1,1}^1$  to mean that the agent is in  $[1,1]$  at time 1,  $L_{2,1}^2$  to mean that the agent is in  $[2,1]$  at time 2, and so on. The orientation and action propositions also need to depend on time. Therefore, the correct sentence is

$$\begin{aligned} L_{1,1}^1 \wedge \text{FacingRight}^1 \wedge \text{Forward}^1 &\Rightarrow L_{2,1}^2, \\ \text{FacingRight} \wedge \text{TurnLeft}^1 &\Rightarrow \text{FacingUp}^2, \end{aligned}$$

and so on. It turns out to be quite tricky to build a complete and correct knowledge base for keeping track of everything in the wumpus world; we will defer the full discussion until Chapter 10. The point we want to make here is that the initial knowledge base will contain sentences like the preceding two examples for every time  $t$ , as well as for every location. That is, for every time  $t$  and location  $[x, y]$ , the knowledge base contains a sentence of the form

$$L_{x,y}^t \wedge \text{FacingRight}^t \wedge \text{Forward}^t \Rightarrow L_{x+1,y}^{t+1}. \quad (7.3)$$

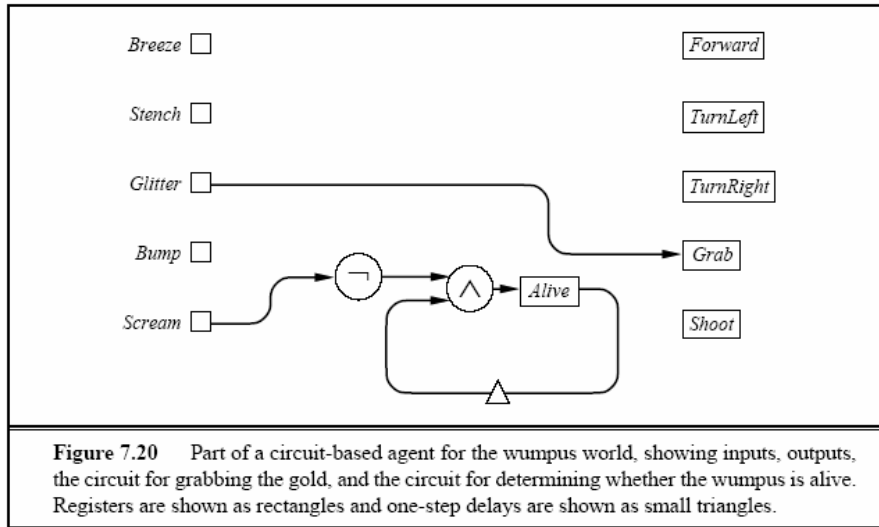
Even if we put an upper limit on the number of time steps allowed—100, perhaps—we end up with tens of thousands of sentences. The same problem arises if we add the sentences “as needed” for each new time step. This proliferation of clauses makes the knowledge base unreadable for a human, but fast propositional solvers can still handle the  $4 \times 4$  Wumpus world with ease (they reach their limit at around  $100 \times 100$ ). The circuit-based agents in the next subsection offer a partial solution to this clause proliferation problem, but the full solution will have to wait until we have developed first-order logic in Chapter 8.

### Circuit-based agents

A **circuit-based agent** is a particular kind of reflex agent with state, as defined in Chapter 2. The percepts are inputs to a **sequential circuit**—a network of **gates**, each of which implements a logical connective, and **registers**, each of which stores the truth value of a single proposition. The outputs of the circuit are registers corresponding to actions—for example,

CIRCUIT-BASED  
AGENT  
SEQUENTIAL  
CIRCUIT  
GATES  
REGISTERS

<sup>13</sup> The observant reader will have noticed that this allowed us to finesse the connection between the raw percepts such as *Breeze* and the location-specific propositions such as  $B_{1,1}$ .



the *Grab* output is set to *true* if the agent wants to grab something. If the *Glitter* input is connected directly to the *Grab* output, the agent will grab the goal whenever it sees it. (See Figure 7.20.)

DATAFLOW

Circuits are evaluated in a **dataflow** fashion: at each time step, the inputs are set and the signals propagate through the circuit. Whenever a gate has all its inputs, it produces an output. This process is closely related to the process of forward chaining in an AND-OR graph such as Figure 7.15(b).

DELAY LINE

We said in the preceding section that circuit-based agents handle time more satisfactorily than propositional inference-based agents. This is because the value stored in each register gives the truth value of the corresponding proposition symbol *at the current time  $t$* , rather than having a different copy for each different time step. For example, we might have an *Alive* register that should contain *true* when the wumpus is alive and *false* when it is dead. This register corresponds to the proposition symbol  $Alive^t$ , so on each time step it refers to a different proposition. The internal state of the agent—i.e., its memory—is maintained by connecting the output of a register back into the circuit through a **delay line**. This delivers the value of the register at the *previous* time step. Figure 7.20 shows an example. The value for *Alive* is given by the conjunction of the negation of *Scream* and the delayed value of *Alive* itself. In terms of propositions, the circuit for *Alive* implements the biconditional

$$Alive^t \Leftrightarrow \neg Scream^t \wedge Alive^{t-1} . \quad (7.4)$$

which says that the wumpus is alive at time  $t$  *if and only if* there was no scream perceived at time  $t$  (from a scream at  $t - 1$ ) *and* it was alive at  $t - 1$ . We assume that the circuit is initialized with *Alive* set to *true*. Therefore, *Alive* will remain true until there is a scream, whereupon it will become false and stay false. This is exactly what we want.



expressions!) When both  $K(B_{4,4})$  and  $K(\neg B_{4,4})$  are false, it means the truth value of  $B_{4,4}$  is unknown. (If both are true, there's a bug in the knowledge base!) Now whenever we would use  $B_{4,4}$  in some part of the circuit, we use  $K(B_{4,4})$  instead; and whenever we would use  $\neg B_{4,4}$ , we use  $K(\neg B_{4,4})$ . In general, we represent each potentially indeterminate proposition with two **knowledge propositions** that state whether the underlying proposition is known to be true and known to be false.

We will see an example of how to use knowledge propositions shortly. First, we need to work out how to determine the truth values of the knowledge propositions themselves. Notice that, whereas  $B_{4,4}$  has a fixed truth value,  $K(B_{4,4})$  and  $K(\neg B_{4,4})$  *do* change as the agent finds out more about the world. For example,  $K(B_{4,4})$  starts out false and then becomes true as soon as  $B_{4,4}$  can be determined to be true—that is, when the agent is in  $[4,4]$  and detects a breeze. It stays true thereafter. So we have

$$K(B_{4,4})^t \Leftrightarrow K(B_{4,4})^{t-1} \vee (L_{4,4}^t \wedge \text{Breeze}^t). \quad (7.6)$$

A similar equation can be written for  $K(\neg B_{4,4})^t$ .

Now that the agent knows about breezy squares, it can deal with pits. The absence of a pit in a square can be ascertained if and only if one of the neighboring squares is known not to be breezy. For example, we have

$$K(\neg P_{4,4})^t \Leftrightarrow K(\neg B_{3,4})^t \vee K(\neg B_{4,3})^t. \quad (7.7)$$

Determining that there *is* a pit in a square is more difficult—there must be a breeze in an adjacent square that cannot be accounted for by another pit:

$$\begin{aligned} K(P_{4,4})^t &\Leftrightarrow (K(B_{3,4})^t \wedge K(\neg P_{2,4})^t \wedge K(\neg P_{3,3})^t) \\ &\vee (K(B_{4,3})^t \wedge K(\neg P_{4,2})^t \wedge K(\neg P_{3,3})^t). \end{aligned} \quad (7.8)$$



While the circuits for determining the presence or absence of pits are somewhat hairy, *they have only a constant number of gates for each square*. This property is essential if we are to build circuit-based agents that scale up in a reasonable way. It is really a property of the wumpus world itself; we say that an environment exhibits **locality** if the truth of each proposition of interest can be determined looking only at a constant number of other propositions. Locality is very sensitive to the precise “physics” of the environment. For example, the minesweeper domain (Exercise 7.11) is nonlocal because determining that a mine is in a given square can involve looking at squares arbitrarily far away. For nonlocal domains, circuit-based agents are not always practical.

There is one issue around which we have tiptoed carefully: the question of **acyclicity**. A circuit is acyclic if every path that connects the output of a register back to its input has an intervening delay element. We require that all circuits be acyclic because cyclic circuits, as physical devices, do not work! They can go into unstable oscillations resulting in undefined values. As an example of a cyclic circuit, consider the following augmentation of Equation (7.6):

$$K(B_{4,4})^t \Leftrightarrow K(B_{4,4})^{t-1} \vee (L_{4,4}^t \wedge \text{Breeze}^t) \vee K(P_{3,4})^t \vee K(P_{4,3})^t. \quad (7.9)$$

The extra disjuncts,  $K(P_{3,4})^t$  and  $K(P_{4,3})^t$ , allow the agent to determine breeziness from the known presence of adjacent pits, which seems entirely reasonable. Now, unfortunately,

breeziness depends on adjacent pits, and pits depend on adjacent breeziness through equations such as Equation (7.8). Therefore, the complete circuit would contain cycles.

The difficulty is not that the augmented Equation (7.9) is *incorrect*. Rather, the problem is that the interlocking dependencies represented by these equations cannot be resolved by the simple mechanism of propagating truth values in the corresponding Boolean circuit. The acyclic version using Equation (7.6), which determines breeziness only from direct observation, is *incomplete* in the sense that at some points the circuit-based agent might know less than an inference-based agent using a complete inference procedure. For example, if there is a breeze in [1,1], the inference-based agent can conclude that there is also a breeze in [2,2], whereas the acyclic circuit-based agent using Equation (7.6) cannot. A complete circuit *can* be built—after all, sequential circuits can emulate any digital computer—but it would be significantly more complex.

## A comparison

The inference-based agent and the circuit-based agent represent the declarative and procedural extremes in agent design. They can be compared along several dimensions:

- *Conciseness*: The circuit-based agent, unlike the inference-based agent, need not have separate copies of its “knowledge” for every time step. Instead, it refers only to the current and previous time steps. Both agents need copies of the “physics” (expressed as sentences or circuits) for every square and therefore do not scale well to larger environments. In environments with many objects related in complex ways, the number of propositions will swamp any propositional agent. Such environments require the expressive power of first-order logic. (See Chapter 8.) Propositional agents of both kinds are also poorly suited for expressing or solving the problem of finding a path to a nearby safe square. (For this reason, PL-WUMPUS-AGENT falls back on a search algorithm.)
- *Computational efficiency*: In the *worst* case, inference can take time exponential in the number of symbols, whereas evaluating a circuit takes time linear in the size of the circuit (or linear in the *depth* of the circuit if realized as a physical device). In *practice*, however, we saw that DPLL completed the required inferences very quickly.<sup>14</sup>
- *Completeness*: We suggested earlier that the circuit-based agent might be incomplete because of the acyclicity restriction. The reasons for incompleteness are actually more fundamental. First, remember that a circuit executes in time linear in the circuit size. This means that, for some environments, a circuit that is complete (i.e., one that computes the truth value of every determinable proposition) must be exponentially larger than the inference-based agent’s KB. Otherwise, we would have a way to solve the propositional entailment problem in less than exponential time, which is very unlikely. A second reason is the nature of the internal state of the agent. The inference-based agent remembers every percept and knows, either implicitly or explicitly, every sentence that follows from the percepts and initial KB. For example, given  $B_{1,1}$ , it knows the disjunction  $P_{1,2} \vee P_{2,1}$ , from which  $B_{2,2}$  follows. The circuit-based agent, on the

<sup>14</sup> In fact, all the inferences done by a circuit can be done in linear time by DPLL! This is because evaluating a circuit, like forward chaining, can be emulated by DPLL using the unit propagation rule.



other hand, forgets all previous percepts and remembers just the individual propositions stored in registers. Thus,  $P_{1,2}$  and  $P_{2,1}$  remain *individually* unknown after the first percept, so no conclusion will be drawn about  $B_{2,2}$ .

- *Ease of construction*: This is a very important issue about which it is hard to be precise. Certainly, this author found it much easier to state the “physics” declaratively, whereas devising small, acyclic, not-too-incomplete circuits for direct detection of pits seemed quite difficult.

In sum, it seems there are *tradeoffs* among computational efficiency, conciseness, completeness, and ease of construction. When the connection between percepts and actions is simple—as in the connection between *Glitter* and *Grab*—a circuit seems optimal. For more complex connections, the declarative approach may be better. In a domain such as chess, for example, the declarative rules are concise and easily encoded (at least in first-order logic), but a circuit for computing moves directly from board states would be unimaginably vast.

We see different points on these tradeoffs in the animal kingdom. The lower animals with very simple nervous systems are probably circuit-based, whereas higher animals, including humans, seem to perform inference on explicit representations. This enables them to compute much more complex agent functions. Humans also have circuits to implement reflexes, and perhaps also **compile** declarative representations into circuits when certain inferences become routine. In this way, a **hybrid agent** design (see Chapter 2) can have the best of both worlds.

COMPILATION

## EXERCISES

7.1 Describe the wumpus world according to the properties of task environments listed in Chapter 2.

7.2 Suppose the agent has progressed to the point shown in Figure 7.4(a), having perceived nothing in [1,1], a breeze in [2,1], and a stench in [1,2]. and is now concerned with the contents of [1,3], [2,2], and [3,1]. Each of these can contain a pit and at most one can contain a wumpus. Following the example of Figure 7.5, construct the set of possible worlds. (You should find 32 of them.) Mark the worlds in which the KB is true and those in which each of the following sentences is true:

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

$\alpha_3$  = “There is a wumpus in [1,3].”

Hence show that  $KB \models \alpha_2$  and  $KB \models \alpha_3$ .

7.3 Consider the problem of deciding whether a propositional logic sentence is true in a given model.

- Write a recursive algorithm  $PL\text{-}TRUE?(s, m)$  that returns *true* if and only if the sentence  $s$  is true in the model  $m$  (where  $m$  assigns a truth value for every symbol in  $s$ ). The algorithm should run in time linear in the size of the sentence. (Alternatively, use a version of this function from the online code repository.)

- b. Give three examples of sentences that can be determined to be true or false in a *partial* model that does not specify a truth value for some of the symbols.
- c. Show that the truth value (if any) of a sentence in a partial model cannot be determined efficiently in general.
- d. Modify your PL-TRUE? algorithm so that it can sometimes judge truth from partial models, while retaining its recursive structure and linear runtime. Give three examples of sentences whose truth in a partial model is *not* detected by your algorithm.
- e. Investigate whether the modified algorithm makes TT-ENTAILS? more efficient.

7.4 Prove each of the following assertions:

- a.  $\alpha$  is valid if and only if  $\text{True} \models \alpha$ .
- b. For any  $\alpha$ ,  $\text{False} \models \alpha$ .
- c.  $\alpha \models \beta$  if and only if the sentence  $(\alpha \Rightarrow \beta)$  is valid.
- d.  $\alpha \equiv \beta$  if and only if the sentence  $(\alpha \Leftrightarrow \beta)$  is valid.
- e.  $\alpha \models \beta$  if and only if the sentence  $(\alpha \wedge \neg\beta)$  is unsatisfiable.

7.5 Consider a vocabulary with only four propositions,  $A$ ,  $B$ ,  $C$ , and  $D$ . How many models are there for the following sentences?

- a.  $(A \wedge B) \vee (B \wedge C)$
- b.  $A \vee B$
- c.  $A \Leftrightarrow B \Leftrightarrow C$

7.6 We have defined four different binary logical connectives.

- a. Are there any others that might be useful?
- b. How many binary connectives can there be?
- c. Why are some of them not very useful?

7.7 Using a method of your choice, verify each of the equivalences in Figure 7.11.

7.8 Decide whether each of the following sentences is valid, unsatisfiable, or neither. Verify your decisions using truth tables or the equivalence rules of Figure 7.11.

- a.  $\text{Smoke} \Rightarrow \text{Smoke}$
- b.  $\text{Smoke} \Rightarrow \text{Fire}$
- c.  $(\text{Smoke} \Rightarrow \text{Fire}) \Rightarrow (\neg\text{Smoke} \Rightarrow \neg\text{Fire})$
- d.  $\text{Smoke} \vee \text{Fire} \vee \neg\text{Fire}$
- e.  $((\text{Smoke} \wedge \text{Heat}) \Rightarrow \text{Fire}) \Leftrightarrow ((\text{Smoke} \Rightarrow \text{Fire}) \vee (\text{Heat} \Rightarrow \text{Fire}))$
- f.  $(\text{Smoke} \Rightarrow \text{Fire}) \Rightarrow ((\text{Smoke} \wedge \text{Heat}) \Rightarrow \text{Fire})$
- g.  $\text{Big} \vee \text{Dumb} \vee (\text{Big} \Rightarrow \text{Dumb})$
- h.  $(\text{Big} \wedge \text{Dumb}) \vee \neg\text{Dumb}$

7.9 (Adapted from Barwise and Etchemendy (1993).) Given the following, can you prove that the unicorn is mythical? How about magical? Horned?

If the unicorn is mythical, then it is immortal, but if it is not mythical, then it is a mortal mammal. If the unicorn is either immortal or a mammal, then it is horned. The unicorn is magical if it is horned.

7.10 Any propositional logic sentence is logically equivalent to the assertion that each possible world in which it would be false is not the case. From this observation, prove that any sentence can be written in CNF.

7.11 Minesweeper, the well-known computer game, is closely related to the wumpus world. A minesweeper world is a rectangular grid of  $N$  squares with  $M$  invisible mines scattered among them. Any square may be probed by the agent; instant death follows if a mine is probed. Minesweeper indicates the presence of mines by revealing, in each probed square, the *number* of mines that are directly or diagonally adjacent. The goal is to have probed every unmined square.

- Let  $X_{i,j}$  be true iff square  $[i, j]$  contains a mine. Write down the assertion that there are exactly two mines adjacent to  $[1,1]$  as a sentence involving some logical combination of  $X_{i,j}$  propositions.
- Generalize your assertion from (a) by explaining how to construct a CNF sentence asserting that  $k$  of  $n$  neighbors contain mines.
- Explain precisely how an agent can use DPLL to prove that a given square does (or does not) contain a mine, ignoring the global constraint that there are exactly  $M$  mines in all.
- Suppose that the global constraint is constructed via your method from part (b). How does the number of clauses depend on  $M$  and  $N$ ? Suggest a way to modify DPLL so that the global constraint does not need to be represented explicitly.
- Are any conclusions derived by the method in part (c) invalidated when the global constraint is taken into account?
- Give examples of configurations of probe values that induce *long-range dependencies* such that the contents of a given unprobed square would give information about the contents of a far-distant square. [Hint: consider an  $N \times 1$  board.]

7.12 This exercise looks into the relationship between clauses and implication sentences.

- Show that the clause  $(\neg P_1 \vee \dots \vee \neg P_m \vee Q)$  is logically equivalent to the implication sentence  $(P_1 \wedge \dots \wedge P_m) \Rightarrow Q$ .
- Show that every clause (regardless of the number of positive literals) can be written in the form  $(P_1 \wedge \dots \wedge P_m) \Rightarrow (Q_1 \vee \dots \vee Q_n)$ , where the  $P$ s and  $Q$ s are proposition symbols. A knowledge base consisting of such sentences is in **implicative normal form** or **Kowalski form**.
- Write down the full resolution rule for sentences in implicative normal form.

IMPlicative  
NORMAL FORM

- 7.13 In this exercise, you will design more of the circuit-based wumpus agent.
- Write an equation, similar to Equation (7.4), for the *Arrow* proposition, which should be true when the agent still has an arrow. Draw the corresponding circuit.
  - Repeat part (a) for *FacingRight*, using Equation (7.5) as a model.
  - Create versions of Equations 7.7 and 7.8 for finding the wumpus, and draw the circuit.
- 7.14 Discuss what is meant by *optimal* behavior in the wumpus world. Show that our definition of the PL-WUMPUS-AGENT is not optimal, and suggest ways to improve it.
- 7.15 Extend PL-WUMPUS-AGENT so that it keeps track of all relevant facts *within* the knowledge base.
- 7.16 How long does it take to prove  $KB \models \alpha$  using DPLL when  $\alpha$  is a literal *already contained in KB*? Explain.
- 7.17 Trace the behavior of DPLL on the knowledge base in Figure 7.15 when trying to prove  $Q$ , and compare this behavior with that of the forward chaining algorithm.

