# Logical Foundations
## of
# Artificial Intelligence

**Michael R. Genesereth and Nils J. Nilsson**
*Stanford University*

# CHAPTER 6
# Nonmonotonic Reasoning

WE HAVE ALREADY SEEN some indication of the power of the first-order predicate calculus as a language for expressing declarative knowledge in AI systems. We can use the predicate calculus to express any conceptualization based on objects and their relations in a domain of discourse. Given what we have presented so far, we might imagine that a typical AI system using first-order logic would work somewhat as follows. Information that the system has about its domain is expressed as a finite set $\Delta$ of first-order formulas. We call $\Delta$ the database or *base set of beliefs* of the system. To answer queries or to take appropriate actions, the system typically will have to decide whether or not some formula $\phi$ is logically entailed by its beliefs. We can imagine that the system will make this decision by performing logical deductions on $\Delta$, perhaps by using resolution on the clause form of $\Delta \wedge \neg \phi$. (Our notation is simplified by letting $\Delta$ also stand for the conjunction of the formulas in the set $\Delta$.)

Even though this model is quite useful for a variety of tasks requiring knowledge about a domain, it has major limitations. The three most important ones are:

(1) Language (probably *any* language) cannot capture all that we want to say about the world. A finite set of sentences can never be more than an approximate description of things as they really are. Any general rule that we might care to frame is subject to an unlimited number of exceptions and qualifications. If we are going to use language to describe the world, we will have to use it in a way that is robust in the face of an ever-expanding set of more highly articulated statements.

(2) The inference rules of ordinary logic (such as modus ponens and the resolution principle) are *sound*. Thus, deductions from a base set of beliefs never produce *new* knowledge about the world. If $\phi$ logically follows from $\Delta$, then all the models of $\Delta$, including our intended interpretation, also are models of $\phi$. Deriving $\phi$ in no way eliminates any of the models and thus $\phi$ tells us nothing about the world that was not already described by $\Delta$. Of course, we want to be able to manipulate our description of the world so that implicit facts about it become represented explicitly, and sound inference rules do just that for us. We also will want, however, to add formulas to $\Delta$ that say new (or revised) things about the world, and ordinary logic gives us no hints about how to do that. We need methods for reasoning with tentative statements because they are the only ones we are ever going to have. These reasoning methods will have to anticipate the possibility of later revisions of the knowledge base.

(3) The logical languages that we have used so far are adequate for expressing only those statements that we are willing to take as being either wholly true or wholly false. Often, we have information about a situation that is known to be *uncertain*. For example, we know that it is likely (but not certain) that it will be sunny in Pasadena on New Year's Day.

In the next few chapters, we will be concerned with confronting and overcoming some of these limitations. An important technique involves the use of *nonsound* inferences of various kinds. That is, from a database $\Delta$, we will allow certain inferences that do not logically follow from $\Delta$. Often, these inferences depend globally on *all* of the sentences in $\Delta$ rather than on a small subset. In particular, we will be introducing inference techniques the application of which depends on certain sentences *not* being in $\Delta$. With such inference rules, if another sentence is added to $\Delta$, an inference may have to be retracted. For this reason, such inference rules are called *nonmonotonic*. Ordinary logical inference rules, on the other hand, are monotonic because the set of theorems derivable from premises is not reduced by adding to the premises.

There are many situations in which it is appropriate for intelligent systems to augment their beliefs by new ones that do not logically follow from their explicit ones. The press of events sometimes compels *some* action before all the relevant facts are at hand. It would be useful perhaps for systems to be able to assume that the beliefs they have presently about a certain subject are all the beliefs that are important about that subject. Natural-language dialogues among humans, for example, depend on both the speaker and the hearer using several general augmenting conventions of this sort. (Example: "He didn't say John was his brother, so I'll assume he is not.")

Also, as we mentioned earlier, there is a sense in which any attempt to capture *all* the knowledge about the real world by a finite set of sentences is fundamentally impossible. As our own knowledge (and that of science) increases, our conceptualization of a subject area changes. Any conceptualization, put forward for certain purposes, is subject to challenge. Consider, for example, the following sentence about birds: "All birds fly." With the obvious intended interpretation, we might express this as ∀x Bird(x) ⇒ Flies(x). This sentence might be useful for certain limited purposes, but if we tried to apply it more generally, we might be confronted by the fact that ostriches, which are birds, do not fly. After having this problem pointed out to us, we could correct it by changing our axiom about birds flying to the following:

$$\forall x \; \text{Bird}(x) \wedge \neg \text{Ostrich}(x) \Rightarrow \text{Flies}(x)$$

Even this sentence does not accurately capture the *real* world, however, because we can imagine several other kinds of birds that do not fly: baby birds, dead birds, wingless birds, and so on. The list of such *qualifications* is very long if not endless, leading us, perhaps, to despair of using language for knowledge representation. This problem often is called the *qualification problem.* Most universally quantified sentences will have to include an infinite number of qualifications if they are to be interpreted as accurate statements about the world. Yet, in our everyday reasoning, we humans use sentences that we assume to be true. What we would seem to need for our machines is an inference rule that permits us to make somewhat temporary, or *default*, assumptions that can later be revised when additional qualifications become important.

There are a number of ways to achieve the appropriate nonmonotonic effects. In this chapter, we explore three methods. In one, we adopt a convention that allows us to infer that, if we cannot prove a ground atom, then we can assume its negation. In another, we show how to compute a formula that can be added to Δ that restricts the objects that satisfy a certain predicate to just those that Δ says must satisfy that predicate. In a third, we introduce nonmonotonic rules of inference called *defaults*, and show how they are used to derive default conclusions.

These methods have several potential applications. Through examples in this and the next chapter, we indicate their utility for making reasonable assumptions about what can be inferred from a finite set of sentences. We see these nonmonotonic techniques as promising candidates for extending the reach of logic beyond some of the limitations we have noted.

## 6.1 The Closed-World Assumption

Recall that a theory $T$ is *complete* if either every ground atom in the language or its negation is in the theory. Thus, the logical closure of

the formula P(A) ∧ (P(A) ⇒ Q(A)) ∧ P(B) is not a complete theory, because neither Q(B) nor ¬Q(B) is in the theory. One convention for augmenting a theory is to complete it.

The most straightforward and simple way to complete a theory is by a convention called the *closed-world assumption (CWA)*. The CWA completes the theory defined by a base set of beliefs, $\Delta$, by including the *negation* of a ground atom in the completed theory whenever that ground atom does not logically follow from $\Delta$. The effect of the CWA is as though we augmented the base belief set with all the negative ground literals the positive versions of which cannot be deduced from $\Delta$. The CWA is nonmonotonic because the set of augmented beliefs would shrink if we added a new positive ground literal to $\Delta$.

We define the effects of the CWA in terms of customary logical notation. We call our belief set, $\Delta$, the *proper axioms* of a theory. The theory, denoted by $T[\Delta]$, is the closure of $\Delta$ under logical entailment. The CWA *augments* $T[\Delta]$ by adding a set, $\Delta_{asm}$, of *assumed beliefs*. The closure under logical implication of the union of these assumed beliefs with $\Delta$ comprises the CWA-augmented set of beliefs CWA$[\Delta]$. The CWA can be stated succinctly as follows:

- The formula $\phi$ (constructed from elements of a predefined predicate-calculus language) is in $T[\Delta]$ if and only if $\Delta \models \phi$. (This is the ordinary definition of a theory $T[\Delta]$ in terms of a base set $\Delta$.)

- $\neg P$ is in $\Delta_{asm}$ if and only if the ground atom $P$ is not in $T[\Delta]$. ($\Delta_{asm}$ is a set of added beliefs assumed under the CWA.)

- $\phi$ is in CWA$[\Delta]$ if and only if $\{\Delta \cup \Delta_{asm}\} \models \phi$. (The augmented theory, CWA$[\Delta]$, is the closure of all beliefs, explicit and assumed.)

In our example, in which $\Delta$ is P(A) ∧ (P(A) ⇒ Q(A)) ∧ P(B), the CWA adds ¬Q(B), since $\Delta$ does not logically entail Q(B).

The CWA often is used with database systems. Suppose we have a database listing pairs of countries that are geographic neighbors:

Neighbor(US,Canada)

Neighbor(US,Mexico)

Neighbor(Mexico,Guatemala)

⋮

For such a database, it would be useful to adopt the convention that countries that are not listed as neighbors are not neighbors. That convention is an example of the CWA. Without such a convention, we would have to list explictly all the nonneighbor pairs if we wanted to be able to answer queries such as, "Are Brazil and Canada neighbors?"

Note that the CWA depends on a syntactic feature of a set of beliefs; namely, whether a *positive* ground literal can be derived. If we systematically replaced each predicate letter $P_i$ by $\neg Q_i$ (defining $P_i \equiv \neg Q_i$) the theory would be the same, but the CWA would give different results with respect to the original predicates. The convention is most efficient when the "positive facts" are few in number compared to the "negative facts." A database designer who uses the CWA will want to conceptualize the domain in a way that matches this expectation.

We might ask whether the CWA always results in a consistent augmented theory, CWA[$\Delta$]. The following example shows that it does not.

Let $\Delta$ contain only the clause P(A) $\vee$ P(B). Then neither P(A) nor P(B) is in $T[\Delta]$, so by the CWA their negations are *both* in CWA[$\Delta$]. Together, however, these negations are not consistent with P(A) $\vee$ P(B).

The source of this difficulty is that $\Delta$ contained a disjunction of ground atoms (positive ground literals) but no way to prove any of them. Thus, the conjunction of their negations, which contradicts the original disjunction, is in the augmented theory. The following theorem links this difficulty with the inconsistency of CWA[$\Delta$].

**THEOREM 6.1** CWA[$\Delta$] *is consistent if and only if, for every positive-ground-literal clause $L_1 \vee L_2 \vee \ldots \vee L_n$ that follows from $\Delta$, there is also entailed by $\Delta$ at least one ground literal $L_i$ that subsumes it. (Equivalently, the CWA augmentation* CWA[$\Delta$] *of a consistent $\Delta$ is inconsistent if and only if there are positive ground literals $L_1, \ldots, L_n$ such that $\Delta \models L_1 \vee L_2 \vee \ldots \vee L_n$ but, for $i = 1, \ldots, n$, $\Delta \not\models L_i$.)*

**Proof** CWA[$\Delta$] can be inconsistent only if $\Delta \cup \Delta_{asm}$ is. Then, by the compactness theorem of logic, there is a finite subset of $\Delta_{asm}$ that contradicts $\Delta$. Let this subset be $\{\neg L_1, \ldots, \neg L_n\}$. Then $\Delta$ implies the negation of the conjunction of these formulas; that is, $\Delta \models L_1 \vee \ldots \vee L_n$. Since each $\neg L_i$ is in $\Delta_{asm}$, by the definition of $\Delta_{asm}$, none of the $L_i$ follow from $\Delta$. The proof in the other direction is obvious. $\square$

The application of Theorem 6.1 depends critically on the terms that we allow as part of the language. For example, if the only object constants in the language are A and B, then the following clauses do not have an inconsistent augmentation (even though one of them is a disjunction of positive literals):

$$P(x) \ \vee \ Q(x)$$

$$P(A)$$

$$Q(B)$$

In this case, the only ground clauses of the form $L_1 \vee L_2 \vee \ldots \vee L_n$ that can be proved from $\Delta$ are P(A) $\vee$ Q(A) and P(B) $\vee$ Q(B) (by universal instantiation). Each of these is subsumed by one of the clauses in $\Delta$. On the other hand, if we also admit the object constant C, we also can prove P(C) $\vee$ Q(C), but we can prove neither P(C) nor Q(C) to subsume it, so the CWA produces an inconsistent augmentation in that case.

In the first case of this example, we limited the object constants of the language to those that occurred in $\Delta$. We sometimes also want to make the assumption that the only objects in the domain are the ones that can be named using the object and function constants occurring in the language. This is called the *domain-closure assumption (DCA)*. If there are no function constants in the language, the DCA can be written as the following axiom, the domain-closure axiom:

$$\forall x \; x = t_1 \vee x = t_2 \vee \ldots$$

where the $t_i$ are the object constants of the language. (If the language contained function constants, there would be an infinite number of terms that could be constructed, and typically the DCA could not be expressed by a first-order formula.) This axiom makes a strong assumption. It allows us, for example, to replace any quantifiers by finite conjunctions and disjunctions; then the belief set would be equivalent to a propositional combination of ground literals.

Another assumption often used in connection with nonmonotonic reasoning is the *unique-names assumption (UNA)*: If ground terms cannot be proved equal, they can be assumed unequal. The UNA is a consequence of the CWA; it is merely an application of the CWA to the equality predicate. The DCA is sometimes used in addition to the CWA to restrict the augmentation further.

Since it may be difficult to test the conditions of Theorem 6.1, the following corollary is important. (Recall that a *Horn clause* is defined to be one that has at most one positive literal.)

**COROLLARY 6.1** *If the clause form of $\Delta$ is Horn and consistent, then the CWA augmentation* CWA[$\Delta$] *is consistent.*

**Proof** Suppose the contrary; that is, that $\Delta$ is Horn and consistent but that CWA[$\Delta$] is inconsistent. Then, according to Theorem 6.1, we can deduce from $\Delta$ a ground clause $L_1 \vee L_2 \vee \ldots \vee L_n$ containing only positive ground literals no one of which is derivable from $\Delta$. Thus, $\Delta \cup \{\neg L_1, \ldots, \neg L_n\}$ is inconsistent. However, because $\Delta$ contains only Horn clauses, it must then be the case that for some $i$, $\Delta \wedge \neg L_i$ is inconsistent (see Exercise 3). Or, for some $i$, $\Delta \models L_i$. But this contradicts the choice of $L_i$. $\square$

Thus we see that an important class of theories—the so-called Horn theories—have consistent CWA augmentations. We see from Theorem 6.1, however, that the condition that $\Delta$ be Horn is not absolutely necessary for the CWA augmentation of $\Delta$ to be consistent.

The CWA is too strong for many applications. We do not always want to assume that *any* ground atom not provable from $\Delta$ is false. Weakening this assumption in the natural way leads to the idea of the CWA *with respect to a predicate P*. Under that convention, ground atoms in some particular predicate, $P$, that are not provable from $\Delta$ are assumed to be false. The assumed beliefs, $\Delta_{asm}$, in that case contain only negative ground literals in $P$.

For example, suppose $\Delta$ is:

$$\forall x \ Q(x) \ \Rightarrow \ P(x)$$

$$Q(A)$$

$$R(B) \ \vee \ P(B)$$

Applying the CWA to $\Delta$ with respect to P allows us to conclude ¬P(B), because P(B) cannot be concluded from $\Delta$. That, in turn, allows us also to conclude R(B) from $\Delta$. (Unconstrained application of the CWA to $\Delta$ would have condoned concluding both ¬R(B) and ¬P(B), which contradict $\Delta$.)

We also can make the CWA with respect to a *set* of predicates. In database applications, this assumption allows us to assume that certain relations in the database are complete and others are not. If the set contains all the predicates in $\Delta$, then we get the same result as we would have got with the ordinary CWA.

It is interesting to note that the CWA with respect to a set of predicates may produce an inconsistent augmentation even when the CWA with respect to each one of the predicates alone produces a consistent augmentation. For example, the CWA with respect to the set {P,Q} is inconsistent with the belief set (PvQ), even though the CWA with respect to either P or Q is consistent with that belief set.

We might be tempted to say that the source of this difficulty is that (P ∨ Q) is not Horn in the set {P,Q}. (We say that a set of clauses is Horn in the predicate $P$ if there is at most one positive occurrence of $P$ in each clause. We say that a set $\Delta$ of clauses is Horn in a set of predicates $\Pi$ if and only if each of the clauses would be Horn in $P$ after substituting the letter $P$ for each occurence of any of the letters in $\Pi$ in the clauses of $\Delta$.) Even if a belief set were Horn in a set of predicates, however, the CWA with respect to the predicates in this set might produce an inconsistent augmentation. Suppose $\Delta$ is {P(A)vQ, P(B)v¬Q}. This set is Horn in the set {P}, and making the CWA with respect to the predicates (just P) in {P} yields both ¬P(A) and ¬P(B). These, together with $\Delta$, are inconsistent.

## 6.2 Predicate Completion

It happens that we often can express in a single sentence of logic the assumption that the only objects that satisfy a predicate are those that *must* do so—given our beliefs. We will describe several of these methods—all related but of increasing generality and power.

First consider the simple case in which P(A) is the *only* formula in $\Delta$. P(A) is equivalent to the following expression:

$$\forall x \ x=A \ \Rightarrow \ P(x)$$

Such a formula could be taken to be the "if" half of a *definition* for P. The assumption that there are no other objects that satisfy P can then be made by writing the "only if" half as

$$\forall x \ P(x) \ \Rightarrow \ x=A$$

This half is called the *completion formula* for P. It makes the explicit information about P in $\Delta$ *complete*.

The conjunction of $\Delta$ with the completion formula is called *the completion of* P *in* $\Delta$, and is denoted by COMP[$\Delta$; P]. In this case, we have

$$COMP[\Delta;P] \equiv (\forall x \ P(x) \ \Rightarrow \ x=A) \wedge \Delta$$
$$\equiv \ \forall x \ P(x) \ \Leftrightarrow \ x=A$$

In this example, predicate completion (if augmented by the UNA) produces the same effect as the CWA with respect to P.

If $\Delta$ contained only two formulas in P, say P(A) and P(B), the completion formula would be

$$\forall x \ P(x) \ \Rightarrow \ x=A \ \vee \ x=B$$

Here again, predicate completion of P (together with unique names) has the same effect as the CWA with respect to P.

If $\Delta$ contains formulas in which a predicate $P$ occurs disjunctively with other predicates or in which $P$ contains variables, predicate completion is more complex. In fact, we define predicate completion for only certain kinds of clauses.

We say that a set of clauses is *solitary* in $P$ if each clause with a positive occurrence of $P$ has at most *one* occurrence of $P$. Note that clauses solitary in $P$ are also Horn in $P$, but not necessarily vice versa. For example Q(A) $\vee$ ¬P(B) $\vee$ P(A) is Horn in P, but not solitary in P.

We define predicate completion of $P$ only for clauses solitary in $P$. Suppose $\Delta$ is a set of clauses solitary in $P$. We can write each of the clauses in $\Delta$ that contains a positive $P$ literal in the following form:

$$\forall y \ Q_1 \wedge \ldots \wedge Q_m \Rightarrow P(t)$$

where $t$ is a tuple of terms, $[t_1, t_2, \ldots, t_n]$, and the $Q_i$ are literals not containing $P$. There may be no $Q_i$, in which case the clause is simply $P(t)$. The $Q_i$ and $t$ may contain variables, say the tuple of variables y.

This expression is equivalent to

$$\forall y \forall x \ (x=t) \wedge Q_1 \wedge \ldots \wedge Q_m \Rightarrow P(x)$$

where x is a tuple of variables not occurring in $t$ and $(x=t)$ is an abbreviation for $(x_1=t_1 \wedge \ldots \wedge x_n=t_n)$. Finally, since the variables y occur in only the antecedent of the implication, this expression is equivalent to

$$\forall x \ (\exists y \ (x=t) \wedge Q_1 \wedge \ldots \wedge Q_m) \Rightarrow P(x)$$

This way of writing the clause is called the *normal form* of the clause. Suppose there are exactly $k$ clauses $(k > 0)$ in $\Delta$ that have a positive $P$ literal. Let the normal forms of these clauses be

$$\forall x \ E_1 \Rightarrow P(x)$$
$$\forall x \ E_2 \Rightarrow P(x)$$
$$\vdots$$
$$\forall x \ E_k \Rightarrow P(x)$$

Each of the $E_i$ will be an existentially quantified conjunction of literals, as in the preceding generic case. If we group together these clauses as a single implication, we obtain

$$\forall x \ E_1 \vee E_2 \vee \ldots \vee E_k \Rightarrow P(x)$$

Here we have an expression that can be taken as the "if" half of a definition for $P$. It suggests the following, "only if," *completion formula* for $P$:

$$\forall x \ P(x) \Rightarrow E_1 \vee E_2 \vee \ldots \vee E_k$$

Since the $E_i$ do not contain $P$, the "if" and "only if" parts together can be thought of as a *definition* for $P$:

$$\forall x \ P(x) \Leftrightarrow E_1 \vee E_2 \vee \ldots \vee E_k$$

Since the "if" part already is entailed by $\Delta$, we define the *completion of P in $\Delta$* as

$$\text{COMP}[\Delta; P] \equiv_{def} \Delta \wedge (\forall x \ P(x) \Leftrightarrow E_1 \vee \ldots \vee E_k)$$

where the $E_i$ are the antecedents of the normal forms of the clauses in $\Delta$, as defined previously.

Let us consider a simple example of predicate completion. Suppose $\Delta$ is

$$\forall x \ \text{Ostrich(x)} \Rightarrow \text{Bird(x)}$$
$$\text{Bird(Tweety)}$$
$$\neg\text{Ostrich(Sam)}$$

(All ostriches are birds; Tweety is a bird; Sam is not an ostrich.) We note that $\Delta$ is solitary in Bird. Let us complete Bird in $\Delta$. Writing those clauses containing Bird in normal form yields

$$\forall x \ \text{Ostrich(x)} \vee \text{x=Tweety} \Rightarrow \text{Bird(x)}$$

The completion of Bird in $\Delta$ is then simply

$$\text{COMP}[\Delta; \text{Bird}] \equiv \Delta \wedge (\forall x \ \text{Bird(x)} \Leftrightarrow \text{Ostrich(x)} \vee \text{x=Tweety})$$

(The only birds are ostriches or Tweety.) With the completion formula added to $\Delta$ (and augmented by the UNA) we could prove, for example, $\neg\text{Bird(Sam)}$.

What is predicate completion doing for us in this case? $\Delta$ tells us that Tweety is a bird, that Sam is not an ostrich, and that all ostriches are birds. Completion of Bird in $\Delta$ is a way of making the assumption that there are no birds other than those about which $\Delta$ tells us. That is, the only birds are Tweety and ostriches. Since Sam is not an ostrich, and since the UNA lets us assume that Sam is not Tweety, we can conclude that Sam is not a bird.

If we did not limit $\Delta$ to clauses solitary in $P$, the completion process might produce circular definitions for $P$, which would not restrict the objects that satisfy $P$ to those that must do so, given $\Delta$. Sometimes, we can formally apply the completion process to clauses Horn (but not solitary) in $P$ and still get meaningful results. Consider the following Horn

clauses that describe the factorial relation (we assume implicit universal quantification):

```
x=0  ⇒  Factorial(x,1)
```

$\forall x\ (\neg0x \land x=\text{?}\omega) \supset \mathcal{B}x$

```
x≠0 ∧ Factorial(Minus(x,1),y) ⇒ Factorial(x,Times(x,y))
```

Writing these expressions in normal form yields

```
x=0 ∧ z=1  ⇒  Factorial(x,z)

(∃y x≠0 ∧ z=Times(x,y) ∧ Factorial(Minus(x,1),y))  ⇒
   Factorial(x,z)
```

Now we formally perform predicate completion on `Factorial` (even though the clauses are not solitary in `Factorial`). The result is

```
Factorial(x,z)  ⇒
   (x=0 ∧ z=1) ∨
   (∃y x≠0 ∧ z=Times(x,y) ∧ Factorial(x-1,y))
```

$\forall x\ \exists$

This result is easily interpreted as a recursive definition of factorial. It illustrates that limiting predicate completion to solitary clauses is perhaps unnecessarily restrictive. Not all definitions of a predicate in terms of itself are circular—some are recursive.

There are two special cases of predicate completion that give rise to interesting forms for the completion formulas. Suppose $\Delta$ is of the form $(\forall x\ P(x))$. Making use of the atom T, we can write this clause as $(\forall x\ T \Rightarrow P(x))$. Completion of $P$ then gives the completion formula $(\forall x\ P(x) \Rightarrow T)$, which is a valid formula and thus does not further restrict our theory. (Restricting the objects that satisfy $P$ to all objects in the domain is no restriction.)

On the other hand, if there are *no* clauses in $\Delta$ that are positive in $P$, we can assume any valid one; e.g., $(\forall x\ F \Rightarrow P(x))$. Completion of $P$ then gives the completion formula $(\forall x\ P(x) \Rightarrow F)$, which is equivalent to $(\forall x\ \neg P(x))$. In this latter case, $\Delta$ says nothing about there being any objects satisfying $P$; therefore, we assume there are none.

Although in simple cases predicate completion and the CWA have the same effect, in general they are different. For example, suppose that $\Delta$ contains the single formula P(A), and that the language also contains the object constant B. Then the CWA extension includes ¬P(B), and the completion formula is $(\forall x\ P(x) \Rightarrow (x=A))$. These two expressions are not equivalent, although ¬P(B) plus the DCA entails $(\forall x\ P(x) \Rightarrow (x=A))$; and $(\forall x\ P(x) \Rightarrow (x=A))$ plus the UNA entails ¬P(B). ([Lifschitz 1985b] derives general conditions relating these two augmenting conventions.)

Predicate completion, like the CWA, is nonmonotonic because, if another clause positive in $P$ were added to $\Delta$, the completion formula for $P$ would be different. In general, it would be weaker; i.e., the augmented theory would allow more objects to satisfy $P$ than the original one did. Thus, some proofs of expressions of the form $\neg P$ could no longer be obtained. In our earlier example about birds, if we were to augment $\Delta$ by including `Penguin(x)`$\Rightarrow$`Bird(x)`, then the new completion formula for `Bird` would be:

$$\text{Bird(x)} \Rightarrow \text{Ostrich(x)} \vee \text{Penguin(x)} \vee \text{x=Tweety}$$

and we could no longer prove `¬Bird(Sam)` as we could earlier. (Sam might be a penguin.)

Extending a set of beliefs by predicate completion preserves their consistency.

**THEOREM 6.2** *If $\Delta$ is a consistent set of clauses solitary in $P$, then the completion of $P$ in $\Delta$ is consistent.*

This theorem follows from stronger results, Theorem 6.7 or Theorem 6.8, given later in the chapter (also without proof).

We also can perform predicate completion of several predicates in parallel. In *parallel predicate completion* of a set of predicates, each predicate in the set is completed separately (without regard for the others), and the conjunction of these separate completion formulas is added to $\Delta$. The completion process for each uses only the original clauses in $\Delta$ and not the formulas added by the completion process for other predicates. Parallel predicate completion allows us to restrict the objects that satisfy any of several predicates to those that are forced to do so by $\Delta$.

For the several completion formulas to avoid circularity, we must impose a condition on the way in which the predicates being completed can occur in $\Delta$. To motivate this condition, consider the following clauses (which are solitary in P, Q, and R):

$$Q(x) \Rightarrow P(x)$$
$$R(x) \Rightarrow Q(x)$$
$$P(x) \Rightarrow R(x)$$

Parallel predicate completion of $\{P,Q,R\}$ would yield:

$$P(x) \Leftrightarrow Q(x) \Leftrightarrow R(x) \Leftrightarrow P(x)$$

which is circular.

Recall that writing clauses that are solitary in $P$ in their normal forms allowed us to combine all clauses in $\Delta$ containing a positive $P$ literal into a single formula of the form

$$\forall x \ E_1 \ \lor \ E_2 \ \lor \ldots \lor \ E_k \ \Rightarrow \ P(x)$$

Denoting the antecedent of this implication simply by $E$ gives us

$$\forall x \ E \ \Rightarrow \ P(x)$$

where $E$ contains *no* occurrences of $P$.

To perform parallel predicate completion of the set $\Pi = \{P_1, P_2, \ldots, P_n\}$ of predicates in $\Delta$, we first write those clauses in $\Delta$ containing members of $\Pi$ in their normal forms, then combine the clauses containing the same $P_i$s into single formulas.

$$\forall x \ E_1 \ \Rightarrow \ P_1(x)$$

$$\forall x \ E_2 \ \Rightarrow \ P_2(x)$$

$$\forall x \ E_3 \ \Rightarrow \ P_3(x)$$

$$\vdots$$

$$\forall x \ E_n \ \Rightarrow \ P_n(x)$$

Parallel predicate completion is then accomplished by adding to $\Delta$ the completion formulas $(\forall x \ P_i(x) \Rightarrow E_i)$ for $i = 1, \ldots, n$. To avoid circular definitions of the $P_i$, we must be able to order the $P_i$ such that each $E_i$ has no occurrences of any of $\{P_i, P_{i+1}, \ldots, P_n\}$ (and has no negative occurrences of any of $\{P_1, \ldots, P_{i-1}\}$). If this ordering can be achieved, we say that the clauses in $\Delta$ are *ordered in* $\Pi$. In the next section, we illustrate parallel predicate completion with an example.

Note that, if $\Delta$ is ordered in $\Pi$, it also is solitary in each of the $P_i$ individually (but not necessarily conversely).

Theorem 6.2, about the consistency of predicate completion, can be generalized to the case of parallel predicate completion.

**THEOREM 6.3** *If $\Delta$ is consistent and ordered in $\Pi$, then the parallel completion of $\Pi$ in $\Delta$ is consistent.*

This theorem is a consequence of extended versions of either Theorem 6.7 or Theorem 6.8, given later in the chapter.

## 6.3  Taxonomic Hierarchies and Default Reasoning

Several AI systems have included simple mechanisms to allow a kind of
reasoning called *reasoning by default.* Since typically birds can fly, for
example, we assume (by default) that an arbitrary bird can fly unless we
know that it cannot. In this section, we describe a technique for stating
the *typical* properties of objects and then show how a variant of parallel
predicate completion can be used to perform default reasoning.

Often, this style of reasoning is applied to taxonomic hierarchies in which
subclasses *inherit* the properties of their superclasses unless these properties
are specifically canceled. Suppose, for example, that our beliefs, $\Delta$, include
the following formulas that define a taxonomic hierarchy:

$$\text{Thing(Tweety)}$$

$$\text{Bird(x)} \Rightarrow \text{Thing(x)}$$

$$\text{Ostrich(x)} \Rightarrow \text{Bird(x)}$$

$$\text{Flying-Ostrich(x)} \Rightarrow \text{Ostrich(x)}$$

(Tweety is a thing; all birds are things; all ostriches are birds; all flying
ostriches are ostriches.)

This subset of $\Delta$ that defines the taxonomic hierarchy will be denoted
by $\Delta_H$.

Suppose we also want to include in $\Delta$ statements that describe some of
the properties of the objects in the taxonomic hierarchy. For example, we
might want to say that no things except birds can fly, that all birds except
ostriches can fly, and that no ostriches except flying ostriches can fly. One
way to do this is with the following formulas:

a. $\text{Thing(x)} \land \neg\text{Bird(x)} \Rightarrow \neg\text{Flies(x)}$
b. $\text{Bird(x)} \land \neg\text{Ostrich(x)} \Rightarrow \text{Flies(x)}$
c. $\text{Ostrich(x)} \land \neg\text{Flying-Ostrich(x)} \Rightarrow \neg\text{Flies(x)}$
d. $\text{Flying-Ostrich(x)} \Rightarrow \text{Flies(x)}$

The subset of $\Delta$ that describes properties of objects in the hierarchy will
be denoted by $\Delta_P$. Whether we regard a predicate as defining a taxonomic
*kind* of object or a nontaxonomic *property* of an object is left to us. In
this example, we choose to think of flying simply as a property that certain
objects have—not as defining a kind of object.

Here, exceptions to general rules are listed explicitly in the rules. If
we had exceptions to birds flying other than ostriches, we would have to
list each of these exceptions in rule b. Of course, a general commonsense
reasoning system would need to know about other common exceptions,
such as penguins and baby birds. As we mentioned earlier in discussing the
qualification problem, there would be no particular difficulty in principle

with listing all known exceptions in the rule. The problem is that the system designer cannot really think of *all* of the exceptions that the system might later confront—exceptions such as wingless eagles, brain-damaged gulls, and roast ducks. Instead of listing all these exceptions, we would prefer some technique that allows us to say that birds (typically) can fly unless they are abnormal in some respect—an abnormality shared by ostriches, penguins, and such. Exceptions that we think of later can then be simply introduced by conferring this same abnormality on the new exceptions. Similarly, we would want to say that things (typically) cannot fly unless they are abnormal in some other respect—an abnormality shared by birds, airplanes, and mosquitoes. A hierarchy of exceptions would thus have to deal with several different kinds of abnormalities. We make these abnormalities part of the taxonomic hierarchy.

The following rule seems to capture what we want to say about things in general:

$$\text{Thing(x)} \land \lnot \text{Ab1(x)} \Rightarrow \lnot \text{Flies(x)}$$

where $\text{Ab1}$ is a predicate that has to do with that particular type of abnormality that must be provably absent if we are to use this general rule to prove that things cannot fly. Thus, our rule states that things do not fly unless they have an abnormality of type 1, say. (We will soon introduce other types of abnormalities.)

Birds are among those objects that have an abnormality of type 1:

$$\text{Bird(x)} \Rightarrow \text{Ab1(x)}$$

We call this rule an *inheritance cancellation rule*. The taxonomic rule $\text{Bird(x)} \Rightarrow \text{Thing(x)}$ ordinarily could be used to conclude that birds *inherit* the traits of things generally—including the inability to fly (if they are not abnormal). Cancellation rules, declaring abnormalities, thus block the inheritance of specified traits. We include them in $\Delta_H$, the formulas that define the taxonomic hierarchy.

The designer of a commonsense reasoning system could put in such information as is available about objects that might have abnormalities of type 1; e.g., airplanes, certain insects, and so on. The important feature of this way of dealing with exceptions is that additional axioms about abnormalities can be added at any time. New knowledge about flying objects can be expressed by *adding* axioms to the belief set instead of by having to *change* them!

Continuing with our example, we express the general knowledge that birds (typically) can fly by the rule:

$$\text{Bird(x)} \land \lnot \text{Ab2(x)} \Rightarrow \text{Flies(x)}$$

The predicate `Ab2` is about a type of abnormality the presence of which in birds prevents us from using this rule to conclude that those birds can fly. Ostriches are among those objects that have this type of abnormality; thus, we have another cancellation rule:

$$\texttt{Ostrich(x)} \Rightarrow \texttt{Ab2(x)}$$

Ordinarily, ostriches cannot fly:

$$\texttt{Ostrich(x)} \wedge \neg\texttt{Ab3(x)} \Rightarrow \neg\texttt{Flies(x)}$$

The predicate `Ab3` is about a type of abnormality the presence of which in ostriches prevents us from using this rule to conclude that those ostriches cannot fly. Flying ostriches (if there are such) are among those objects having this type of abnormality:

$$\texttt{Flying-Ostrich(x)} \Rightarrow \texttt{Ab3(x)}$$

Using this approach, we now have the following formulas in $\Delta_H$ defining the taxonomic hierarchy:

$$\texttt{Flying-Ostrich(x)} \Rightarrow \texttt{Ostrich(x)}$$
$$\texttt{Flying-Ostrich(x)} \Rightarrow \texttt{Ab3(x)}$$
$$\texttt{Ostrich(x)} \Rightarrow \texttt{Bird(x)}$$
$$\texttt{Ostrich(x)} \Rightarrow \texttt{Ab2(x)}$$
$$\texttt{Bird(x)} \Rightarrow \texttt{Thing(x)}$$
$$\texttt{Bird(x)} \Rightarrow \texttt{Ab1(x)}$$
$$\texttt{Thing(Tweety)}$$

(We include the information that Tweety is a "thing" to illustrate how our approach can be used to reason nonmonotonically about the properties of Tweety.)

This taxonomy is represented graphically by the network of Figure 6.1. Note that our taxonomy does not have to be a tree. (To allow us to use parallel predicate completion, which we will do, our taxonomy does have to be a partial order.)

The properties of objects in the hierarchy are described by the following formulas in $\Delta_P$:

$$\texttt{Thing(x)} \wedge \neg\texttt{Ab1(x)} \Rightarrow \neg\texttt{Flies(x)}$$
$$\texttt{Bird(x)} \wedge \neg\texttt{Ab2(x)} \Rightarrow \texttt{Flies(x)}$$

$$\text{Ostrich(x)} \wedge \neg\text{Ab3(x)} \Rightarrow \neg\text{Flies(x)}$$

$$\text{Flying-Ostrich(x)} \Rightarrow \text{Flies(x)}$$

We now perform parallel predicate completion on the set {Ab1, Ab2, Ab3, Flying-Ostrich, Ostrich, Bird, Thing} in just $\Delta_H$ to make the assumption that the only objects that are things, birds, ostriches, flying ostriches, or objects that are abnormal in any respect are those that are forced to be so by $\Delta_H$. The clauses in $\Delta_H$ have an ordering in the set {Ab1, Ab2, Ab3, Flying-Ostrich, Ostrich, Bird, Thing}, so parallel predicate completion will not result in circular definitions.

In this simple example, we obtain the following completion clauses (by completing {Ab1, Ab2, Ab3, Flying-Ostrich, Ostrich, Bird, Thing}, in $\Delta_H$):

1. Thing(x) ⇒ Bird(x) ∨ x=Tweety
2. Bird(x) ⇒ Ostrich(x)
3. Ostrich(x) ⇒ Flying-Ostrich(x)
4. ¬Flying-Ostrich(x)
5. Ab1(x) ⇒ Bird(x)
6. Ab2(x) ⇒ Ostrich(x)
7. Ab3(x) ⇒ Flying-Ostrich(x)

The only object mentioned is Tweety, and it is a thing, so these clauses appropriately tell us that there are no things other than Tweety, and no birds, ostriches, or flying ostriches at all. Also, there are no objects
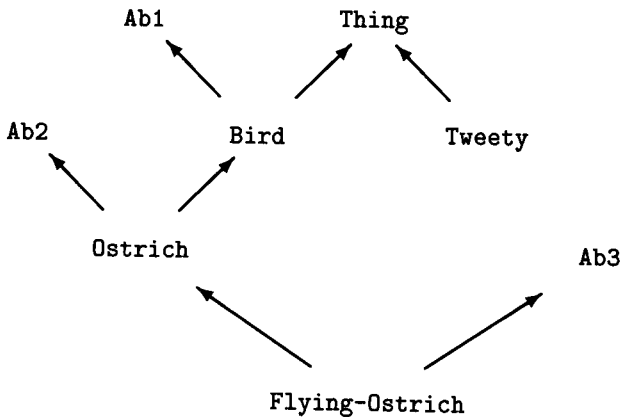


**Figure 6.1   A taxonomic hierarchy with abnormalities.**

abnormal in any respect. Using the properties mentioned in $\Delta_P$, we can deduce ¬Flies(Tweety) after first proving ¬Flying-Ostrich(Tweety), ¬Ostrich(Tweety), ¬Bird(Tweety), and ¬Ab1(Tweety).

If we were to add Bird(Tweety) to our taxonomic hierarchy, completion formula 2 would change to Bird(x)⇒Ostrich(x)∨(x=Tweety). We would still be able to prove ¬Ab2(Tweety) (but not ¬Ab1(Tweety)), so we would then be able to conclude Flies(Tweety), and so on. As the reasoning system learns about other objects and other ways in which objects might have various types of abnormalities, the taxonomy changes, the predicate completion formulas are recomputed as appropriate, and the conclusions that the system can derive change correspondingly.

We call this process of completing predicates in a subset of $\Delta$ *delimited completion*. It is important to note that delimited completion of a set of predicates is not generally the same as completing these same predicates in the full $\Delta$. (The reader should work out full completion in this case as an exercise.) Delimited completion typically produces a stronger augmenting assumption than would be produced by completion of the same predicates in the full $\Delta$, but it is an assumption that is often justified and useful. One must be careful, however, because delimited completion might produce an inconsistent augmentation (see Exercise 6 at the end of the chapter). Later, we shall discuss a more general and robust procedure for augmenting beliefs with default assumptions of this kind.

## 6.4 Circumscription

Reviewing what we have said about augmentation conventions so far, we have seen that the CWA augments a belief set to include the negations of ground positive literals when these ground atoms cannot be proved true. Predicate completion is defined for belief sets consisting of clauses solitary in a predicate and augments such belief sets by formulas that state that the only objects satisfying the predicate are those that must do so, given the belief set.

Both of these augmenting ideas are based on a kind of *minimization* principle. In the case of predicate completion, the minimality idea is particularly clear. If that part of $\Delta$ containing a predicate $P$ to be completed is written as $(\forall x)[E⇒P(x)]$, then $P$ is *completed* by the formula $(\forall x)[P(x)⇒E]$. That is, no objects have property $P$ unless $\Delta$ says they must.

We might like to use the same minimality assumption (i.e., that the only objects satisfying $P$ are those that must, given $\Delta$) in cases where $\Delta$ cannot be written as a set of clauses solitary in a predicate. For example, suppose $\Delta$ consists of the single formula $(\exists y\ P(y))$. What can be said about the smallest set of objects satisfying P in this case? This formula is not a clause, so we cannot use predicate completion. We know that there must be at least *one* object satisfying P, but there is nothing in $\Delta$ that

says that there need be more than one. We can assume there is only one by adding the formula $(\exists y \forall x \ (x=y) \iff P(x))$.

Now, suppose $\Delta$ consists of the single clause $(P(A) \lor P(B))$. This clause is not solitary in P, so we cannot use predicate completion here either. Intuitively, however, it would seem that the formula

$$(\forall x \ P(x) \iff x{=}A) \ \lor \ (\forall x \ P(x) \iff x{=}B)$$

says what we want to say about $P$ being minimized.

To derive augmentations of this sort for arbitrary belief sets, we must delve into the minimization process in more detail. In doing so, we define a process called *circumscription* that, like predicate completion, involves computing a certain formula that, when conjoined to $\Delta$, says that the only objects satisfying a predicate are those that have to do so, given $\Delta$.

Circumscription is based on the idea of *minimal* models. Let $M[\Delta]$ and $M^*[\Delta]$ be two models of $\Delta$. (You may want to refer to Chapter 2 to review the definition of a model in logic.) We say that $M^*[\Delta]$ is no larger than $M[\Delta]$ in predicate $P$, writing $M^*[\Delta] \preceq_P M[\Delta]$, if (1) $M$ and $M^*$ have the same domain of objects, (2) all other relation and function constants in $\Delta$ besides $P$ have the same interpretations in $M$ and $M^*$, but (3) the extension of (the relation corresponding to) $P$ in $M^*$ is a subset of the extension of $P$ in $M$. That is, if $M^* \preceq_P M$, then the set of objects satisfying $P$ in $M^*$ is a subset of those satisfying $P$ in $M$. We write $M^* \prec_P M$ if $M^* \preceq_P M$ and $M \npreceq_P M^*$.

There may be models of $\Delta$ that are *minimal in P* according to the ordering $\preceq_P$. $M_m$ is *P-minimal* if, for any $M \preceq_P M_m$, $M = M_m$. (As we shall see later, minimal models do not always exist.) If a model $M_m$ of $\Delta$ is $P$-minimal, then no objects satisfy the extension of $P$ except those that have to do so, given $\Delta$. We would like to find a sentence $\phi_P$ such that, for any $M$ that is a model of $\Delta \land \phi_P$, there is no $M^*$ that is a model of $\Delta$ with $M^* \prec_P M$. That is, the models of $\Delta \land \phi_P$ are $P$-minimal models of $\Delta$. The sentence $\phi_P$ when conjoined to $\Delta$ asserts that there are no objects that satisfy $P$ except those that have to do so, given $\Delta$. We call this conjunction the *circumscription of P in* $\Delta$.

To find an expression for $\phi_P$ in terms of $P$ and $\Delta$, we reason as follows. Let $P^*$ be a relation constant of the same arity as $P$, and let $\Delta(P^*)$ be $\Delta$ with each occurrence of the relation constant $P$ in $\Delta$ replaced by $P^*$. We note that any model of

$$(\forall x \ P^*(x) \Rightarrow P(x)) \ \land \ \neg(\forall x \ P(x) \Rightarrow P^*(x)) \ \land \ \Delta(P^*)$$

is *not* a $P$-minimal model of $\Delta$ because, in such a model, the extension corresponding to $P^*$ would be a strict subset of the extension corresponding

to $P$ (and $P^*$ satisfies $\Delta$). (Again, for brevity, we allow $\mathbf{x}$ to be possibly a tuple of variables.) Therefore, any model of

$$\neg((\forall\mathbf{x}\ P^*(\mathbf{x}) \Rightarrow P(\mathbf{x})) \wedge \neg(\forall\mathbf{x}\ P(\mathbf{x}) \Rightarrow P^*(\mathbf{x})) \wedge \Delta(P^*))$$

*is* a $P$-minimal model of $\Delta$.

Since $P^*$ in the preceding expression can be *any* relation constant of the same arity as $P$, the $\phi_P$ that we seek is the second-order logic formula obtained by universally quantifying over the relation variable $\mathbf{P*}$:

$$\forall\mathbf{P*}\ \neg((\forall\mathbf{x}\ \mathbf{P*}(\mathbf{x}) \Rightarrow P(\mathbf{x})) \wedge \neg(\forall\mathbf{x}\ P(\mathbf{x}) \Rightarrow \mathbf{P*}(\mathbf{x})) \wedge \Delta(\mathbf{P*}))$$

We call this expression the *circumscription formula* for $P$ in $\Delta$. Any model of the circumscription formula is a $P$-minimal model of $\Delta$. Conjoining the circumscription formula with $\Delta$ itself gives us the *circumscription* of $P$ in $\Delta$:

$$\mathrm{CIRC}[\Delta; P] \equiv_{def} \Delta \wedge \forall\mathbf{P*}\ \neg((\forall\mathbf{x}\ \mathbf{P*}(\mathbf{x}) \Rightarrow P(\mathbf{x})) \wedge$$
$$\neg(\forall\mathbf{x}\ P(\mathbf{x}) \Rightarrow \mathbf{P*}(\mathbf{x})) \wedge \Delta(\mathbf{P*}))$$

Although the use of a second-order formula is unsettling (since we have not described inference techniques for second-order logics), we shall see that, in many important cases, this formula can be reduced to an equivalent first-order formula.

Before discussing methods for simplifying the second-order formula for circumscription, we shall first rewrite it in some alternative forms.

Distributing the negation over the three conjuncts in the circumscription formula and then rewriting the resulting disjunction as an implication yields the usual form for circumscription:

$$\mathrm{CIRC}[\Delta; P] \equiv \Delta \wedge (\forall\mathbf{P*}\ (\Delta(\mathbf{P*}) \wedge (\forall\mathbf{x}\ \mathbf{P*}(\mathbf{x}) \Rightarrow P(\mathbf{x}))) \Rightarrow$$
$$(\forall\mathbf{x}\ P(\mathbf{x}) \Rightarrow \mathbf{P*}(\mathbf{x})))$$

We can get an additional perspective by deriving another form of circumscription. Since the preceding circumscription formula is universally quantified over $\mathbf{P*}$, it holds in particular for $P \wedge P'$ substituted for $\mathbf{P*}$ (where $P'$ is a relation constant of the same arity as $P$):

$$\Delta(P \wedge P') \wedge (\forall\mathbf{x}\ P(\mathbf{x}) \wedge P'(\mathbf{x}) \Rightarrow P(\mathbf{x})) \Rightarrow$$
$$(\forall\mathbf{x}\ P(\mathbf{x}) \Rightarrow P(\mathbf{x}) \wedge P'(\mathbf{x}))$$

This formula reduces to

$$\Delta(P \wedge P') \Rightarrow (\forall\mathbf{x}\ P(\mathbf{x}) \Rightarrow P'(\mathbf{x}))$$

Since $P'$ is arbitrary, this formula says that $P$ is circumscribed if and only if any alleged strengthening of $P$ (say to $P \wedge P'$) that still satisfies $\Delta$ is no real strengthening because $P$ would then already imply $P'$ anyway.

It is convenient to abbreviate $(\forall x \; P^*(x) \Rightarrow P(x))$ by the expression $P^* \leq P$. We also use the abbreviations $(P^* < P)$ for $((P^* \leq P) \wedge \neg(P \leq P^*))$ and $(P^* = P)$ for $((P^* \leq P) \wedge (P \leq P^*))$. These abbreviations also help us to remember that, when $(\forall x \; P^*(x) \Rightarrow P(x))$, an extension corresponding to $P^*$ is a subset of an extension corresponding to $P$.

In terms of these abbreviations, we can write the circumscription formula as

$$\forall P* \; (\Delta(P*) \wedge (P* \leq P)) \; \Rightarrow \; (P \leq P*)$$

which is equivalent to

$$\forall P* \; \Delta(P*) \; \Rightarrow \; \neg(P* < P)$$

or

$$\neg(\exists P* \; \Delta(P*) \wedge (P* < P))$$

This last form of the circumscription formula makes the intuitively helpful statement that there can be no P* that, when substituted for $P$ in $\Delta$, still satisfies $\Delta$ and that has a corresponding extension that is a strict subset of the extension corresponding to $P$.

There are several cases in which circumscription can be simplified. The following theorem often is useful:

**THEOREM 6.4** *Given a predicate $P$, an arbitrary belief set $\Delta(P)$ (containing the predicate $P$), and any $P'$ of the same arity as $P$ but not defined in terms of $P$, then, if $\Delta(P) \models \Delta(P') \wedge (P' \leq P)$, $\mathrm{CIRC}[\Delta; P] \equiv \Delta(P) \wedge (P = P')$.*

We first discuss the importance of this theorem and then give its proof and an example of its use. The theorem states that, if we have some predicate $P'$ of the same arity as $P$ but not containing $P$, and if we can prove $\Delta(P') \wedge (P' \leq P)$, given $\Delta$, then $(P = P')$ is equivalent to the circumscription formula for $P$ in $\Delta$. The theorem is most often used to confirm conjectures about circumscription formulas. $P'$ may contain bound predicate variables, so the circumscription formula may still be a second-order formula; in many cases of interest, however, it will be a first-order one.

***Proof*** Assume the conditions of the theorem; namely,

$$\Delta[P] \models \Delta[P'] \land (P' \leq P)$$

*Left to right*: Assume CIRC$[\Delta; P]$; that is, assume

$$\Delta(P) \land (\forall P* \ \Delta[P*] \land (P* \leq P) \ \Rightarrow \ (P \leq P*))$$

Using the condition of the theorem, we then have

$$\Delta(P') \land (P' \leq P)$$

Universal specialization on the circumscription formula yields

$$\Delta(P') \land (P' \leq P) \ \Rightarrow \ (P \leq P')$$

Modus ponens applied to these last two expressions yields

$$(P \leq P')$$

This result, together with $(P' \leq P)$, yields $(P = P')$.

*Right to left*: If the circumscription formula did not follow from the conditions of the theorem, we would have some $P*$ such that $\Delta(P*) \land (P* < P)$. Assuming $P = P'$ (the right side of the equivalence in the theorem), we then have $\Delta(P*) \land (P* < P')$. The conditions of the theorem, however, state that $\Delta(P*)$ logically entails $(P' \leq P*)$—a contradiction. $\square$

As an example of how this theorem is used, let $\Delta$ be P(A) $\land$ ($\forall$x Q(x)$\Rightarrow$P(x)). If we were going to perform predicate completion on these clauses, we would rewrite $\Delta$ as ($\forall$x Q(x) $\lor$ (x=A)$\Rightarrow$P(x)). Predicate completion would yield the completion formula ($\forall$x P(x)$\Rightarrow$Q(x) $\lor$ (x=A)). Since predicate completion was motivated as a technique for minimizing the objects that satisfy a predicate, we might suspect that it gives the same result as does circumscription. We can use Theorem 6.4 to show that it does in this case.

We let the $P'$ of the theorem be the consequent of the completion formula Q(x) $\lor$ (x=A). Strictly speaking, we should write $P'$ as ($\lambda$x Q(x) $\lor$ (x=A)), a *lambda* expression. Now, to use the theorem, we must prove that $\Delta$ logically entails $\Delta(P') \land (P' \leq P)$.

Substituting ($\lambda$x Q(x) $\lor$ (x=A)) for $P$ in $\Delta$ yields

$$\Delta(P') \equiv (\forall x \ Q(x) \ \Rightarrow \ Q(x) \ \lor \ x=A) \land (Q(A) \ \lor \ A=A)$$

We see that $\Delta(P')$ is trivially valid. It remains to show that $\Delta$ logically implies $(P' \leq P)$; that is, $(\forall x \; Q(x) \; \vee \; (x=A) \Rightarrow P(x))$. The latter formula, however, is just the normal form of $\Delta$. Thus, the conditions of the theorem are satisfied, and the theorem confirms that $\text{CIRC}[\Delta; P]$ is $(\forall x \; Q(x) \; \vee \; (x=A) \; \Leftrightarrow \; P(x))$.

This example can be generalized to show that predicate completion gives the same results as circumscription when $\Delta$ consists of clauses solitary in $P$.

In many cases of interest in AI, $\text{CIRC}[\Delta; P]$ "collapses" to a first-order formula. (We give some examples later in which it does not.) The simplest case in which circumscription is collapsible is when all occurrences of $P$ in $\Delta$ are *positive*. (A formula has a *positive occurrence of $P$* if $P$ occurs positively in its clause form; it has a *negative occurrence of $P$* if $P$ occurs negatively in its clause form.)

As an example, consider the case in which $\Delta$ is $(\exists y \; P(y))$. $P$ occurs only positively in $\Delta$. By manipulating expressions in second-order logic, it can be shown that the circumscription of $P$ in $\Delta$ is $(\exists y \forall x \; (x=y) \; \Leftrightarrow \; P(x))$. Circumscribing $P$ in this case limits the extension of $P$ to a minimal nonempty set; i.e., a singleton.

An important case in which circumscription is collapsible can best be thought of as a simple generalization of the solitary condition used in defining predicate completion. Earlier, we defined what is meant for clauses to be *solitary* in a predicate $P$. Recall that a clause is solitary in $P$ if, whenever it contains a positive occurrence of $P$, it contains only one occurrence of $P$. Generalizing this definition, we say that a *formula* is *solitary in $P$* if and only if it can be written in the following *normal form*:

$$N[P] \wedge (E \leq P)$$

where $N[P]$ is a formula that contains no positive occurrences of $P$, $E$ is a formula that contains no occurrences of $P$, and $E \leq P$ is our usual abbreviation for $(\forall x \; E(x) \Rightarrow P(x))$ (where x might be a tuple of variables).

Note that the *normal form* of a conjunction of clauses solitary in $P$ is in the form $E \leq P$. Thus, solitary clauses are a special case of solitary formulas.

For solitary formulas in general, we have the following theorem.

**THEOREM 6.5** $\text{CIRC}[N[P] \wedge (E \leq P); P] \equiv N[E] \wedge (E = P)$, *where* $N[E]$ *is* $N[P]$ *with each occurrence of $P$ replaced by $E$.*

**Proof** This theorem follows directly from Theorem 6.4. First note that, since $N[P]$ has no positive occurrences of $P$, $N[P] \wedge (E \leq P)$ logically entails $N[E]$. (This entailment can be thought of as a kind of "generalized resolution.") Thus, the conditions of Theorem 6.4 are satisfied. $\square$

Thus, for solitary formulas, circumscription is collapsible to a first-order formula. We see that circumscription gives the same result as predicate completion does in the special case of clauses solitary in $P$. Using Theorem 6.5, we can compute circumscription for theories not in clause form, as long as they can be written in normal form.

We illustrate with an example. Let $\Delta$ be given by

$$\exists x \ \neg On(A,x) \ \wedge \ On(A,B)$$

We want to compute the circumscription of On in $\Delta$. We can write $\Delta$ in normal form to show that it is solitary in On:

$$(\exists x \ \neg On(A,x)) \ \wedge \ (\forall x \forall y \ x{=}A \ \wedge \ y{=}B \ \Rightarrow \ On(x,y))$$

We identify the first conjunct of this expression as $N[On]$ (it has no positive occurrences of On), and the second as $(E \leq On)$ with $E(x,y) \equiv (x{=}A) \wedge (y{=}B)$. ($E$ has no occurrences of On.) Thus, by the theorem, $CIRC[\Delta; On]$ is

$$(\forall x \forall y \ On(x,y) \ \Leftrightarrow \ x{=}A \ \wedge \ y{=}B) \ \wedge \ (\exists x \ \neg(x{=}B))$$

(The only thing "on" something is the object denoted by A, and it is on the object denoted by B, and there is at least one object that is not the same as the one denoted by B.)

Some interesting complications arise when attempting circumscription in formulas slightly more general than solitary ones are. Consider an example in which $\Delta$ is

$$Ostrich(x) \ \Rightarrow \ Bird(x)$$

$$Bird(Tweety) \ \vee \ Bird(Sam)$$

We cannot use Theorem 6.5 to compute the circumscription of Bird in $\Delta$, because $\Delta$ is not solitary in Bird. Before we attempt circumscription, however, let us see whether we can guess what sort of augmenting statement circumscription might make about Bird. From $\Delta$, we can guess that there might be two alternative minimizations of Bird; namely,

- $\forall x \ Bird(x) \ \Leftrightarrow \ Ostrich(x) \ \vee \ x{=}Tweety$
- $\forall x \ Bird(x) \ \Leftrightarrow \ Ostrich(x) \ \vee \ x{=}Sam$

The belief set is not sufficiently "definite" to allow us to determine which one of these might be appropriate. This indefiniteness makes it impossible to give a single minimizing Bird. Instead, we can merely say something *about* the minimization of Bird; namely, that it has to be one or the other

of these two expressions. That is, all we can apparently say about the minimizing Bird is the following:

$$(\forall x \; Bird(x) \Leftrightarrow Ostrich(x) \lor x{=}Tweety) \lor$$
$$(\forall x \; Bird(x) \Leftrightarrow Ostrich(x) \lor x{=}Sam)$$

Indeed, we can use circumscription to derive this formula. The general circumscription formula for Bird in $\Delta$ is

$$\forall Bird* \; \Delta(Bird*) \land (\forall x \; Bird*(x) \Rightarrow Bird(x)) \Rightarrow$$
$$(\forall x \; Bird(x) \Rightarrow Bird*(x))$$

Let us first substitute Ostrich(x)$\lor$(x=Tweety) for Bird*(x). This yields (after reduction)

$$(\forall x \; Ostrich(x) \lor x{=}Tweety \Rightarrow Bird(x)) \Rightarrow$$
$$(\forall x \; Bird(x) \Rightarrow Ostrich(x) \lor x{=}Tweety)$$

Next, we substitute Ostrich(x)$\lor$(x=Sam) for Bird*(x). This yields

$$(\forall x \; Ostrich(x) \lor x{=}Sam \Rightarrow Bird(x)) \Rightarrow$$
$$(\forall x \; Bird(x) \Rightarrow Ostrich(x) \lor x{=}Sam)$$

Neither of these formulas has an antecedent that follows from $\Delta$, but the disjunction of the antecedents does. That is, from $\Delta$, we can prove

$$(\forall x \; Ostrich(x) \lor x{=}Sam \Rightarrow Bird(x)) \lor$$
$$(\forall x \; Ostrich(x) \lor x{=}Tweety \Rightarrow Bird(x))$$

(In fact, $\Delta$ itself can be written in this form. We first rewrite Bird(Tweety) and Bird(Sam) as $(\forall x \; (x{=}Tweety) \Rightarrow Bird(x))$ and $(\forall x \; (x{=}Sam) \Rightarrow Bird(x))$, respectively. Then, by distributivity, we write the conjunction of these formulas and $(\forall x \; Ostrich(x) \Rightarrow Bird(x))$ in the form shown previously.)

Since the disjunction of the antecedents of instances of the circumscription formulas follows from $\Delta$, the disjunction of the consequents does also. The disjunction of the consequents, however, is exactly the formula we guessed would be an appropriate statement to make about the minimizing Bird in this case.

The interesting point about this example is that we can derive a slightly more restrictive statement about Bird from the circumscription formula. In this example, $\Delta$ does not force us to accept a formula about Bird as general as we have guessed; the perceptive reader may have already noted that the disjunction of definitions can be tighter. The formula we guessed, although true in all Bird-minimal models, does allow some non-Bird-minimal models

as well—in particular when *both* Tweety and Sam are birds. We will return to this example after describing how circumscription collapses for a class of formulas more general than solitary ones are.

Next we consider a more general class of formulas—those we call *separable*. A *formula* is *separable* with respect to a predicate $P$ if any one of the following conditions is satisfied:

(1) It has no *positive* occurrences of $P$.

(2) It is of the form $(\forall x \; E(x) \; \Rightarrow \; P(x))$, where x is a tuple of variables and $E(x)$ is a formula that does not contain $P$ (we again abbreviate to $E \leq P$).

(3) It is composed of conjunctions and disjunctions of separable formulas.

Note that this definition implies that formulas that are solitary with respect to $P$ also are separable with respect to $P$. Also, as we shall show, all quantifier-free formulas are separable.

Positive occurrences of $P$ in belief sets of this form are separated into isolated components, and this separation makes possible a collapsed version of circumscription—as we shall see.

First, we point out that a rather wide class of formulas can be put in separable form. In the following pairs of equivalent formulas, the one preceded by a dot is a version in which separability (by the above definition) is obvious. (In the first two cases, the formulas also are solitary in $P$.)

(1) P(A)

     • $\forall$x x=A $\Rightarrow$ P(x)

(2) $\forall$y P(F(y))

     • $\forall$x$\exists$y x=F(y) $\Rightarrow$ P(x)

(3) Bird(Tweety) $\lor$ Bird(Sam)

     • ($\forall$x x=Tweety $\Rightarrow$ Bird(x)) $\lor$ ($\forall$x x=Sam $\Rightarrow$ Bird(x))

(4) <any unquantified formula>

     • <move negations in and use the method of example 1. to rewrite any positive occurrences of $P$>

(5) ($\forall$u P(u,A)) $\lor$ ($\forall$u P(u,B))

     • ($\forall$u$\forall$x x=A $\Rightarrow$ P(u,x)) $\lor$ ($\forall$u$\forall$x x=B $\Rightarrow$ P(u,x))

However, ($\forall$u P(u,A) $\lor$ P(u,B)) is *not* separable with respect to P because it cannot be written as a propositional combination of separable formulas.

Although our definition of separability can easily be used to check whether or not a formula is separable (taking into account equivalences such

as those in the preceding list), it is not obvious what this definition might have to do with circumscription. It happens that there is a *normal form* for separable formulas—somewhat like the normal form used in defining solitary formulas. We describe this normal form next and then show how it is used to compute circumscription.

From the definition of separability, it is straightforward to show that every formula separable with respect to $P$ is equivalent to a formula in the following normal form with respect to $P$:

$$\bigvee_i [N_i[P] \wedge (E_i \leq P)]$$

where each $E_i$ is a formula that has no occurrences of $P$, and each $N_i[P]$ is a formula that has no positive occurrences of $P$.

We obtain this standard form from any conjunction or disjunction of (separable) formulas by distributivity and by applications of the following rules:

$$(\phi \Rightarrow P) \wedge (\psi \Rightarrow P) \equiv (\phi \vee \psi) \Rightarrow P$$
$$(\phi \Rightarrow P) \vee (\psi \Rightarrow P) \equiv (\phi \wedge \psi) \Rightarrow P$$
$$(\phi \Rightarrow P) \equiv \mathrm{T} \wedge (\phi \Rightarrow P)$$
$$\phi \equiv \phi \wedge (\mathrm{F} \Rightarrow P)$$

(The last two rules are sometimes needed to ensure that each disjunct in the normal form has $N_i$ and $E_i$ terms. Use of these rules gives a T for $N_i$ and an F for $E_i$. In this case, when writing $(E_i \leq P)$ in its expanded form, we write $(\forall \mathbf{x} \ \mathbf{F} \ \Rightarrow \ P(\mathbf{x}))$.)

If $\Delta$ is in normal form for $P$, then the circumscription of $P$ in $\Delta$ is collapsible to a first-order formula, as defined in the following theorem.

**THEOREM 6.6** *Suppose $\Delta$ is separable with respect to $P$ and has normal form with respect to $P$ given by*

$$\bigvee_i [N_i[P] \wedge (E_i \leq P)]$$

*Then the circumscription of $P$ in $\Delta$ is equivalent to*

$$\bigvee_i [D_i \wedge (P = E_i)]$$

*where $D_i$ is*

$$N_i[E_i] \wedge \bigwedge_{j \neq i} \neg[N_j[E_j] \wedge (E_j < E_i)]$$

*and each $N[E]$ is $N[P]$ with all occurrences of $P$ replaced by $E$.*

(Recall that the expanded version of $(E_j < E_i)$ is $[(E_j \leq E_i) \wedge \neg(E_i \leq E_j)]$ which, further expanded, is $(\forall x\ E_j(x) \Rightarrow E_i(x)) \wedge \neg(\forall x\ E_i(x) \Rightarrow E_j(x))$.)

To demonstrate that circumscription implies a formula of the form $\bigvee_i[N_i[E_i] \wedge (P = E_i)]$ requires only a simple generalization of the proof of Theorem 6.5. To show that the extra conjuncts can be included in $D_i$ is somewhat more difficult. These extra conjuncts essentially allow us to omit from the disjunction of definitions for $P$ those disjuncts that would be redundant under certain conditions and given the other disjuncts. (The theorem is proved in [Lifschitz 1987b].)

We will illustrate the role of the $D_i$ by some examples later.

In computing circumscription, the result of Theorem 6.6 simplifies in some special cases. If the normal form has only one disjunct, then we have the special case of a formula solitary in $P$, and $D$ is $N[E]$. Or, if all the $N_i$ are T, then $D_i$ becomes

$$\bigwedge_{j \neq i}(E_i \leq E_j) \vee \neg(E_j \leq E_i)$$

Suppose, for example, that $\Delta$ is P(A) ∨ P(B). We rewrite this in normal form for P:

$$(T \wedge (\forall x\ x{=}A \Rightarrow P(x))) \vee (T \wedge (\forall x\ x{=}B \Rightarrow P(x)))$$

Here, the normal form has two disjuncts. $D_1$ and $D_2$ are, respectively,

$$(\forall x\ x{=}A \Rightarrow x{=}B) \vee (\exists y\ y{=}B \wedge \neg(y{=}A))$$

and

$$(\forall x\ x{=}B \Rightarrow x{=}A) \vee (\exists y\ y{=}A \wedge \neg(y{=}B))$$

which are both true, so the circumscription formula is equivalent to

$$(\forall x\ P(x) \Leftrightarrow x{=}A) \vee (\forall x\ P(x) \Leftrightarrow x{=}B)$$

(In computing the $D_i$ for this case, it is helpful to use the equivalence $(\forall x\ (x{=}A) \Rightarrow P(x)) \Leftrightarrow P(A)$.)

In the last example, the $D_i$ "disappeared," and we were left with a simple disjunction of definitions for $P$. To illustrate how the $D_i$ can act to restrict these disjunctions, consider the following example. Let $\Delta$ be given by P(A) ∨ (P(B) ∧ P(C)). In normal form, $\Delta$ is

$$(T \wedge (\forall x\ x{=}A \Rightarrow P(x))) \vee (T \wedge (\forall x\ x{=}B \vee x{=}C \Rightarrow P(x)))$$

Thus,

$$N_1 \equiv N_2 \equiv \text{T}$$
$$E_1 \equiv (\lambda x \ \text{x=A})$$
$$E_2 \equiv (\lambda x \ \text{x=B} \ \lor \ \text{x=C})$$
$$D_1 \equiv \text{T}$$
$$D_2 \equiv \text{A=B=C} \ \lor \ (\text{A≠B} \ \land \ \text{A≠C})$$

and Theorem 6.6 gives

$$\text{CIRC}[\Delta; P] \equiv (\forall x \ \text{P(x)} \ \Leftrightarrow \ \text{x=A}) \ \lor$$
$$((\forall x \ \text{P(x)} \ \Leftrightarrow \ \text{x=B} \ \lor \ \text{x=C}) \ \land$$
$$(\text{A=B=C} \ \lor \ (\text{A≠B} \ \land \ \text{A≠C})))$$

When (A=B=C), the first disjunct suffices, so this formula simplifies to

$$\text{CIRC}[\Delta; P] \equiv (\forall x \ \text{P(x)} \ \Leftrightarrow \ \text{x=A}) \ \lor$$
$$((\forall x \ \text{P(x)} \ \Leftrightarrow \ \text{x=B} \ \lor \ \text{x=C}) \ \land \ (\text{A≠B} \ \land \ \text{A≠C}))$$

This example nicely illustrates the role of the $D_i$. In this case, they tighten the definitions by taking into consideration the possibility that A may be equal to one of B,C. (If A were equal to one of B,C, then $\Delta \equiv$ P(A), and circumscription would yield simply ($\forall x$ P(x) $\Leftrightarrow$ (x=A)).)

Let us now reconsider the example that we discussed earlier when we attempted to guess the result of circumscription. In that example, $\Delta$ was given by

$$(\forall x \ \text{Ostrich(x)} \ \Rightarrow \ \text{Bird(x)}) \ \land \ (\text{Bird(Tweety)} \ \lor \ \text{Bird(Sam)})$$

The normal form is

$$(\text{T} \ \land \ (\forall x \ \text{Ostrich(x)} \ \lor \ \text{x=Tweety} \ \Rightarrow \ \text{Bird(x)})) \ \lor$$
$$(\text{T} \ \land \ (\forall x \ \text{Ostrich(x)} \ \lor \ \text{x=Sam} \ \Rightarrow \ \text{Bird(x)}))$$

Here again, the $D_i$ do not disappear. After some manipulation, we derive

$$D_1 \equiv \text{Sam=Tweety} \ \lor \ \lnot\text{Ostrich(Sam)} \ \lor \ \text{Ostrich(Tweety)}$$

Which, after the making the UNA, becomes

$$\lnot\text{Ostrich(Sam)} \ \lor \ \text{Ostrich(Tweety)}$$

and

$$D_2 \equiv \text{Tweety=Sam} \ \lor \ \lnot\text{Ostrich(Tweety)} \ \lor \ \text{Ostrich(Sam)}$$

which, again after the UNA is

$$\neg\text{Ostrich(Tweety)} \lor \text{Ostrich(Sam)}$$

Using these values in Theorem 6.6 yields

$$
\begin{aligned}
\text{CIRC}[\Delta;\text{Bird}] \equiv\ & ((\forall x\ \text{Bird}(x) \Leftrightarrow \text{Ostrich}(x) \lor x{=}\text{Tweety}) \land \\
& (\neg\text{Ostrich(Sam)} \lor \text{Ostrich(Tweety)})) \lor \\
& ((\forall x\ \text{Bird}(x) \Leftrightarrow \text{Ostrich}(x) \lor x{=}\text{Sam}) \land \\
& (\neg\text{Ostrich(Tweety)} \lor \text{Ostrich(Sam)}))
\end{aligned}
$$

The circumscription is more restrictive than is the formula we had guessed earlier. Circumscription says that there are two alternative "minimal definitions" of Bird; namely, either something is a bird if it is an ostrich or Tweety (and that definition need be a possibility only if Sam is not an ostrich or if Tweety is an ostrich), or something is a bird if it is an ostrich or Sam (and that definition need be a possibility only if Tweety is not an ostrich or if Sam is an ostrich). In our guess earlier, we did not narrow our definition as much as we might have; e.g., in the case in which Sam was an ostrich and Tweety was not. In that case, a minimal definition of Bird does not really need to include the possibility of conferring "birdhood" on Tweety (in order to satisfy $\Delta$) because Bird(Tweety) $\lor$ Bird(Sam) is satisfied by virtue of Sam being an ostrich.

In all the cases considered so far, we were able to construct a first-order formula the addition of which to $\Delta$ achieved the effect of circumscribing a predicate in $\Delta$. There are cases, however, in which circumcription does not collapse to a first-order formula. Here is an example: Suppose $\Delta$ contains the single formula

$$(\forall u \forall v\ Q(u,v) \Rightarrow P(u,v)) \land (\forall u \forall v \forall w\ P(u,v) \land P(v,w) \Rightarrow P(u,w))$$

In this case, the difficulty in saying that $\Delta$ expresses all and the only information about P is that $\Delta$ says quite a bit about P. It says that P is (at least) the transitive closure of Q. To circumscribe P in $\Delta$ would require saying that P is identically the transitive closure of Q, and this cannot be expressed in a first-order formula. One way of saying this, of course, is by the circumscription formula itself:

$$
\begin{aligned}
(\forall P*)\ & (\forall u \forall v\ Q(u,v) \Rightarrow P*(u,v)) \\
& \land (\forall u \forall v \forall w\ P*(u,v) \land P*(v,w) \Rightarrow P*(u,w)) \\
& \land (\forall u \forall v\ P*(u,v) \Rightarrow P(u,v)) \\
& \Rightarrow (\forall u \forall v\ P(u,v) \Rightarrow P*(u,v))
\end{aligned}
$$

Besides involving a second-order quantifier, this formula is not in the form of a definition for P. We can use Theorem 6.4 to rewrite the circumscription formula in its equivalent definitional form. We leave it to the reader to verify that the following expression for P' satisfies the conditions of Theorem 6.4:

```
P'(x,y) ⟺ (∀P* (∀u∀v Q(u,v) ⟹ P*(u,v))
   ∧(∀u∀v∀w (P*(u,v) ∧ P*(v,w) ⟹ P*(u,w)) ⟹ P*(x,y)))
```

Theorem 6.4 then states that circumscription is equivalent to the following definition for P:

$$\forall u \forall v \ P(u,v) \iff P'(u,v)$$

Another example of the insufficiency of a first-order formula to express circumscription comes from the algebraic axioms for natural numbers. Suppose $\Delta$ is

$$NN(0) \ \wedge \ (\forall x \ NN(x) \ \Rightarrow \ NN(S(x)))$$

That is, 0 is a nonnegative integer and the successor of every nonnegative integer is a nonnegative integer. Circumscribing NN in $\Delta$ produces an expression equivalent to the usual second-order formula for induction:

```
∀NN* (NN*(0) ∧ (∀x NN*(x) ⟹ NN*(S(x))))
   ∧ (∀x NN*(x) ⟹ NN(x))
   ⟹ (∀x NN(x) ⟹ NN*(x))
```

Replacing NN*(x) in this expression by [NN'(x) ∧ NN(x)] allows us to write

```
∀NN' NN'(0) ∧ (∀x NN'(x) ⟹ NN'(S(x)))
   ⟹ (∀x NN(x) ⟹ NN'(x))
```

which is closer to the usual induction formula.

The preceding two examples had belief sets that were neither positive nor separable in the predicates being circumscribed, and therefore it is not surprising that circumscription was not collapsible to first-order formulas in those cases.

It is also possible for $\Delta$ to have no minimal models. Consider the following set of formulas:

$$\exists x \ NN(x) \ \wedge \ (\forall y \ NN(y) \ \Rightarrow \ \neg(x=S(y)))$$

$$\forall x \ NN(x) \ \Rightarrow \ NN(S(x))$$

$$\forall x \forall y \ S(x)=S(y) \ \Rightarrow \ x=y$$

One interpretation of these formulas involves the claims that there is a number that is not the successor of any number, that every number has a successor that is a number, and that two numbers are equal if their successors are equal. One interpretation for NN is that any integer greater than $k$ satisfies NN. A "smaller" interpretation is that any integer greater than $k + 1$ satisfies NN—and so on. Thus, there is no NN-minimal model for $\Delta$. Since there is no NN-minimal model, we might expect that the circumscription of these formulas in NN is inconsistent, and indeed it is. (If the circumscription formula had a model, that model would be a minimal model of the formulas.)

Various sufficient conditions have been established for the circumscription of a consistent belief set to be consistent. We state some of the results here without proof.

**THEOREM 6.7** *If a belief set $\Delta$ is consistent and universal, then the circumscription of $P$ in $\Delta$ is consistent. (A set of formulas is said to be universal if either it is a set of clauses or if the conjunctive normal form of each of its formulas contains no Skolem functions.)*

**THEOREM 6.8** *If a belief set $\Delta$ is consistent and separable with respect to $P$, then the circumscription of $P$ in $\Delta$ is consistent.*

Since sets of clauses are universal and since circumscription of $P$ reduces to predicate completion of $P$ in clauses solitary (and thus separable) in $P$, Theorem 6.2 follows from either Theorem 6.7 or Theorem 6.8. (Theorem 6.3 follows from versions of these theorems extended to a more general case of circumscription, discussed in Section 6.7.)

Theorems 6.7 and 6.8 apply to two different kinds of formulas; namely, universal and separable ones. These two classes are instances of a common general class—the *almost universal formulas*. A formula is *almost universal with respect to $P$* if it has the form $(\forall x)\varphi$, where $x$ is a tuple of object variables, and $\varphi$ does not contain positive occurrences of $P$ in the scope of quantifiers. Every universal formula is, of course, almost universal with

respect to any $P$. It is not difficult to show that any formula that is separable with respect to $P$ also is almost universal with respect to $P$.

Theorems 6.7 and 6.8 are thus each a special case of Theorem 6.9.

**THEOREM 6.9** *If a belief set $\Delta$ is consistent and almost universal with respect to $P$, then the circumscription of $P$ in $\Delta$ is consistent.*

## 6.5 More General Forms of Circumscription

There are more general versions of circumscription that produce stronger results. First, instead of minimizing only a single predicate, we can minimize a set of predicates. The *parallel circumscription* of $\{P_1, P_2, \ldots, P_N\}$ in $\Delta$ is given by the same formula as before, except that we now let $P$ stand for a tuple of predicates:

$$\text{CIRC}[\Delta; P] \equiv \Delta(P) \land \neg(\exists\text{P*} \; \Delta(\text{P*}) \land (\text{P*} < P))$$

where P* is a tuple of predicate variables of the same arities as $P$, (P* $< P$) is an abbreviation for (P* $\leq P$) $\land \neg(P \leq$ P*), and (P* $\leq P$) is (P*$_1 \leq P_1$) $\land \ldots \land$ (P*$_N \leq P_N$).

Rewriting yields:

$$\text{CIRC}[\Delta; P] \equiv \Delta(P) \land (\forall\text{P*} \; (\Delta(\text{P*}) \land (\text{P*} \leq P)) \; \Rightarrow \; (P \leq \text{P*}))$$

Parallel circumscription is not really any more difficult to compute than is ordinary, single-predicate circumscription. Theorem 6.4, for example, is easily generalized to deal with this case. In fact, when all the predicates in the tuple $P$ occur positively in $\Delta$, we have Theorem 6.10.

**THEOREM 6.10** *If all occurrences of $P_1, P_2, \ldots, P_N$ in $\Delta$ are positive, then* $\text{CIRC}[\Delta; P]$ *is equivalent to*

$$\bigwedge_{i=1}^{N} \text{CIRC}[\Delta; P_i]$$

(This theorem is stated without proof in [Lifschitz 1986c], and is proved in [Lifschitz 1987b].)

As an example, we apply Theorem 6.10 to compute the parallel circumscription of {P1,P2} in ($\forall$x P1(x) $\lor$ P2(x)). Each of P1 and P2 occur positively in $\Delta$, so the parallel circumscription is just the conjunction of the individual circumscriptions of P1 and P2. Since each of $\text{CIRC}[\Delta; \text{P1}]$ and $\text{CIRC}[\Delta; \text{P2}]$ is ($\forall$x P1(x) $\Leftrightarrow$ $\neg$P2(x)), their conjunction is also.

The definitions of formulas solitary or separable in $P$ extend naturally to the case in which $P$ is a tuple of predicates. For example, a formula $\Delta$ is *solitary in a tuple of predicates* $P$ if it can be written in the form $N[P] \wedge (E \leq P)$, where $N[P]$ has no positive occurrences of any member of $P$, and no member of $E$ has any occurrences of any member of $P$. Theorems 6.5 and 6.6 (interpreting $P$ as a tuple of predicates) can then also be used to compute parallel circumscription.

For parallel circumscription, we can state a stronger result than would be obtained by extending Theorem 6.5 to formulas *solitary* in a tuple of predicates. Generalizing on the definition given in Section 6.2 of *clauses* ordered in $P$, we say that a *formula is ordered* in $P = \{P_1, P_2, \ldots, P_N\}$ if it can be written in the form:

$$N[P] \wedge (E_1 \leq P_1) \wedge (E_2 \leq P_2) \wedge \ldots \wedge (E_N \leq P_N)$$

where $N[P]$ has no positive occurrences of any of the predicates in $P$, and each $E_i$ has no occurrences of any of $\{P_i, P_{i+1}, \ldots, P_N\}$ and has no negative occurrences of any of $\{P_1, \ldots, P_{i-1}\}$.

Using this definition, we have the following theorem.

**THEOREM 6.11** *Suppose $\Delta$ is ordered in $P$ and can be written in the form $N[P] \wedge (E_1 \leq P_1) \wedge (E_2[P_1] \leq P_2) \wedge \ldots \wedge (E_N[P_1, P_2, \ldots, P_{N-1}] \leq P_N)$ (with no positive occurrences of the $P_i$ in $N$ and no occurrences of $P_i, \ldots, P_N$ in $E_i$).*
*Then the parallel circumscription of $P$ in $\Delta$ is given by*

$$\text{CIRC}[\Delta; P] \equiv N[E_1, \ldots, E_N] \wedge (P_1 = E_1) \wedge (P_2 = E_2[E_1]) \wedge \ldots$$
$$\wedge (P_N = E_N[E_1, \ldots, E_{N-1}])$$

The proof is analogous to that of Theorem 6.5 and, like it, is based on Theorem 6.4.

Note that parallel predicate completion for *clauses* ordered in $P$ is a special case of parallel circumscription.

Another generalization of circumscription allows other predicates to "vary" in addition to those being minimized. That is, we suppose that the extensions of these variable predicates can be changed during the minimization process to allow the predicates being circumscribed to have extensions even smaller than they would have otherwise. This means that an object will be allowed to satisfy one of the variable predicates (to satisfy $\Delta$), so that it does not have to satisfy a predicate being minimized (to satisfy $\Delta$). Which predicates should be allowed to vary depends on the purpose of the circumscription process; the decision is part of what we call the *circumscription policy*. The usual situation is that we are interested in knowing what is the effect of the circumscription of one predicate, $P$, (or

set of predicates) on some other variable predicate, $Z$, (or set of predicates). We want circumscription to minimize the number of objects satisfying $P$, even at the expense of allowing more or different objects to satisfy $Z$, the variable predicate. After defining circumscription with variable predicates, we will give examples of how this process is used.

Suppose $P$ is a tuple of predicates to be minimized and $Z$ (disjoint from $P$) is a tuple of predicates. The parallel circumscription of $P$ in $\Delta(P; Z)$, with predicates $Z$ also allowed to vary, is

$$\text{CIRC}[\Delta; P; Z] \equiv \Delta(P; Z) \wedge \neg(\exists \text{P}*\exists \text{Z}* \ \Delta(\text{P}*; \text{Z}*) \wedge (\text{P}* < P))$$

where P* and Z* are tuples of predicate variables (of the same arities as $P$ and $Z$, respectively), and $\Delta(\text{P}*; \text{Z}*)$ is the belief set expressed as a single wff with all occurrences of $P$ replaced by P* and all occurrences of $Z$ replaced by Z*.

Rewriting yields

$$\text{CIRC}[\Delta; P; Z]$$
$$\equiv \Delta(P; Z) \wedge (\forall \text{P}*\forall \text{Z}* \ (\Delta(\text{P}*; \text{Z}*) \wedge (\text{P}* \leq P)) \ \Rightarrow \ (P \leq \text{P}*))$$
$$\equiv \Delta(P; Z) \wedge (\forall \text{P}* \ (\exists \text{Z}* \ (\Delta(\text{P}*; \text{Z}*) \wedge (\text{P}* \leq P)) \ \Rightarrow \ (P \leq \text{P}*)))$$
$$\equiv \Delta(P; Z) \wedge \text{CIRC}[(\exists \text{Z}* \ \Delta(P; \text{Z}*)); P]$$

From this form, we see that the parallel circumscription of $P$ in $\Delta(P; Z)$ with $Z$ allowed to vary during the minimization is the same as ordinary parallel circumscription of $P$ in $(\exists \text{Z}* \ \Delta[P; \text{Z}*])$. The main difficulty now is how to handle the second-order quantifier in $(\exists \text{Z}* \ \Delta[P; \text{Z}*])$.

This problem can be dealt with if $\Delta$ is solitary, separable, or ordered in $Z$. Recall, for example, that if $\Delta$ is solitary in $Z$, it can be written as $N[Z] \wedge (E \leq Z)$, where $N[Z]$ is a formula that contains no positive occurrences of (any elements of) $Z$, and $E$ is a formula that contains no occurrences of (any elements of) $Z$. It is straightforward to show that $(\exists \text{Z}* \ N[\text{Z}*] \wedge (E \leq \text{Z}*)) \equiv N[E]$, where $N[E]$ is $N[\text{Z}*]$ with $E$ substituted for Z*.

These results establish the following theorem for the case in which $\Delta$ is solitary in $Z$.

**THEOREM 6.12**

$$\text{CIRC}[N(Z) \wedge (E \leq Z); P; Z] \equiv N(Z) \wedge (E \leq Z) \wedge \text{CIRC}[N(E); P]$$

*where $N$ has no positive occurrences of $Z$, and $E$ has no occurrences of $Z$. $E$, $P$, and $Z$ can be tuples of predicates.*

## *COROLLARY 6.2*

$$\mathrm{CIRC}[E_1 \wedge (E_2 \leq Z); P; Z] \equiv E_1 \wedge (E_2 \leq Z) \wedge \mathrm{CIRC}[E_1; P]$$

*where neither $E_1$ nor $E_2$ has any occurrences of $Z$. (That is, in this case varying $Z$ lets us drop the clause $(E_2 \leq Z)$ from $\Delta$ when computing the circumscription formula.)*

A simple default reasoning example will illustrate the effects on circumscription of letting a predicate vary. Let $\Delta$ be

$$\forall x \; \mathrm{Bird(x)} \wedge \neg\mathrm{Ab(x)} \Rightarrow \mathrm{Flies(x)}$$

$$\forall x \; \mathrm{Ostrich(x)} \Rightarrow \mathrm{Ab(x)}$$

Ordinary circumscription of Ab in $\Delta$ gives us

$$\mathrm{CIRC}[\Delta; \mathrm{Ab}]$$
$$\equiv \Delta \wedge (\forall x \; \mathrm{Ab(x)} \Leftrightarrow \mathrm{Ostrich(x)} \vee (\mathrm{Bird(x)} \wedge \neg\mathrm{Flies(x)}))$$

(The only abnormal things are either ostriches or birds that do not fly.)

A tighter characterization of Ab can be obtained by letting Flies vary. We see that we can use Corollary 6.2 to write

$$\mathrm{CIRC}[\Delta; \mathrm{Ab}; \mathrm{Flies}] \equiv \Delta \wedge \mathrm{CIRC}[(\forall x \; \mathrm{Ostrich(x)} \Rightarrow \mathrm{Ab(x)}); \mathrm{Ab}]$$
$$\equiv \Delta \wedge (\forall x \; \mathrm{Ab(x)} \Leftrightarrow \mathrm{Ostrich(x)})$$

(The only abnormal things are ostriches. Allowing Flies to vary lets us rule out the possibility that birds might not fly.)

As a more complex example, consider the taxonomic hierarchy we used earlier to illustrate delimited predicate completion. We repeat the formulas used in that example:

$$\mathrm{Flying\text{-}Ostrich(x)} \Rightarrow \mathrm{Ostrich(x)}$$

$$\mathrm{Flying\text{-}Ostrich(x)} \Rightarrow \mathrm{Ab3(x)}$$

$$\mathrm{Ostrich(x)} \Rightarrow \mathrm{Bird(x)}$$

$$\mathrm{Ostrich(x)} \Rightarrow \mathrm{Ab2(x)}$$

$$\mathrm{Bird(x)} \Rightarrow \mathrm{Thing(x)}$$

$$\mathrm{Bird(x)} \Rightarrow \mathrm{Ab1(x)}$$

$$\mathrm{Thing(Tweety)}$$

$$\mathrm{Ostrich(x)} \wedge \neg\mathrm{Ab3(x)} \Rightarrow \neg\mathrm{Flies(x)}$$

$$\text{Thing}(x) \wedge \neg \text{Ab1}(x) \Rightarrow \neg \text{Flies}(x)$$

$$\text{Bird}(x) \wedge \neg \text{Ab2}(x) \Rightarrow \text{Flies}(x)$$

Default reasoning can be accomplished by parallel circumscription of all predicates except Flies, but letting Flies vary. We let Flies vary because we are willing to let it take on any value required in the minimization of the other predicates. In being unconcerned about the value of Flies, we can use the entire belief set $\Delta$ in the minimization process (rather than just the taxonomic part used in delimited completion) to achieve the default assumptions that we want.

Thus, we circumscribe $\Delta$ (as before) in the predicates {Flying-Ostrich, Ostrich, Ab3, Bird, Ab2, Thing, Ab1}, letting Flies vary as well. In applying the procedure for parallel circumscription, we note first that $\Delta$ is solitary in Flies. To see this, observe that all the preceding clauses except the last one have no positive occurrences of Flies, and the antecedent of the last clause has no occurrence of Flies. Thus, we can use Theorem 6.12 and substitute Bird(x) $\wedge$ ¬Ab2(x) for Flies(x) in all but the last clause of the preceding set to yield

$$\text{Flying-Ostrich}(x) \Rightarrow \text{Ostrich}(x)$$

$$\text{Flying-Ostrich}(x) \Rightarrow \text{Ab3}(x)$$

$$\text{Ostrich}(x) \Rightarrow \text{Bird}(x)$$

$$\text{Ostrich}(x) \Rightarrow \text{Ab2}(x)$$

$$\text{Bird}(x) \Rightarrow \text{Thing}(x)$$

$$\text{Bird}(x) \Rightarrow \text{Ab1}(x)$$

$$\text{Thing}(\text{Tweety})$$

$$\text{Ostrich}(x) \wedge \neg \text{Ab3}(x) \Rightarrow \neg(\text{Bird}(x) \wedge \neg \text{Ab2}(x))$$

$$\text{Thing}(x) \wedge \neg \text{Ab1}(x) \Rightarrow \neg(\text{Bird}(x) \wedge \neg \text{Ab2}(x))$$

The last two clauses are subsumed by the fourth and sixth clauses, so the last two clauses can be eliminated. The circumscription we want can then be obtained by ordinary parallel circumscription of {Flying-Ostrich, Ostrich, Ab3, Bird, Ab2, Thing, Ab1} in the conjunction of the first seven clauses (without any predicates variable).

Since these clauses are ordered in {Flying-Ostrich, Ostrich, Ab3, Bird, Ab2, Thing, Ab1}, we can circumscribe by parallel predicate completion to obtain the following completion clauses (just as before):

1. Thing(x) $\Rightarrow$ Bird(x) $\vee$ x=Tweety
2. Bird(x) $\Rightarrow$ Ostrich(x)
3. Ostrich(x) $\Rightarrow$ Flying-Ostrich(x)

4. ¬Flying-Ostrich(x)
5. Ab1(x) ⇒ Bird(x)
6. Ab2(x) ⇒ Ostrich(x)
7. Ab3(x) ⇒ Flying-Ostrich(x)

## 6.6  Default Theories

We also can deal with the problem of nonmonotonic reasoning by defining a logic that uses nonstandard, nonmonotonic inference rules. Let us call these inference rules *default rules*, and the resulting theories *default theories*.

A default rule is an inference rule used to augment $\Delta$ under certain specified conditions, which we will soon describe. If $D$ is a set of such rules, we denote an augmentation of $\Delta$ (there may be more than one), with respect to $D$, by $\mathcal{E}[\Delta, D]$. (As before, augmentations include $\Delta$ and are closed under ordinary deduction.) Default rules are written in the form

$$\frac{\alpha(x) : \beta(x)}{\gamma(x)}$$

where $x$ is a tuple of individual constant schema variables, and $\alpha$, $\beta$, and $\gamma$ are wff schemata. (In running text, we write this rule as $\alpha(x) : \beta(x)/\gamma(x)$.)

The expressions above the line specify conditions on $\mathcal{E}[\Delta, D]$ that, if satisfied, permit (roughly speaking) the inclusion in $\mathcal{E}[\Delta, D]$ of the *consequent* below the line. A default rule is interpreted as follows: If there is an instance, $X_0$, of $x$ for which the ground instance $\alpha(X_0)$ follows from $\mathcal{E}[\Delta, D]$ and for which $\beta(X_0)$ is consistent with $\mathcal{E}[\Delta, D]$, then $\mathcal{E}[\Delta, D]$ includes $\gamma(X_0)$.

The rules are called *default rules* because they are useful in expressing beliefs about statements that are typically, but not necessarily always, true. For example, we might express the belief that birds typically fly by the default rule Bird(x) : Flies(x) / Flies(x). That is, if $x$ is a bird, and if it is consistent to believe that $x$ can fly, then it can be believed that $x$ can fly (or $x$ can fly "by default.") If $\Delta$ contained only the formulas Bird(Tweety) and Ostrich(x)⇒¬Flies(x), then $\mathcal{E}[\Delta, D]$ would contain Flies(Tweety). If we added to $\Delta$ the formula Ostrich(Tweety), use of the default rule would be *blocked* because Flies(Tweety) is not consistent with the new $\Delta$. Thus, default theories are nonmonotonic.

Our description of how default rules augment a theory may have been misleadingly simple because default theories may have more than one default rule, and rules may interact. The precise definition of $\mathcal{E}[\Delta, D]$ in terms of $\Delta$ and a set of default rules, $D$, must take into account the contributions of all of the default rules plus the closure of $\mathcal{E}[\Delta, D]$ under ordinary deduction. As we shall see, these interactions operate such as sometimes to warrant the existence of more than one augmentation.

Something like the CWA with respect to a predicate $P$ can be formulated using a default inference rule as follows:

$$\frac{: \neg P(x)}{\neg P(x)}$$

That is, if it is consistent to believe an instance of $\neg P(x)$, then that instance of $\neg P(x)$ may be believed. There is a difference, however, between the effects of the CWA with respect to $P$ and a default theory with this default. The CWA permits inferring an instance of $\neg P(x)$ if that instance is consistent with $\Delta$. The default rule permits it only if the instance is consistent with $\mathcal{E}[\Delta, D]$. Since there may be other default rules contributing to $\mathcal{E}[\Delta, D]$, the two techniques can lead to different augmentations.

Most applications of default rules involve a special case in which default rules take the form $\alpha(x) : \gamma(x)/\gamma(x)$. These are called *normal* default rules, and theories using them are called *normal* default theories. The CWA-type default rule we mentioned is an example of a normal default rule.

(It also is possible to define more general default rules. Consider one of the form $\alpha(x) : \beta_1(x), \beta_2(x), \ldots, \beta_n(x)/\gamma(x)$, the interpretation of which is that, if a ground instance $\alpha(X_0)$ follows from $\mathcal{E}[\Delta, D]$ and if *each* of the $\beta_i(X_0)$ is separately consistent with $\mathcal{E}[\Delta, D]$, then $\gamma(X_0)$ is in $\mathcal{E}[\Delta, D]$. This rule is different from one of the form $\alpha(x) : \beta_1(x) \wedge \beta_2(x) \wedge \ldots \wedge \beta_n(x)/\gamma(x)$, because the conjunction may be inconsistent with $\mathcal{E}[\Delta, D]$, whereas each conjunct may be separately consistent.)

Default theories have a number of interesting properties. (Some of these are particular to normal default theories.) We state, without proof, some of the more important properties here, and illustrate them with some examples.

(1) Just as circumscription sometimes does not produce a unique definition for a predicate, a default theory may have more than one augmentation. Consider, for example, the following (normal) default rules:

$$:\neg A/\neg A$$

$$:\neg B/\neg B$$

Now, if $\Delta$ is simply {AvB}, then there are two possible augmentations of $\Delta$; namely, {AvB, ¬A} and {AvB, ¬B}. Whereas the CWA, with respect to both A and B, would have produced an inconsistent augmentation, our default rules in this example give us two choices. We may select either one as an appropriate augmentation of our beliefs.

(2) The union of the two augmentations of the preceding example is inconsistent. In fact, we have the following result: If a

normal default theory has distinct augmentations, they are mutually inconsistent.

(3) There are default theories with no augmentations. Consider the default :A/¬A. If Δ is empty, so is $\mathcal{E}[\Delta, D]$. However, see (4).

(4) Every normal default theory has an augmentation.

(5) A default theory can have an inconsistent augmentation if and only if Δ itself is inconsistent. Since one can prove anything from an inconsistent augmentation and since augmentations (like theories) are closed under ordinary deduction, if a default theory has an inconsistent augmentation, this is its only augmentation.

(6) If $D$ and $D'$ are sets of normal default rules such that $D' \subseteq D$, then for any $\mathcal{E}[\Delta, D']$ there is an $\mathcal{E}[\Delta, D]$ such that $\mathcal{E}[\Delta, D'] \subseteq \mathcal{E}[\Delta, D]$. Thus, we say that normal default theories are *semimonotonic*. Adding new normal default rules does not require the withdrawal of beliefs, even though adding new beliefs might.

After having specified a set of default rules, how are we to use them to perform the kind of nonmonotonic reasoning inherent in their definitions? Typically, we must decide whether or not the initial beliefs Δ and the default rules $D$ warrant including some arbitrary formula $\phi$ among the augmented beliefs. That is, we must determine whether or not there is an augmentation $\mathcal{E}[\Delta, D]$ that includes the formula $\phi$.

We limit our definition of default proofs to the case of normal default theories. (The computation of augmentations for nonnormal default theories can be extremely complex; in fact, it is not even known yet what might be an appropriate proof theory for nonnormal defaults.) Informally, a default proof of $\phi$ given Δ and $D$ is just like an ordinary proof of $\phi$ from Δ, except that the (normal) default rules can be used as rules of inference. The use of a default rule must, in strict accordance with its definition, make the necessary consistency check before inferring its consequent. This test can be performed at the time of rule application in *forward* proofs. *Backward* proofs might best be done in two passes—first ignoring the consistency checks to find potential chains of inferences, and then performing the consistency checks, in forward order, on the default rules in a chain.

Suppose, for example, that $D$ consists of the following two default rules: Bird$(x)$:Flies$(x)$/Flies$(x)$ (by default, birds can fly), and FC$(x)$:Bird$(x)$/Bird$(x)$ (by default, feathered creatures are birds). If Δ contains only the statement FC(Tweety), then there is a default proof of Flies(Tweety). If, however, Δ contains instead the statements Ostrich(Tweety), Ostrich(x) ⟹ ¬Flies(x), and Ostrich(x) ⟹ FC(x), then there is no default proof of Flies(Tweety), because the rule instance Bird(Tweety):Flies(Tweety)/Flies(Tweety) cannot be used consistently.

Because default rules interact in complex ways, one must be careful about how knowledge is expressed. An illustrative instance of the sort of knowledge-representation difficulty that might arise is provided by the fact that default rules can be transitive. Suppose, for example, we have $D = \{\text{D}(x):\text{A}(x)/\text{A}(x), \text{A}(x):\text{E}(x)/\text{E}(x)\}$. We might interpret these expressions as saying: Typically, high-school dropouts are adults; and, typically, adults are employed. One of the consequences of these two rules also could be obtained by the combined rule $\text{D}(x):\text{E}(x)/\text{E}(x)$, the corresponding interpretation of which would be: Typically, high-school dropouts are employed. Although we might assent to the first two rules, we would not necessarily want to use the combination.

This unwelcome transitivity can be blocked in one of two ways. We could change the second default rule to the nonnormal rule $\text{A}(x):$ $[\neg\text{D}(x) \wedge \text{E}(x)]/\text{E}(x)$. Nonnormal defaults do not enjoy the simple and desirable properties of normal defaults, however. We often can block transitivity by a more careful formulation using only normal defaults. In our example, because, typically, adults are not high-school dropouts, we could use the following normal defaults: $\{\text{D}(x):\text{A}(x)/\text{A}(x), [\text{A}(x)$ $\wedge \neg\text{D}(x)]:\text{E}(x)/\text{E}(x), \text{A}(x):\neg\text{D}(x)/\neg\text{D}(x)\}$. Now we would not conclude that some particular high-school dropout was employed.

## 6.7  Bibliographical and Historical Remarks

Almost all interesting applications of AI require some form of nonmonotonic reasoning because the knowledge that AI systems contain about their domains is always subject to change and elaboration, yet it is still necessary for AI systems to reason as best they can with the knowledge they currently have. Reiter gave an excellent summary of nonmonotonic reasoning and its applications in AI [Reiter 1987b]. An important, typical application is to the problem of equipment diagnosis [Reiter 1987a]. McCarthy discussed several applications of one type of nonmonotonic reasoning [McCarthy 1986].

The closed-world assumption (CWA) is an important convention of database design. Reiter [Reiter 1978] was the first to describe and prove several of its properties. Theorem 6.1 is taken from [Shepherdson 1984]. The domain-closure assumption (DCA) and unique-names assumption (UNA) were discussed by Reiter [Reiter 1980b].

The *qualification problem* was discussed by McCarthy [McCarthy 1980]. It often is cited as one of the reasons that a strictly logical approach to AI will not work, and it has motivated much of the work in nonmonotonic reasoning.

Predicate completion of a set of clauses was first described by Clark [Clark 1978]. Parallel predicate completion suggests itself by analogy with parallel circumscription. Taxonomic hierarchies are ubiquitous in AI applications. Several frame-based systems have facilities for property

inheritance and default reasoning in these hierarchies [Stefik 1986]. Our use of the Ab predicate in this connection is based on a suggestion by McCarthy [McCarthy 1986].

Circumscription as a method of nonmonotonic reasoning was first proposed by McCarthy [McCarthy 1980]. Our notation follows that of Lifschitz [Lifschitz 1985a]. (The alternative form of circumscription—which says that any alleged strengthening of $P$ by $P'$ is no real strengthening because the circumscription of $P$ already implies $P'$—was attributed by Minker and Perlis [Minker 1984] to Reiter.) The circumscription formula is a formula in second-order logic. Although our treatment of circumscription in this book is essentially limited to those cases that collapse to first order, the reader may want to look at Enderton's chapter on second-order logic [Enderton 1972].

Theorems 6.4 through 6.6 all were developed by Lifschitz; their proofs are given in [Lifschitz 1987b]; Theorems 6.5 and 6.6 are stated without proof in [Lifschitz 1985a]. That CIRC$[\Delta; P]$ is collapsible to a first-order formula if all occurrences of $P$ are positive in $\Delta$ follows immediately from results in Lifschitz [Lifschitz 1986c], and also was proved in [Lifschitz 1987b].

Etherington, Mercer, and Reiter showed that the circumscription of a formula having no minimal model is inconsistent, and they also proved a sufficient condition on the consistency of circumscription (Theorem 6.7) [Etherington 1985]. Theorem 6.8 was developed by Lifschitz [Lifschitz 1986b]. Theorems 6.7 and 6.8 are each special cases of Theorem 6.9, also developed by Lifschitz [Lifschitz 1986b]. Perlis and Minker [Perlis 1986] also have done work relating properties of circumscription to minimal models.

Parallel circumscription is a natural extension of ordinary circumscription. Theorem 6.10, developed by Lifschitz [Lifschitz 1986c, Lifschitz 1987b], often is useful in computing parallel circumscription. Otherwise, for *ordered* formulas, parallel circumscription can be computed using Theorem 6.11. (Ordered formulas and Theorem 6.11 are original here.) Theorem 6.12, developed by Lifschitz [Lifschitz 1987b], is useful for computing circumscription with variable predicates.

Etherington [Etherington 1986] and Lifschitz [Lifschitz 1986b] independently extended Theorem 6.7 to the case of parallel circumscription with variable predicates. That is, the parallel circumscription of universal theories (even with variable predicates) is consistent if the theory itself is consistent.

Several authors have been concerned with the relationship between circumscription and other nonmonotonic reasoning methods. For example, there are conditions under which parallel circumscription and the CWA augment a belief set identically. Lifschitz [Lifschitz 1985b] showed that the CWA applied to the belief set produces the same effect as does the parallel circumscription of all predicates in the belief set if (1) the CWA can be applied to a belief set consistently, (2) constant terms in the belief set name all possible objects in the domain (the DCA), and (3) different

constant terms in the belief set denote different objects in the domain (the UNA). Gelfond, Przymusinska, and Przymusinski studied the relationships among various generalizations of the CWA and circumscription [Gelfond 1986]. Reiter was the first person to show that predicate completion is a special case of circumscription (using an argument similar to that used in the proof of Theorem 6.4) [Reiter 1982].

Przymusinski [Przymusinski 1986] proposed a method for deciding whether or not there is a minimal model of a theory $T$ that also satisfies a formula $\phi$; it can be used to answer queries in circumscriptive theories.

Imielinski and Grosof each investigated relationships between default logics and circumscription [Imielinski 1985, Grosof 1984].

Default logic was originally proposed and analyzed in a paper by Reiter [Reiter 1980a]. Our discussion of default theories is based on that paper. He shows that default logic is not even semidecidable, but describes a resolution theorem prover that can be used for top-down or backward searches for default proofs. Reiter and Criscoulo [Reiter 1983] gave examples of default-rule formulations of typical commonsense reasoning problems and showed how to avoid various pitfalls without using nonnormal defaults.

Other methods of nonmonotonic reasoning also have been proposed. McDermott and Doyle [McDermott 1980, McDermott 1982a] defined a logic having a *modal operator* **M**. (Modal operators are discussed in Chapter 9.) In the semantics for such a logic, the formula **M**$P$ has value true just in case $P$ is consistent (with the theory based on $\Delta$). Any derivations of **M**$P$ or its consequences are nonmonotonic because the meaning of **M** depends globally on the theory. If we add another formula to $\Delta$, **M**$P$ may no longer be consistent. With a somewhat different application in mind, Moore [Moore 1985b] proposed a variant, which he called *autoepistemic logic*, and compared it with McDermott and Doyle's nonmonotonic logic. Konolige [Konolige 1987] analyzed the connections between default theories and autoepistemic logic.

Several additional papers appear in the proceedings of a workshop on nonmonotonic reasoning [Nonmonotonic 1984].

---

## Exercises

1. *Idempotence.* We denote the CWA augmentation of $\Delta$ by CWA[$\Delta$]. Prove that

$$\text{CWA[CWA[}\Delta\text{]]} = \text{CWA[}\Delta\text{]}$$

   (Assume CWA[$\Delta$] is consistent.)

2. *Insensitivity to negative clauses.* Suppose $\Delta$ is Horn and consistent. Show that deleting a negative clause (i.e., one without any positive literals) from $\Delta$ does not change the CWA augmentation of $\Delta$.

3. *Inconsistencies.* Prove that, if a consistent $\Delta$ contains only Horn clauses and is inconsistent with $\neg L_1 \wedge \neg L_2$ (where $L_1$ and $L_2$ are positive literals), then either $\Delta \wedge \neg L_1$ or $\Delta \wedge \neg L_2$ is inconsistent also.

4. *Even and odd.* Compute the completion of EVEN in the conjunction of the following formulas:

$$\forall x \ \mathtt{ODD(x)} \ \wedge \ \mathtt{x>0} \ \Rightarrow \ \mathtt{EVEN(Succ(x))}$$

$$\forall x \ \mathtt{ODD(x)} \ \wedge \ \mathtt{x>0} \ \Rightarrow \ \mathtt{EVEN(Pred(x))}$$

5. *Integers.* Compute the completion of INT in $\mathtt{INT(0)} \ \wedge \ (\mathtt{INT(x)} \ \Rightarrow \ \mathtt{INT(Succ(x))})$.

6. *Delimited predicate completion.* Discuss how *delimited* predicate completion might produce an inconsistent augmentation.

7. *Completion.* Compute the completion of P in the following clauses:

$$\mathtt{Q1(x)} \ \wedge \ \mathtt{Q2(x)} \ \Rightarrow \ \mathtt{P(F(x))}$$

$$\mathtt{Q3(x)} \ \Rightarrow \ \mathtt{P(G(x))}$$

8. *Is there a Q that is not a P?* Express in words what the effect would be of circumscribing Q in P < Q.

9. *Parallel.* Compute $\mathrm{CIRC}[(\forall x \ \mathtt{Q(x)} \ \Rightarrow \ \mathtt{P1(x)} \ \vee \ \mathtt{P2(x)});\mathtt{P1},\mathtt{P2}]$.

10. *Knights and Knaves.* Let $\Delta$ be the conjunction of the following formulas:

$$\forall x \ \mathtt{KNIGHT(x)} \ \Rightarrow \ \mathtt{PERSON(x)}$$

$$\forall x \ \mathtt{KNAVE(x)} \ \Rightarrow \ \mathtt{PERSON(x)}$$

$$\forall x \ \mathtt{KNAVE(x)} \ \Rightarrow \ \mathtt{LIAR(x)}$$

$$\exists x \ \neg\mathtt{LIAR(x)} \ \wedge \ \neg\mathtt{KNAVE(x)}$$

$$\mathtt{LIAR(MORK)}$$

$$\mathtt{KNAVE(BORK)}$$

   a. Compute $\mathrm{CIRC}[\Delta;\mathtt{LIAR}]$.

   b. Compute $\mathrm{CIRC}[\Delta;\mathtt{LIAR},\mathtt{KNAVE}]$.

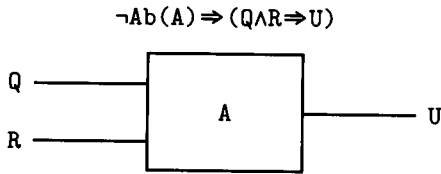$$\neg Ab(A) \Rightarrow (Q \wedge R \Rightarrow U)$$



**Figure 6.2   An AND gate.**

11. *AND gate.*   The AND gate, A, shown in Figure 6.2, is described by the following formula:

$$\neg Ab(A) \Rightarrow (Q \wedge R \Rightarrow U)$$

(Unless A is abnormal, Q and R imply U. Q denotes the proposition "input 1 is on;" R denotes the proposition "input 2 is on;" U denotes the proposition "the output is on.")

   a. Suppose, in fact, that Q, R, and U are true. Use the circumscription of Ab in this theory to prove that nothing is abnormal. The "theory" in this case is

$$Q \wedge R \wedge U \wedge (\neg Ab(A) \Rightarrow (Q \wedge R \Rightarrow U))$$

   b. Suppose now, instead, that Q and R are true, but U is false. Use the circumscription of Ab in *this* theory to show that the *only* abnormal thing is A.

12. *Both P and Q.*   Let $\Delta$ consist of these two formulas:

$$\forall x\ R(x) \Rightarrow P(x)$$
$$\forall x\ R(x) \Rightarrow Q(x)$$

Show how to use circumscription to assert that the only objects that satisfy both P and Q also satisfy R. Hint: Use $(\forall x\ (P(x) \wedge Q(x)) \Rightarrow Ab(x))$ and minimize Ab with P and Q variable.